**Università degli Studi di Padova**

DEPARTMENT OF INFORMATION ENGINEERING

*Ph.D. Course in* INFORMATION ENGINEERING
SCIENCE AND INFORMATION TECHNOLOGY CURRICULUM
XXXII SERIES

# End-to-End Design and Evaluation of mmWave Cellular Networks

*Coordinator*
PROF. ANDREA NEVIANI
*Supervisor*
PROF. MICHELE ZORZI

*Ph.D. Candidate*
MICHELE POLESE

ACADEMIC YEAR 2019/2020

Dec. 24, 1968, 9:48 a.m. UTC
Jerry Carr, CapCom: *"Apollo 8, Houston. One minute to LOS. All systems GO."*
Jim Lovell, Apollo 8: *"We'll see you on the other side"*

# Abstract

The next generation of cellular networks (5G) is being designed to provide unprecedented performance in mobile scenarios, with an increase in capacity, ultra-low latency and a massive number of connections. This will require the integration of novel technologies in more advanced and complex networks. Millimeter wave (mmWave) communications are considered as a key enabler for ultra-high datarates and low latency, thanks to the massive amount of available bandwidth at such high frequencies. Nonetheless, there are a number of challenges that must be solved before this technology can be deployed, mainly related to the high propagation loss, the need for directional communications and the blockage.

This thesis provides system-level solutions to make mmWave mobile networks more reliable, robust and better performing. Notably, we consider mmWave links as a part of more complex, end-to-end networks, in which the quality of experience that the end user perceives is the result of the interaction among the variability and unreliability of the mmWave channel, the full protocol stack, and the deployment strategy of the wireless network. To this end, we develop and describe a tool for the simulation of end-to-end mmWave cellular networks that, combined with analysis and experimental results, makes it possible to consistently evaluate how these systems behave in their entirety.

The main research areas that this thesis explores are the design and evaluation of (i) architectures for mmWave systems, in terms of mobility and beam management, and wireless backhaul solutions; (ii) protocols for end-to-end connectivity over mmWave networks; and (iii) intelligent and data-driven optimizations in cellular networks. Among other results, we highlight the importance of multi connectivity for mmWave systems, in the access network and at the transport layer, discuss the tradeoffs of beam management in 3GPP NR, propose how to update protocols at the transport layer for an improved end-to-end performance, and evaluate practical approaches for the integration of intelligent techniques in 5G networks.

# Sommario

La progettazione della prossima generazione di reti cellulari (5G) ha l'obiettivo di garantire prestazioni senza precedenti in scenari mobili, con un aumento nella capacità, latenze molto basse e un elevato numero di connessioni. Questo richiederà l'integrazione di nuove tecnologie in reti più complesse ed avanzate. Le comunicazioni a frequenze millimetriche (mmWave) sono considerate una componente fondamentale per raggiungere altissime velocità dati e bassa latenza, grazie alla smisurata quantità di banda disponibile a frequenze così elevate. Ci sono tuttavia numerosi problemi che vanno risolti prima di poter dispiegare questa tecnologia, principalmente legati all'elevata perdita di potenza in propagazione, alla necessità di comunicazioni direzionali e all'ostruzione dei segnali da parte di comuni ostacoli (ad esempio, il corpo umano stesso).

L'obiettivo di questa tesi è proporre soluzioni di sistema per rendere le reti mobili mmWave più affidabili, robuste e con prestazioni migliori. In particolare, consideriamo il collegamento a frequenze millimetriche come una sola parte di reti più complesse, in cui la qualità sperimentata dall'utente finale è il risultato dell'interazione della variabilità e inaffidabilità del canale mmWave, dell'intero stack protocollare, e dell'architettura della rete mobile. Pertanto, sviluppiamo e descriviamo uno strumento per simulazione di reti cellulari mmWave che considera le prestazioni tra i due capi della rete, e che, combinato con analisi e risultati sperimentali, consente di valutare come questi sistemi si comportino nella loro interezza.

Le aree di ricerca principali che questa tesi esplora sono la progettazione e la valutazione di (i) architetture per sistemi a onde millimetriche, per quel che riguarda la mobilità e la gestione delle trasmissioni direzionali, e soluzioni di backhaul senza fili; (ii) protocolli per la connessione tra due capi della rete usando almeno una connessione mmWave; e (iii) sistemi che ottimizzano reti cellulari usando i dati che le stesse reti generano. Tra i vari risultati, sottolineiamo l'importanza della disponibilità di connessioni multiple per sistemi mmWaves, sia nel collegamento di accesso che al livello di trasporto, discutiamo i compromessi della gestione della direzionalità, proponiamo come aggiornare i protocolli al livello di traporto per una migliore prestazione globale, e valutiamo approcci pratici per l'integrazione di tecniche intelligenti in reti 5G.

# Contents

## II The Architecture: System Level Design of 5G mmWave Networks 47

# Listing of figures

xvii

# Listing of tables

# Acronyms

**3GPP** 3rd Generation Partnership Project
**5G** 5th generation
**5GC** 5G Core

**ADC** Analog to Digital Converter
**AGWN** Additive White Gaussian Noise
**AI** Artificial Intelligence
**AIMD** Additive Increase Multiplicative Decrease
**AM** Acknowledged Mode
**AMC** Adaptive Modulation and Coding
**AMF** Access and Mobility Management Function
**AQM** Active Queue Management
**AVC** Advanced Video Coding

**BALIA** Balanced Link Adaptation Algorithm
**BBU** Base Band Unit
**BDP** Bandwidth-Delay Product
**BF** Beamforming
**BLER** Block Error Rate
**BRR** Bayesian Ridge Regressor
**BSR** Buffer Status Report

**CA** Carrier Aggregation
**CB** Code Block
**CC** Carrier Component
**CC** Congestion Control
**CDF** Cumulative Distribution Function
**CN** Core Network
**CoDel** Controlled Delay Management
**CQI** Channel Quality Information
**CRAN** Cloud RAN
**CRS** Cell Reference Signal
**CSI** Channel State Information
**CSI-RS** Channel State Information - Reference Signal
**CU** Central Unit

**DAC** Digital to Analog Converter
**DAG** Directed Acyclic Graph
**DASH** Dynamic Adaptive Streaming over HTTP
**DC** Dual Connectivity
**DCE** Direct Code Execution
**DCI** Downlink Control Information
**DL** Downlink
**DMR** Deadline Miss Ratio
**DMRS** DeModulation Reference Signal
**DU** Distributed Unit

**E2E** end-to-end

**ECN** Explicit Congestion Notification
**EDF** Earliest Deadline First
**eMBB** Enhanced Mobile Broadband
**eNB** evolved Node Base
**EPC** Evolved Packet Core
**ETA** Estimated Time of Arrival

**FDD** Frequency Division Duplexing
**FDM** Frequency Division Multiplexing
**FDMA** Frequency Division Multiple Access
**FR2** Frequency Range 2
**FS** Fast Switching
**FW** Flow Window

**gNB** Next Generation Node Base
**GOP** Group of Pictures
**GPR** Gaussian Process Regressor
**GTP** GPRS Tunneling Protocol

**HARQ** Hybrid Automatic Repeat reQuest
**HH** Hard Handover
**HOL** Head-of-Line
**HQF** Highest-quality-first
**HTTP** HyperText Transfer Protocol

**IA** Initial Access
**IAB** Integrated Access and Backhaul
**IETF** Internet Engineering Task Force
**IMSI** International Mobile Subscriber Identity
**IMT** International Mobile Telecommunication
**IoT** Internet of Things
**ITU** International Telecommunication Union

**KPI** Key Performance Indicator

**LOS** Line-of-Sight
**LTE** Long Term Evolution

**M2M** Machine to Machine
**MAC** Medium Access Control
**MCS** Modulation and Coding Scheme
**MEC** Mobile Edge Cloud
**MI** Mutual Information
**MIB** Master Information Block
**MIMO** Multiple Input, Multiple Output
**ML** Machine Learning
**MLR** Maximum-local-rate
**MME** Mobility Management Entity
**mMTC** Massive Machine-Type Communications
**mmWave** millimeter wave
**MPTCP** Multipath TCP
**MR** Maximum Rate
**MSE** Mean Square Error
**MSS** Maximum Segment Size

xxiv

**MT** Mobile Termination
**MTD** Machine-Type Device
**MTU** Maximum Transmission Unit

**NALU** Network Abstraction Layer Unit
**NB-IoT** Narrow Band IoT
**NFV** Network Function Virtualization
**NLOS** Non-Line-of-Sight
**NOW** Non Overlapping Window
**NR** New Radio
**NSA** Non Stand Alone

**O2I** Outdoor to Indoor
**OFDM** Orthogonal Frequency Division Multiplexing
**OLIA** Opportunistic Linked Increase Algorithm

**PA** Position-aware
**PBCH** Physical Broadcast Channel
**PDCCH** Physical Downlonk Control Channel
**PDCP** Packet Data Convergence Protocol
**PDSCH** Physical Downlink Shared Channel
**PDU** Packet Data Unit
**PF** Proportional Fair
**PGW** Packet Gateway
**PHY** Physical
**PPP** Poisson Point Process
**PRB** Physical Resource Block
**PSNR** Peak Signal to Noise Ratio
**PSS** Primary Synchronization Signal
**PUCCH** Physical Uplink Control Channel
**PUSCH** Physical Uplink Shared Channel

**QoS** Quality of Service
**QUIC** Quick UDP Internet Connections

**RACH** Random Access Channel
**RAN** Radio Access Network
**RAT** Radio Access Technology
**RED** Random Early Detection
**RF** Radio Frequency
**RFC** Request for Comments
**RFR** Random Forest Regressor
**RLC** Radio Link Control
**RLF** Radio Link Failure
**RLNC** Random Linear Network Coding
**RMSE** Root Mean Squared Error
**RR** Round Robin
**RRC** Radio Resource Control
**RRM** Radio Resource Management
**RSRP** Reference Signal Received Power
**RSS** Received Signal Strength
**RTT** Round Trip Time
**RX** Receiver

**SA** standalone

**SACK** Selective Acknowledgment
**SAP** Service Access Point
**SCH** Secondary Cell Handover
**SCOOT** Split Cycle Offset Optimization Technique
**SCTP** Stream Control Transmission Protocol
**SDAP** Service Data Adaptation Protocol
**SDM** Space Division Multiplexing
**SDMA** Spatial Division Multiple Access
**SDN** Software Defined Networking
**SGW** Service Gateway
**SI** Study Item
**SIB** Secondary Information Block
**SINR** Signal to Interference plus Noise Ratio
**SM** Saturation Mode
**SNR** Signal-to-Noise-Ratio
**SON** Self-Organizing Network
**SPTCP** Single Path TCP
**SRB** Service Radio Bearer
**SRS** Sounding Reference Signal
**SS** Synchronization Signal
**SSS** Secondary Synchronization Signal
**ST** Spanning Tree
**SVC** Scalable Video Coding

**TB** Transport Block
**TCP** Transmission Control Protocol
**TDD** Time Division Duplexing
**TDM** Time Division Multiplexing
**TDMA** Time Division Multiple Access
**TfL** Transport for London
**TM** Transparent Mode
**TR** Technical Report
**TS** Technical Specification
**TTI** Transmission Time Interval
**TTT** Time-to-Trigger
**TX** Transmitter

**UAV** Unmanned Aerial Vehicle
**UDP** User Datagram Protocol
**UE** User Equipment
**UL** Uplink
**UM** Unacknowledged Mode
**UML** Unified Modeling Language
**UPA** Uniform Planar Array
**URLLC** Ultra Reliable and Low Latency Communications
**UTC** Urban Traffic Control

**VM** Virtual Machine
**VR** Virtual Reality
**VSS** Video Streaming Server

**WBF** Wired Bias Function
**WF** Wired-first

# 1
## Introduction

The 5th generation (5G) of cellular networks is being designed and deployed to address the traffic demands and new use cases of the digital society beyond 2020 [1]. In particular, the International Telecommunication Union (ITU) has defined in the International Mobile Telecommunication (IMT) 2020 framework a set of requirements that any 5G network will have to satisfy [2]. More specifically, 5G deployments should support: (i) a user experienced rate of at least 100 Mbps, with a peak data rate in ideal conditions of 20 Gbps, and three times higher spectral efficiency with respect to 4G; (ii) ultra-low latency, i.e., 1 ms round-trip over the air; (iii) support for mobility, with communications at up to 500 km/h; (iv) an area capacity of 10 Mbps/m$^2$ with up to $10^6$ connections per km$^2$; and (v) a 100x increase in energy efficiency with respect to 4G networks.

These requirements have been mapped into a number of different use cases and services that end users will benefit from when connected to 5G networks. The ITU report in [2] introduces three main categories for the usage scenarios:

- Enhanced Mobile Broadband (eMBB), to address the need for higher datarates and better coverage in human communications in mobile contexts. The ITU distinguishes between the hotspot case, in which a high volume of traffic needs to be served in a small but densely populated area, with low mobility, and the wide area coverage case, in which seamless connectivity in medium and high mobility is guaranteed, with a datarate that may be smaller than in the hotspot case;

- Ultra Reliable and Low Latency Communications (URLLC), in which the requirements on reliability and low-latency (with a high-enough throughput) are combined to address new use cases for wireless networks, such as industrial automation [3] or remote surgery [4];

- Massive Machine-Type Communications (mMTC), where a large number of devices connect to the network but generate low traffic, for example for monitoring and periodic reporting.

Paper [5] identifies similar use cases, with 5G connectivity being able to provide very good service even in crowded areas, ubiquitous connectivity and real-time communications.

The research on solutions to support these scenarios and requirements has been particularly active in the last 5 years [6], leading to the quick development of a new set of 3rd Generation Partnership Project (3GPP) specifications (i.e., NR) [7] and to a number of new groundbreaking technologies in the wireless communication and networking domains. Notably, papers [6,8] have identified some key innovations that are now part of the 5G network specifications:

- mmWave communications, i.e., the usage of the spectrum in the 30-300 GHz band[1] in the RAN, in order to exploit the large chunks of free spectrum available at such high frequencies and reach higher datarates at the physical layer;

- massive Multiple Input, Multiple Output (MIMO), i.e., the deployment of antennas with a larger number of elements than the number of users served by the base stations [11], which yields higher spectral efficiency and a smoother channel response [8];

- flexible RAN architecture, i.e., the possibility of splitting the uplink and the downlink, and/or the control and user planes over different links;

- caching and improved device capabilities, e.g., through device-to-device communications, support of multi-connectivity through multiple technologies and advanced interference rejection [6];

- Network Function Virtualization (NFV) and Software Defined Networking (SDN) for the design of the core network, and a disaggregation of the higher layers and lower layers of the base stations in different networking entities [7], towards a Cloud RAN (CRAN) deployment paradigm [12];

- a smarter network, thanks to the integration of data-driven and machine-learning-based optimizations [13].

The research of this thesis develops in this context, addressing some of the technology challenges related to the design and evaluation of end-to-end solutions of mmWave cellular networks and the deployment of machine learning and intelligence in the network. In this regard, we keep a close relation with the development of 3GPP NR, with most of our proposals and solutions that can be seamlessly integrated in NR-compliant deployments, but, at the same time, the generality of our approaches and the validity of the results applies also to future evolutions of cellular networks into beyond 5G and 6G architectures. In the remainder of this introduction, we will provide some details on 3GPP NR specifications in Sec. 1.1, according to our previous work [428], and discuss the potentials and challenges of mmWave communications in Sec. 1.2. Finally, we will highlight the main contributions of this thesis in Sec. 1.3 and present its structure in Sec. 1.4.

## 1.1   3GPP NR: the Set of Specifications for 5G Networks

Long Term Evolution (LTE) is the set of specifications that the 3GPP has introduced in 2009 and evolved since then to satisfy the current 4G requirements. The evolutions of LTE will match some of the next generation requirements in specific deployment scenarios [14], but they will not be able to effectively address all the 5G use cases. For example, as we discuss in [428], LTE operates with a maximum of 20 MHz per carrier, thus limiting the achievable data rate, and has a rigid frame structure that makes it difficult to reduce the round-trip latency below 1 ms. Moreover, LTE has not been designed to account for energy efficiency (e.g., pilot signals are always-on) and to support a massive number of connections (even though this is targeted by the recent Narrow Band IoT (NB-IoT) evolution).

In order to overcome these limitations of LTE networks and address the use cases and requirements of 5G networks previously mentioned, the 3GPP has recently defined a new Radio Access

---

[1]Notice that the industry loosely refers to mmWave communications when considering the spectrum above 10 GHz, with the 24.25–27.5 GHz band being candidate for early deployments of 5G networks in Europe [9, 10].

**Figure 1.1:** Graphical representation of the main novelties introduced in 3GPP NR.

Technology (RAT), i.e., 3GPP NR[2], that introduces novel designs and technologies to comply with the 5G requirements. NR has been standardized by 3GPP with a first set of specifications[3] (Release 15) in December 2017 and a complete one published in June 2018. Release 16 for NR is expected to be completed in December 2019, and will be composed of a set of specifications that match the ITU 5G requirements [16].

Fig. 1.1 represents the main novelties that NR has introduced in cellular networks specifications. In particular, following the aforementioned technology trends, NR exploits a new spectrum, i.e., it is the first set of specifications for cellular networks to support the millimeter wave (mmWave) band, and features new techniques such as massive MIMO, flexibility in terms of frame structure, to target different use cases, and multiple deployment options for the RAN. Moreover, a new core network design (i.e., the 5G Core (5GC)) has been introduced to offer network slicing and virtualization, and different deployment options and inter-networking with LTE have been specified. In the following paragraphs, we will provide an overview of the main design innovations of NR, focusing on those that will be the most relevant for this thesis. Moreover, a deeper discussion on mmWaves and their support in NR is left for Sec. 1.2 and Chapter 4.

**A Flexible Physical Layer –**  The main characteristic of the NR physical layer is its flexibility: the standard, indeed, provides a general technology framework designed to address the different and, in some cases, conflicting 5G requirements [2] and to be forward compatible, so that it can

---

[2]While NR was originally meant as the acronym for "New Radio" [15], according to the latest 3GPP specifications [7] it has lost its original meaning and it now refers to the 5G Radio Access Network.

[3]The specifications for NR are in the *38 series* of the 3GPPs Technical Specifications (TSs), together with the Technical Reports (TRs) that contain related studies. Other relevant RAN specifications can be found in the *36* (LTE) and *37* (LTE-NR inter-networking) series.

accommodate future applications and use cases.

Both LTE and NR use the Orthogonal Frequency Division Multiplexing (OFDM) modulation, which divides the available *time* resources in frames of 10 ms with subframes of 1 ms, and *frequency* resources in subcarriers with spacing $\Delta f$. Moreover, subframes are further divided in slots and symbols, where the combination of a single OFDM symbol and a single subcarrier constitutes the smallest physical resource in NR. While with LTE the symbol duration and the subcarrier spacing are fixed, with NR it is possible to configure different OFDM *numerologies*[4] on a subframe basis, i.e., every subframe is self-contained and can be characterized by a different numerology [18]. This makes it possible to address different 5G use cases with a single RAT: for example, a shorter OFDM symbol duration, combined with a higher subcarrier spacing, can be used for high-data-rate and low-latency traffic, while lower subcarrier spacing can be used for low-frequency narrowband communications for machine-generated traffic [19]. Fig. 1.1 illustrates an example of NR frame structure with two different possible subcarrier spacings.

Another main NR novelty with respect to LTE is the support for ultra-low latency communications [20], to target the sub-1 ms round-trip latency requirement of 5G. First of all, the usage of larger subcarrier spacings and shorter symbols has the potential to reduce the transmission time with respect to the basic LTE frame structure. Moreover, control information related to modulation and resource allocation can be added at the beginning of data packets, allowing the devices to start decoding as soon as they start receiving data [18]. This also translates into tighter processing constraints in 5G NR devices, which must be able to process a received packet in just a few hundreds of microseconds (the actual constraints depend on the subcarrier spacing, as discussed in [14]). Another consequence is that the devices will be able to transmit the Hybrid Automatic Repeat reQuest (HARQ) acknowledgment after just one slot, making it possible to reduce the round-trip latency below 1 ms.[5] Moreover, latency-sensitive data does not need to wait for a new slot to be transmitted, but the base station may decide to transmit it as soon as possible using *mini-slots*, i.e., groups of at least 2 OFDM symbols that can be allocated to a data transmission and do not need to be aligned with the beginning of a standard slot [22].

Finally, in order to increase the flexibility and the energy efficiency of the RAN, NR limits the number of always-on reference signals, thereby configuring them to match the deployment scenario and increase the energy efficiency [23]. Moreover, the self-contained subframe and the minimization of always-on signals make the NR design forward-compatible, i.e., they enable the evolution of the NR RAT to support unforeseen use cases with novel technologies and solutions without compromising the support for legacy devices [22].

**Massive MIMO –**   While the combination of extreme cell densification, increased system bandwidth, and more flexible spectrum usage (e.g., by resource sharing) represents a feasible and sustainable solution to meet 5G performance requirements, MIMO techniques have also emerged in modern wireless networks to improve reliability and spectral efficiency. The main concept is to use multiple transmit and receive antennas to exploit multipath propagation. Among the possible antenna array designs, the most suitable approach is the use of Uniform Planar Arrays (UPAs) where the antenna elements are evenly spaced on a two-dimensional plane and a 3D beam can be synthesized by adapting both azimuth and elevation planes.

As previously discussed, 5G networks will rely on massive MIMO techniques, i.e., with large antenna arrays in which the number of available antennas is much larger than the number of users that are being served by the base station [11]. At sub-6 GHz, the usage of massive MIMO

---

[4]The term numerology refers to a set of parameters for the OFDM waveform, such as subcarrier spacing and symbol duration [17].

[5]In LTE (Release 14), the round-trip latency was fixed to 3 ms [21].

provides *channel hardening*, i.e., the combined usage of a massive number of antennas decreases the channel variability by averaging the small-scale fading [24].

However, massive MIMO comes with its own set of challenges, mainly related to the need for a precise channel estimation in a dynamic radio environment. Therefore, the main issues are (i) hardware impairments, which may introduce non-reciprocity in the uplink and downlink channels, making channel estimation more complex, and (ii) pilot contamination, due to the limited number of Channel State Information (CSI) pilots available and that must consequently be re-used, causing interference.

For NR, support for massive MIMO is introduced by using high-resolution CSI feedback and uplink Sounding Reference Signals (SRSs) targeting the utilization of channel reciprocity (e.g., twelve orthogonal demodulation reference signals are specified for multi-user MIMO transmission operations). Additionally, NR focuses on the support of *distributed MIMO*, through which the NR devices can receive multiple independent Physical Downlink Shared Channels (PDSCHs) per slot to enable simultaneous transmissions from multiple points to the same receiver.

**Towards a disaggregated and virtualized network –** NR has been designed with flexibility in mind, to address the different 5G use cases. This has an impact also on the possible cellular network deployment architectures [25, 26], which follow two recent emerging technology trends: disaggregation and CRAN [27], and virtualization [26].

The LTE RAN and the associated core network (Evolved Packet Core (EPC)) are characterized by the deployment of standalone pieces of equipment and servers, e.g., the evolved Node Bases (eNBs), and the core elements such as the Packet Gateways (PGWs) and Mobility Management Entities (MMEs). With NR, instead, the Next Generation Node Base (gNB) can be split into separate physical units, i.e., the Distributed Unit (DU), which contains the lower layers of the protocol stack and is deployed in the field, and the Central Unit (CU) incorporating complete stack functionalities, which can be co-located with the DU or hosted in a data center facility, according to the CRAN paradigm. As discussed in [28], this allows network operators to deploy the 5G RAN according to the use cases they want to serve, e.g., an ultra-dense small cell deployment with low utilization but high peak rate can rely on the CRAN CU/DU split to maximize the statistical multiplexing gain and enable a centralized control of the RAN, while a rural low-density deployment for the support of Internet of Things (IoT) applications can feature complete gNB nodes. Moreover, as shown in the top-right part of Fig. 1.1, in order to smooth the transition between the different network generations and reuse the widely deployed LTE and EPC infrastructure, the NR specifications foresee a Non Stand Alone (NSA) deployment, in which NR gNBs are connected to the EPC, possibly with a Dual Connectivity (DC) setup aided by LTE [29]. The other option is a standalone (SA) deployment, in which both the RAN and the core network respect the 5G specifications.

Finally, the 5G core network has been redesigned with respect to the 4G core following a service-based approach [26]: the 5G core is composed of multiple network functions, that provide mobility, authentication and routing support, that can be dynamically instantiated in data centers according to the load and traffic demands of the network. For example, while in LTE/EPC networks the control plane for the mobility of the user was handled by a single server (e.g., the MME), with the 5GC multiple network functions concur to offer the same set of services, but can be deployed in different data center locations and quickly turned off and on to decrease resource utilization. Moreover, the 5GC supports network slicing [30], i.e., the resources of the network can be split to serve different portions of traffic, that have different Quality of Service (QoS) requirements (e.g., IoT and mobile broadband traffic). The service-based 5GC architecture is an important enabler of network slicing in 5G, given that network functions can be provisioned dynamically to serve new network slices without the need to use separate servers,

as would happen with the EPC.

## 1.2 The Potentials and Challenges of mmWave Communications

Communication at mmWave frequencies has recently emerged as a possible Physical (PHY) layer technology to provide ultra-high data rates in mobile scenarios, and, as discussed in Sec. 1.1, has been integrated in 3GPP NR, which supports carrier frequencies up to 52.6 GHz in Release 15.[6] In particular, NR identifies as Frequency Range 2 (FR2) the spectrum between 24.25 GHz and 52.6 GHz.

The main potential of the mmWave bands is the availability of large chunks of untapped spectrum [32], which can be allocated to RAN deployments, making it possible to increase the physical layer capacity with respect to systems operating below 6 GHz [33]. The spectrum below 6 GHz, indeed, is heavily congested, and typical LTE [14] or IEEE 802.11a/b/g/n/ac/ax [34, 35] networks generally operate with a 20 MHz bandwidth per carrier. On the other hand, at mmWave frequencies it is possible to allocate very large contiguous chunks of spectrum, making it possible to exploit the 400 MHz bandwidth that is supported by 3GPP NR [10, 428].

The main drawbacks of the communication at such high frequencies, that, so far, have limited the applicability of this technology to mobile cellular networks, are [36]:

- the high isotropic propagation loss, which is proportional to the square of the carrier frequency, resulting in an over 30 dB higher loss at mmWaves than in conventional cellular systems at typical distances between the transmitter and the receiver [37]. However, it is possible to compensate this loss by using large antenna arrays and beamforming techniques, which increase the link budget by focusing the energy of the communication over narrow beams. Moreover, given the small wavelength at such high frequencies, it is feasible to pack many antenna elements in a small area: for example, at 30 GHz the wavelength $\lambda$ is approximately 1 cm, thus a rectangular array with 16 antennas (4 by 4) spaced by $\lambda/2$ would fit in a package with area smaller than a 2 cm by 2 cm square, and can be installed into a modern smartphone or Virtual Reality (VR) headset [419];

- the small wavelength, however, can be easily blocked by common materials, such as brick and mortar [38, 39], and by the human body as well [40]. For example, the authors in [41] have experimentally captured the impact of the hand and the human body blockage on the mmWave signal propagation from a hand-held device, showing that a median loss of 15 dB is incurred by the hand even in the most pessimistic scenario of a hard hand grip, and that the time-scales at which the mmWave signals are disrupted by blockage are on the order of a few hundreds of milliseconds or more. This makes mmWave communications more affected by shadowing than sub-6 GHz systems. Nonetheless, recent measurement campaigns have demonstrated that communications in Non-Line-of-Sight (NLOS) are possible, given a highly reflective surrounding environment [39], even though with a massive decrease in received power and thus available capacity;

- additionally, a reliable performance at mmWave frequencies is particularly challenging in truly mobile scenarios [42]. First, the Doppler spread is linear with the frequency, thus at mmWaves the rate at which the channel changes is much faster than in conventional systems. Then, the appearance of mobile obstacles can cause much wider variations in the received power, and, finally, the deployment of small cells (due to the high propagation loss) results in a relatively short interval between consecutive handovers.

---

[6]A Study Item is ongoing to determine whether NR will support higher frequency bands in future releases [31].

Consequently, the characteristics of the mmWave channel have an impact from the device design to the higher layers of the protocol stack. Therefore, it is necessary to introduce innovations throughout the whole protocol stack in order to support reliable connectivity when using these frequency bands.

From a device standpoint, the need to support large antenna arrays introduces tradeoffs between the performance, the energy efficiency and the cost, according to the beamforming architecture being considered. For example, while analog beamforming has a simple implementation and requires a single Radio Frequency (RF) chain (i.e., Analog to Digital Converters (ADCs), Digital to Analog Converters (DACs) and power amplifiers) for all the antenna elements, it can steer a beam in a single direction at any given time [43]. On the contrary, digital beamforming is complex in terms of modem and antenna design, since it requires an RF chain for each antenna element, but allows the simultaneous transmission and reception of signals in multiple directions [44]. A discussion on beamforming architectures (including hybrid beamforming) and their impact on the performance of the network will be given in Chapter 4.

At the physical layer, the issues are related to the choice of an efficient waveform, processing and coding design for such large bandwidths [17]. Moreover, the high datarates that will be supported by the large bandwidth require a fast and low-latency signal processing, which needs to handle a transport block (i.e., a packet in the NR physical layer) in a few microseconds. Additionally, novel techniques need to be developed in order to provide a practical and efficient channel estimation over large bandwidths, as in [45].

At the Medium Access Control (MAC) layer, the main challenge is given by the introduction of directional communications, which require a fine alignment of the beam at the two endpoints of the link. This calls for the introduction of a number of new operations at the MAC layer [46]. For example, the omnidirectional signals used in sub-6 GHz networks to broadcast synchronization information and perform the initial access would severely limit the coverage of mmWave deployments, thus directional signals are needed. Moreover, the two transceivers need to track the optimal set of beam pairs to be used as they move, and a beam management framework should take care of the recovery of the communication after a disruptive blockage event, with a re-alignment procedure or through a new initial access. Additionally, directionality has also an impact on interference management: on the one hand it generally reduces the interference, given that signals are no longer omnidirectional [47], on the other it may introduce unpredictable and strong interferers (due to the sidelobes, or the main lobe of a neighboring device) [48, 49]. Finally, the directionality may increase the complexity of channel sensing techniques, thus making accurate channel measurement more difficult and limiting the effectiveness of uncoordinated channel access schemes [50].

From a network perspective, the main issues are related to the efficient design and deployment of a mmWave network architecture that can enable a reliable and ubiquitous service [42]. This generally implies a multi-tier heterogeneous architecture, with an ultra-dense deployment, that introduce challenges in terms of (i) availability of fiber backhaul [51]; (ii) mobility management, e.g., how to manage and combine the beam management framework at the MAC layer with the handovers and updates in the serving base station [36]; (iii) energy efficiency, given that at such high densities the network load will be small and power saving techniques can be deployed [52].

Finally, from the transport and application layers, the traditional adaptation mechanisms (e.g., congestion control for TCP [53], or dynamic adaptation for video streaming [54]) need to be able to cope with the sudden fluctuations in available bandwidth at the lower layers, or with outages caused by extended blockages.

## 1.3 End-to-end Design and Evaluation of mmWave Cellular Networks

The goal of this thesis is to address some of the aforementioned challenges to improve the reliability, robustness and overall performance of mmWave systems, especially in the context of 5G cellular networks. In this regard, the principle that guides the research presented in this manuscript is the optimization of the *end-to-end* performance, i.e., of metrics such as throughput, latency, and reliability measured at a mobile device and at the remote server with which the device is communicating. Therefore, we will not focus on the evaluation of the mmWave radio access only, but consider this as part of a complex system that must be jointly studied, designed and improved.

The main rationale behind this choice is that the quality of service that the end users will experience in 5G and beyond networks with a mmWave RAN does not depend *only* on the characteristics of the access link (i.e., ultra-high capacity with an erratic channel behavior), but *also* on the interaction of the RAN with the other parts of the network, its architecture, the protocols that run on top of it, and the algorithms and optimizations that network operators and device manufacturers introduce. Therefore, the success or failure of this new technology will not depend only on the design of the wireless protocol stack (as in 3GPP NR), but also on how it will be deployed, on how higher layers (e.g., in the TCP/IP stack) will collaborate with it, and on how the system as a whole will be optimized.

Before this thesis, most of the studies on mmWave wireless communications were focused on the design of physical and/or link-layer solutions [17, 45], or on system-level evaluations of the achievable capacity in mmWave networks [37, 55, 56], with a few notable exceptions [46, 57]. MmWave systems, however, are just a part of the overall network, and thus we deem important to adopt a novel approach in this domain, which examines the system as a whole, with the complex and, in some cases, unpredictable interactions that may emerge from the combination of the highly variable mmWave channel and the rest of the network.

Therefore, this thesis adopts a system-level, end-to-end approach in the design and evaluation of three macro-areas, that will also be the main parts into which this manuscript is organized:

- **The Architecture: System Level Design of 5G mmWave Networks** − in this first part, we study how to design and efficiently deploy a mmWave network that provides reliable performance also in challenging mobility conditions. The novel contributions in this area are multiple. We propose to exploit an innovative multi-connectivity architecture to implement a reliable and efficient mobility management framework, that, thanks to the tight integration among different technologies in the RAN and new mobility procedures we introduced, improves the latency of the connection, reduces the throughput variability and enables a seamless service in mobile scenarios. Moreover, we study beam management in 3GPP NR cellular networks and for highly mobile scenarios with Unmanned Aerial Vehicles (UAVs), highlighting which are the main tradeoffs related to the performance of directional initial access and tracking according to the latest 3GPP specifications, and providing for the first time insights on how to design and deploy a network that needs to support directional beam management operations. Finally, we address the challenge of backhaul connectivity for ultra-dense mmWave deployments, by studying the performance of wireless backhaul solutions tightly integrated with the access network, as proposed by the 3GPP;

- **The Protocols: End-to-End and Cross-Layer Analysis of 5G mmWave Networks** − in the second part, we analyze how existing networking protocols perform in an end-to-end network with mmWave links in the radio access, and how to improve their performance,

and propose novel transport layer alternatives. Our analysis is the first to consider end-to-end performance (throughput and latency) in 3GPP scenarios with a mmWave RAN, and highlights how TCP suffers from the sudden changes in capacity that affect the mmWave physical layer much more than at sub-6 GHz. The congestion control algorithms, indeed, are too slow to react to an available datarate drop after a Line-of-Sight (LOS) to NLOS transition, and this causes latency spikes due to excessive buffering. Moreover, in some cases, a timeout is triggered, and TCP makes an inefficient usage of the resources after the recovery. Therefore, we propose some cross-layer strategies (i.e., an in-network proxy and a new congestion control algorithm for uplink flows) to improve the user experience when TCP is used as transport, and investigate possible replacements that exploit multi-connectivity at the transport layer, such as Multipath TCP (MPTCP) and a novel UDP- and network-coding-based protocol. Our results show that the proposed solutions improve the end-to-end performance of the network, making a case for the introduction of these optimizations in future 5G mmWave mobile networks and devices;

- **The Intelligence: Data-Driven 5G Networks Optimization** − in the third part, we introduce some strategies to integrate data-driven optimization in 5G networks, not necessarily focusing on mmWaves, but with approaches that can benefit the end users and the network operations independently of the frequency range considered. For this part, we use real-world datasets, provided by the Transport for London authority (i.e., the data collected by sensor traffic throughout the whole city) and a major U.S. telecommunications operator (with events from hundreds of base stations from the San Francisco and Palo Alto areas). We design a novel and practical strategy to deploy intelligence in 3GPP NR networks, exploiting edge controllers to collect data from the network, analyze it, and actuate optimizations based on the learned behavior of the network. We show the benefits of the proposed approach with two examples of applications, i.e., a self-organizing data-driven approach for the clustering of base stations under a certain controller, and the medium-term prediction of the number of users in the base stations of the network. Finally, we present an autonomous scaling of virtual network functions serving as mobility management entities in the London area, based on the expected number of handovers caused by vehicular traffic.

The tools that we will use for this thesis reflect the need to characterize the performance at the system level, and considering end-to-end metrics. We combine analysis, experiments and simulation, focusing mostly on the latter, for which we also contribute to the development of an open source mmWave module for the network simulator ns-3. This tool, which features the 3GPP channel model for mmWave frequencies and a 3GPP-like protocol stack, and benefits from the integration with the TCP/IP stack of ns-3, allowed us to run the first end-to-end performance evaluation campaigns of such complex systems, and to consider in the overall results the effects that emerge from the interactions of the different elements of the network.

We believe that the consistent methodology, the proposal of novel approaches and strategies, and the system-level approach that characterize this thesis make a valid and innovative contribution to the design of future wireless systems.

## 1.4   Thesis Structure

The rest of this thesis is organized into four main parts. The first, with Chapter 2 (which reports materials from our papers [390, 402, 406, 414, 419, 424]), describes the simulation tool we develop and use in the remainder of the thesis.

The second describes the research on architectures for mmWave networks. Chapter 3 (based on [388, 402]) introduces the multi-connectivity mobility management frameworks. Chapter 4 (derived from [394, 395, 416, 426]) discusses the beam management frameworks for mmWave mobile connectivity in cellular and UAV networks. Finally, Chapter 5 (which combines [399, 417, 418]) presents the contributions on wireless relaying architectures for 3GPP NR at mmWaves.

The third part illustrates our results and contributions on transport protocols. Chapter 6 (based on [389, 393, 404]) analyzes the performance of TCP on mmWave networks. Chapter 7 (which incorporates [407, 411]) proposes novel mechanisms to improve the performance of TCP, while Chapter 8 (derived from [389, 404, 413]) discusses possible alternatives to TCP.

Finally, the fourth part, with Chapter 9 (based on [392, 398, 420]), describes the proposals and evaluation of the integration of data-driven and intelligent strategies in 5G networks. Chapter 10 concludes the thesis, and suggests future research directions.

# Part I

# The Tool: End-to-End Performance Evaluation of 5G mmWave Networks

# 2

## Open Source, End-to-end Simulation of 5G mmWave Networks

## 2.1 Introduction

As discussed in Chapter 1, mmWave communications are emerging as a key technology in 5G cellular wireless systems due to their potential to achieve the massive throughputs required by future networks [6, 32, 39, 42, 58]. Due to the unique propagation characteristics of mmWave signals and the need to transmit in beams with much greater directionality than previously used in cellular systems, much of the recent work in mmWave communications has focused on channel modeling, beamforming and other physical layer procedures. However, as previously mentioned, the design of end-to-end (E2E) cellular systems that can fully exploit the high-throughput, low-latency capabilities of mmWave links will require innovations not only at the physical layer, but also across all layers of the communication protocol stack [36, 46, 57].

Discrete-event network simulators are fundamental and widely used tools for the development of new protocols and the analysis of complex networks. Importantly, most network simulators enable *full-stack simulation*, meaning that they model all layers of the protocol stack as well as applications running over the network. This full-stack capability will play a critical role in the development of 5G mmWave systems. The unique characteristics of the underlying mmWave channel have wide ranging effects throughout the protocol stack. For example, the use of highly directional beams increases the complexity of a number of basic MAC-layer procedures such as synchronization, control signaling, cell search and initial access, which in turn affect delay and robustness [46]. MmWave signals are also highly susceptible to blockage [32, 40, 59, 60], which results in high variability of the channel quality. This erratic behavior complicates the design of rate adaptation algorithms and signaling procedures, requiring advanced solutions for multi-connectivity, fast handover and connection re-establishment [61–63, 388]. New transport layer mechanisms may also be required in order to utilize the large capacity, when available, and to react promptly to rapid fading to avoid congestion [57, 64, 389, 404]. The need for ultra-low latency applications [20, 32, 65] may require solutions based on edge computing and distributed architectures that will determine a considerable departure from current cellular core network designs.

To better capture these design challenges, this chapter will describe the tools that were developed to support the E2E design and analyses that constitute the core contribution of this theses,

focusing in particular on the ns-3 mmWave module jointly developed by NYU Wireless and the University of Padova, which can be used to evaluate cross-layer and E2E performance. The ns–3 mmWave module was first presented in [66,67]. The content of this chapter was introduced in [390, 402, 406, 414, 419, 424]. All the code developed for the simulations in this thesis is open source, and links to the repositories will be provided.

This mmWave simulation tool is developed as a new module within the widely used ns–3 network simulator [68]. ns–3 is an open-source platform, that currently implements a wide range of protocols in C++, making it useful for cross-layer design and analysis. The new mmWave module presented here is based on the architecture and design patterns of the LTE LENA module [69,70] and implements all the necessary Service Access Points (SAPs) needed to leverage the robust suite of LTE/EPC protocols provided by LENA. The code is highly modular and customizable to help researchers to design and test novel 5G protocols.

The rest of the chapter is organized as follows. In Sec. 2.2, we describe the main challenges related to the design of a mmWave cellular network simulator. Then, in Sec. 2.3, we introduce ns-3, the network simulator on which our mmWave module is developed, and in Sec. 2.4 we present the overall architecture of the mmWave module. We then take a closer look at each component, starting with the suite of MIMO channel models in Sec. 2.5. In addition to an implementation of the 3GPP "above 6 GHz" model [71], several custom channel models are also provided. Sec. 2.6 discusses the features of the OFDM-based PHY layer, which has a customizable frame structure for evaluating different numerologies and parameters. In Sec. 2.7, we provide a MAC-layer discussion that includes our proposed flexible/variable Transmission Time Interval (TTI) Time Division Multiple Access (TDMA) MAC scheme, which is supported by several schedulers, and a carrier aggregation implementation. Sec. 2.8 presents the enhancements that we introduced to the LTE Radio Link Control (RLC) layer. The dual-connectivity architecture is reported in Sec. 2.9. In Sec. 2.10, we show how the module can be used for cross-layer evaluation of multi-user cellular networks through a number of representative examples, and provide pointers to a large set of general results that have been obtained so far with this module. The integration of native Linux Transmission Control Protocol (TCP) implementations, performed through the ns–3 Direct Code Execution (DCE) framework, is discussed in Sec. 2.11. In Sec. 2.12, we provide details on our future plans for the simulator and suggest possible research areas in which it could be used. Finally, we conclude this tutorial paper in Sec. 2.13.

## 2.2 Potentials and Challenges of System-level Simulations of mmWave Networks

An end-to-end network simulator for mmWave cellular networks is an invaluable tool that can help address these challenges by allowing the evaluation of the impact of the channel and of the PHY layer technology on the whole protocol stack. However, given the characteristics of mmWave communications described in the previous paragraphs, in order to have accurate results it is of paramount importance to model in detail the behavior of the different elements that interact in a cellular system. In the following paragraphs we will introduce and discuss some of the most important elements that need to be considered when designing a mmWave cellular system simulation, and show how they depend on one another:

- The channel model is the fundamental component of every wireless simulation. Given the harsh propagation conditions at mmWaves, the channel is one the main elements that affect the end-to-end network performance. Firstly, it has to account for the different LOS and NLOS states for the propagation loss and the fading [71]. Moreover, beamforming

14

should be applied on top of the channel to accurately model directional transmissions, which have an impact on the link budget, the interference, and the control procedures. Finally, the Doppler effect is particularly relevant at mmWave frequencies, especially with high mobility [42]. An important consideration related to the channel model is the trade off between the accuracy and the computational complexity: very accurate models that require the computation of the complete channel matrix are usually also computationally intensive [406, 415, 424].

- The users' mobility and the network deployment have an important impact on the communication performance, intertwined with that of the channel model. Given the small range of the mmWave cells, the deployment will be dense and will require frequent access point updates, which should be simulated for a realistic performance assessment [409]. Moreover, mobility affects the performance of beam tracking algorithms [62]. Therefore, when simulating a mmWave network it is important to use realistic deployments and mobility models.

- The level of detail when modeling the protocol stack of the mmWave links and of the end devices is another important parameter for network simulations. A simplified model of the protocol stack can be accurate enough for studies that involve limited interplay between different layers, but cannot capture the behaviors that emerge from complex interactions among them, and therefore may not be sufficient to generate realistic results for end-to-end performance evaluations. For example, at mmWave frequencies, it has been shown that the channel behavior has an impact on the TCP performance [57, 393, 404], therefore a model of the TCP/IP stack is needed when analyzing the data rate that an application can reach in an end-to-end mmWave network.

To the best of our knowledge, when this project started, there were no open source simulators capable of thoroughly modeling the mmWave channel along with the cellular network protocol stack as well as other protocols (e.g., the TCP/IP stack), realistic scenarios and mobility. There exists an ns–3-based simulator for IEEE 802.11ad in the 60 GHz band [72–76], which however cannot be used to simulate cellular and 3GPP-like scenarios. Other papers [77–80] report results from system level simulations, with custom (often MATLAB-based and not publicly available) simulators which are not able to capture the complexity of the whole stack with a very high level of detail. This is what originally motivated us to develop an open source cellular mmWave module for the ns–3 simulator. The ns-3 community has then picked up the development of 5G cellular modules. For example, the authors in [81] extend the core network model to implement the 5G NFV architecture. The module described in [19, 82, 83], instead, enhances the NYU/UNIPD mmWave and the ns-3 LTE modules by re-implementing the PHY and MAC layers to be compliant with the numerologies of 3GPP NR. However, it lacks some features, such as the support for dual connectivity, which are instead supported by the NYU/UNIPD mmWave module which we will describe in the following sections.

## 2.3  ns–3

The ns–3 discrete-event network simulator [68, 84] is a very powerful tool available to communication and networking researchers for developing new protocols and analyzing complex systems. It is the successor to ns–2, a well-tested tool that has been in use by the networking community for over a decade in the design and validation of network protocols. ns–3 is open source, and can be downloaded from the website of the project.[1] An active community of researchers from both

---

[1] http://www.nsnam.org

industry and academia has enriched the basic core of the simulator with several modules, and ns–3 can now be used to simulate a wide variety of wireless and wired networks, protocols and algorithms. There is a complete documentation[2] on the models in the ns–3 website, in terms of both the design of the models and what a user can do with the models. Moreover, a complete tutorial on how to install ns–3, set up ns–3 scenarios and topologies, handle the collection of statistics and log useful messages is provided in the documentation.[3] The tutorial is a good starting point for a researcher who approaches ns–3 for the first time.

The ns–3 simulator is organized into multiple folders. The `src` folder provides a collection of C++ classes, which implement a wide range of modular simulation models and network protocols. The different modules can be aggregated and instantiated to build diverse simulated network scenarios, making ns-3 especially useful for cross-layer design and analysis. The modularity and use of object-oriented design patterns also allow for new algorithms to be easily incorporated into the network stack and experimented with. Each module is itself organized into multiple subfolders, which contain the documentation and the source code of the model itself, the helpers, the examples and the tests. The helpers associated with each model have a very important role. They are classes which hide to the final user the complexity involved in setting up a complete scenario, for example by automatically assigning IP addresses, or connecting the different classes of a protocol stack. The `build` folder contains the binaries of the simulator. Finally, the `scratch` folder is a special folder in which scripts with examples and scenarios can be built on-the-fly.

Besides the core module, which provides the basic structure of the simulator, there are modules for networking protocols (e.g., the TCP/IP stack protocols [85]), wireless protocols (LTE [69], Wi-Fi [86], WiMAX [87]), routing algorithms [88], mobility, embedding obstacles and buildings in the simulation scenarios, and data collection. All the modules are listed in the model library.[4]

In the following sections, we will describe in detail the mmWave module for ns–3, following the same approach which is used for the other ns–3 modules. We will first describe the model in terms of implementation of the different components of a mmWave cellular network and protocol stack, and then the examples and scenarios that can be simulated with it and how they can be set up.

## 2.4    mmWave Module Overview

The ns–3 mmWave module is designed to perform end-to-end simulations of 3GPP-style cellular networks. Fig. 2.1 depicts a high level overview of the different components of the protocol stack and the end-to-end network architecture that this module makes it possible to simulate. As shown in Fig. 2.2, the architecture builds upon the ns–3 LTE module (LENA) [69, 70]. It leverages the detailed implementation of LTE/EPC protocols, and implements custom PHY and MAC layers. Additionally, it is possible to connect the module to a patched version of Direct Code Execution [89], a tool that allows the Linux stack TCP/IP implementation to run as the TCP/IP stack of ns–3 nodes, as well as to execute POSIX socket-based applications (i.e., wget, iPerf, etc). Fig. 2.2 also depicts the high-level composition of the `MmWaveEnbNetDevice` and `MmWaveUeNetDevice` classes, which represent the mmWave eNB/gNB[5] and UE radio stacks.

The ns–3 mmWave module also includes a `McUeNetDevice`, which is a `NetDevice` with a dual stack (LTE and mmWave), i.e., a device capable of connecting to both technologies. Moreover,

---

**Figure 2.1:** Representation of the E2E protocol stack and architecture of the ns-3 mmWave module, including the public Internet, the core network and the RAN.

**Figure 2.2:** Simplified UML class diagram for the end-to-end mmWave module.

the protocol stack can be configured to support carrier aggregation. More details will be given in Sec. 2.9.

The `MmWaveEnbMac` and `MmWaveUeMac` MAC layer classes implement the LTE module Service Access Point (SAP) *provider* and *user* interfaces, which enable the inter-operation with the LTE RLC layer. Support for RLC Transparent Mode (TM), Saturation Mode (SM), Unacknowledged Mode (UM), Acknowledged Mode (AM) is built into the MAC and scheduler classes (i.e., `MmWaveMacScheduler` and derived classes). The MAC scheduler also implements a SAP for configuration at the LTE Radio Resource Control (RRC) layer (`LteEnbRrc`). Hence, every component required to establish Evolved Packet Core (EPC) connectivity is available.

The `MmWavePhy` classes handle directional transmission and reception of the Downlink (DL) and Uplink (UL) data and control channels based on control messages from the MAC layer. Similar to the LTE module, each PHY instance communicates over the channel (i.e., `SpectrumChannel`) via an instance of the `MmWaveSpectrumPhy` class, which is shared for both the DL and the UL (since our design of the mmWave PHY layer is based on Time Division Duplexing (TDD), as detailed in Sec. 2.6.1). Instances of `MmWaveSpectrumPhy` encapsulate all PHY-layer models: interference calculation (`MmWaveInterference`), Signal to Interference plus Noise Ratio (SINR) calculation (`MmWaveSinrChunkProcessor`), the Mutual Information (MI)-based error model (`MmWaveMiErrorModel`), which computes the packet error probability, as well as the HARQ PHY-layer entity (`MmWaveHarqPhy`) to perform soft combining.

Since the structure, high-level functions and naming scheme of each class closely follow the LTE LENA module, the reader is also referred to the LENA project documentation for more information [90].

## 2.5 Channel and MIMO Modeling

### 2.5.1 Channel Models

The ns–3 mmWave module allows the user to choose among different channel models, which provide a trade-off between computational complexity, flexibility and accuracy of the results. The most flexible and detailed channel model is the one described in detail in [406], which is based on the official 3GPP channel model for the 6-100 GHz frequency band [71]. It accounts also for spatial consistency of mobility-based simulations and provides a random blockage model, as well as the modeling of outdoor to indoor communications. The second model is based on traces from measurements or third-party ray-tracing software. This makes the channel model detailed and realistic, but constrains the simulation to limited measurements/ray-tracing routes. The third is the statistical channel model introduced in [66] and based on MATLAB traces, which makes the computation less demanding, but is available only for the 28 and 73 GHz frequencies. In the following paragraphs we will provide architectural details of all the available channel models.

#### 3GPP Statistical Channel Model

The 3GPP model for the 6-100 GHz band, described in [71], is applicable for bandwidths up to 10% of the carrier frequency and accounts for mobility. It provides several optional features that can be plugged into the basic model, in order to simulate, for example, spatial consistency (i.e., the radio environment conditions of close-by users are correlated) and random blockage. The model defines different scenarios, which describe different possible cellular network deployments: urban (with macrocells and microcells), rural and indoor.

**Figure 2.3:** Typical realizations of the 3GPP pathloss model described in [71]. We consider three outdoor scenarios (UMi stands for Urban Micro, UMa for Urban Macro, and RMa for Rural Macro) and two indoor scenarios (InH stands for Indoor Hotspot, either in the office or in a shopping mall). The realizations for each single 3GPP scenario differ because of the channel condition between the UE and the base station: it can be a LOS, NLOS or O2I channel. Moreover, 3GPP specifies an additional set of equations for the NLOS channel in the UMi, UMa and InH-Office scenarios, which are marked as optional (opt).

**Pathloss:** The pathloss of the propagation channel is implemented in the `MmWave3gppPropagationLossModel` class. The model provides a statistical LOS/NLOS condition characterization, as well as pathloss computation considering outdoor to indoor penetration loss, as described in [71, Sec. 7.4]. The `MmWave3gppBuildingPropagationLossModel` class, instead, determines the LOS condition according to the relative position of the UE and the eNB and to the presence of buildings or obstacles in the scenario. These classes also optionally apply an additional shadowing component to the pathloss. For a moving UE, the shadowing is correlated in space. Given the distance $\Delta d_{2D} > 0$ on the horizontal plane from the last position in which the shadowing was computed, the exponential correlation parameter is computed as $R(\Delta d_{2D}) = e^{-\Delta d_{2D}/d_{cor}}$, where $d_{cor}$ is the correlation distance. In our implementation, pathloss and shadowing (if enabled) are updated at every transmission. Fig. 2.3 shows the pathloss in dB for the 3D distance from the smallest value supported in each scenario to 1 km for outdoor and 100 m for indoor.

**Small-scale fading:** The small-scale fading model is implemented in the `MmWave3gppChannel` class, and follows the step by step approach of [71, Sec. 7.5]. Small-scale fading is the bottleneck of this channel model implementation, since it is very detailed and computationally demanding. The fading is generated following the 3D statistical spatial approach originally proposed in [91]. The channel is described by a channel matrix $\mathbf{H}(t, f)$, where $t$ is the time and $f$ is the frequency, of size $U \times S$, where $U$ and $S$ are the number of antennas at the receiver and the transmitter. Each entry depends on $N \leq 20$ different multipath components, called *clusters*, which have different delays and received powers, according to an exponential power delay profile. A cluster is itself a combination of $M = 20$ *rays*, each with a slightly different arrival and departure angle in the vertical and horizontal planes.

The `MmWave3gppChannel` class has a method that generates the channel matrix, and stores the coefficient for each transmit element $s$, receive element $u$ and cluster $n$ in a data structure, that can be accessed by other methods in order to update the channel matrix or compute the beamforming gain. We introduced some assumptions with respect to the 3GPP model, in order to decrease the computational overhead introduced by the high level of detail of the channel. For example, we consider only antennas with vertical polarization, and the speed-dependent Doppler effect is not computed for each ray, but only for the central angle of each cluster. Further details on this implementation are given in [406].

**Spatial consistency:** The basic channel model described in the previous paragraphs can be used for drop-based simulations with limited mobility, i.e., for UEs that move in an area

in which the channel is very correlated and the fading parameters do not change. However, for simulations in which mobility is an important factor, the spatial consistency of the channel throughout the path on which the UE moves can be simulated by enabling this option in the `MmWave3gppChannel` class. In the current implementation, we support spatial consistency with Procedure A of [71, Sec 7.6.3.2] for both LOS and NLOS communications. It is possible to set the period of update $t_{PER}$, and every $t_{PER}$ the cluster delays, powers and departure and arrival angles are updated with a transformation that accounts for the speed of the UE and for the distance traveled on the horizontal plane.

**Blockage:** This optional feature can be used to model the attenuation in certain clusters, according to their angle of arrival. The attenuation can be caused by the human body that holds the UE, or by external elements such as for example cars, other human bodies, trees. The blockage model is implemented in the `MmWave3gppChannel` class and can be optionally activated. In our implementation we consider blockage model A, which only distinguishes between self-blocking and non-self-blocking, and is generic and computationally efficient [71]. In particular, this model randomly generates $K + 1$ blocking regions, one for self-blocking, with different parameters according to the orientation of the UE (i.e., portrait or landscape mode), and $K$ for non-self-blocking. The attenuation is 30 dB for self-blocking, whereas it depends on the scenario and on the horizontal and vertical angles of arrival for non-self-blocking. Moreover, the blocking of a certain cluster is correlated in both space and time, according to the UE mobility, the blocker speed and the simulation scenario. Notice that, if both the blockage and the spatial consistency options are used, then the update of the channel with both features is synchronized, i.e., the cluster blockage is updated before the channel coefficients are recomputed with the spatial consistency procedure.

### Ray-tracing or Measurement Trace Model

`MmWaveChannelRaytracing` uses software-generated or measurement traces to model the channel in ns–3, for pathloss and fading. The trace samples need to contain the number of paths and the propagation loss, delay, angle of arrival and angle of departure for each path. The following trace files have been tested in our implementation and are available in `mmwave/model/Raytracing/`.

**Ray-tracing:** Any ray-tracing software (e.g., WinProp [92]) can be used to generate the channel information for a specific route. This means that the simulation scenario must be chosen *a priori*, and cannot be random since it has to be given as input to the ray-tracing software. An example of ray-tracing route[6] is shown in Fig. 2.4b.

**QuaDRiGa:** The Quasi Deterministic Radio Channel Generator model [94], supports consistent user mobility and massive MIMO at several frequencies (10, 28, 43, 60, 82 GHz). It also adds some time evolution characterization on top of the statistical channel to capture user mobility, which makes it suitable for system level simulations.

### NYU Statistical Model

This channel model is based on the approach described in [37] and implemented in our previous work [66]. A MATLAB implementation of the same channel model is also available in [95–97]. It provides two pathloss models, which differ in how they capture the LOS/NLOS condition. The first, `MmWavePropagationLossModel`, is based on a statistical characterization of the LOS state, while the second, `BuildingsObstaclePropagationLossModel`, leverages the ns–3 buildings module in order to decide whether there is an obstacle between the UE and the eNB or

---

[6]The ray tracing data was provided by the Communication Systems and Networks Group, University of Bristol, UK [57, 93].

**(a)** 3GPP statistical channel model  **(b)** Ray-tracing Trace Model  **(c)** NYU Statistical Model

**Figure 2.4:** Example of average SINR plots for the three channel models.

not. In particular, it is possible to deploy – deterministically or randomly – objects of different sizes to mimic humans, cars, and buildings. A virtual line is drawn between the transmitter and the receiver: If this line intersects any object, the state is NLOS, otherwise it is LOS. In both classes, once the channel state is selected, the propagation loss is computed as in [37].

**Channel configuration:** Since the channel matrices and optimal beamforming vectors do not depend on the distance between the UE and the eNB, they are pre-generated in MATLAB to reduce the computational overhead in ns–3. At the beginning of each simulation we load 100 instances of the spatial signature matrices, along with the beamforming vectors. Moreover, in order to simulate realistic channels with large-scale fading, the channel matrices are updated periodically and independently (*block fading*). Currently, no results are available for modeling how the large-scale statistics of the mmWave channel change over time for a mobile user, thus it should be noted that the accuracy of this method is not verified at this time. The matrix update can take place at some fixed intervals, specified by the `LongTermUpdatePeriod` attribute of the `MmWaveBeamforming` class. The small-scale fading, instead, is calculated at every transmission, where we obtain the speed of the user directly from the mobility model. The remaining parameters that depend on the environment are assumed to be constant over the entire simulation time.

**Semi-empirical feature:** Finally, as shown in Fig. 2.4c, the soft transition between LOS/NLOS conditions can be modeled in a "semi-empirical" fashion, meaning that we overlay the statistical channel with blockage measurements performed in our lab [98]: Waving a hand in front of the receiver (hand blockage), walking between the transmitter and the receiver (human blockage), and placing a metal plate between the transmitter and the receiver to emulate an obstacle, like a car or a building.

## 2.5.2 Beamforming Gain

For the long-term statistical channel model, the beamforming vectors are directly loaded from MATLAB generated files. For the other channel models, two methods are implemented to compute beamforming vectors, i.e., the *long-term covariance matrix method* and the *beam search method*. Currently, the only available beamforming architecture for data transmission is *analog*, meaning that devices can transmit or receive in only one direction at a time. As part of our future work, we plan to integrate *hybrid* and *digital* transceiver designs.

In the long-term covariance matrix method, we assume that the transmitter estimates the

22

spatial correlation matrix $\mathbf{Q}_{tx} = \mathbb{E}[\mathbf{H}^\dagger(t, f)\mathbf{H}(t, f)]$, where the expectation is taken over the frequencies and some interval of time. An analogous operation is done for the receiver. In practice, the Transmitter (TX) and Receiver (RX) would estimate the spatial covariance matrix from reference or synchronization signals and beam scanning. Estimation of this covariance matrix is discussed in [99]. We do not, however, model the covariance estimation directly; instead we simply assume that the TX and RX know the correct long-term channel with some configurable delay. Beamforming vectors can then be computed from the maximal eigenvectors of the covariance matrices [100]. A computationally simple procedure is to use the power method [101]. The algorithm selects a random initial beamforming vector and iteratively multiplies it with the spatial correlation matrix $\mathbf{Q}_{tx}$, normalizing the results at each iteration. Finally, the output will converge to the correct eigenvector. The computation for the receiver is done in the same way, starting from $\mathbf{Q}_{rx} = \mathbb{E}[\mathbf{H}(t, f)\mathbf{H}^\dagger(t, f)]$.

In the beam search method, we assume that the TX and the RX scan a discrete number of beams from a pre-designed codebook [102]. Codebook design is discussed in detail in [103]. The beamforming vector is selected as the one with the highest power, possibly with some time-averaging.

Additionally, in [419], we introduced the possibility of deploying multiple antenna arrays at each base station and UE, thus allowing a sectorized deployment. Moreover, we added the possibility of simulating non-isotropic antenna patterns for each single antenna elements in each array, following the 3GPP specifications.

### 2.5.3 Interference

MmWave systems may be interference- or power-limited. Albeit potentially less significant for directional mmWave signals, which are generally assumed to be power-limited, there are still some cases (i.e., high deployment density) where interference is non-negligible [47, 104]. Additionally, although intra-cell interference (i.e., from devices of the same cell) can be neglected in TDMA or Frequency Division Multiple Access (FDMA) operation, it does need to be explicitly calculated in the case of Spatial Division Multiple Access (SDMA)/Multi-User MIMO, where users are multiplexed in the spatial dimension but operate in the same time-frequency resources. Therefore, we propose an interference computation scheme that takes into account the beamforming vectors associated with each link.

As an example, we compute the SINR between nodes $eNB_1$ and $UE_1$ in the presence of an interferer, $eNB_2$. To do so, we first need to obtain the beamforming gains associated with both the desired and interfering signals, i.e.,

$$
\begin{aligned}
G_{11} &= |\mathbf{w}_{rx_{11}}^\dagger \mathbf{H}(t, f)_{11} \mathbf{w}_{tx_{11}}|^2, \\
G_{21} &= |\mathbf{w}_{rx_{11}}^\dagger \mathbf{H}(t, f)_{21} \mathbf{w}_{tx_{22}}|^2,
\end{aligned}
\tag{2.1}
$$

where $\mathbf{w}_{rx_{i,j}}$ is the beamforming vector of receiver $i$ towards transmitter $j$, and $\mathbf{w}_{tx_{i,j}}$ is the beamforming vector of transmitter $i$ towards receiver $j$. The SINR is then computed as:

$$
SINR_{11} = \frac{\frac{P_{Tx,11}}{PL_{11}} G_{11}}{\frac{P_{Tx,22}}{PL_{21}} G_{21} + BW \times N_0},
\tag{2.2}
$$

where $P_{Tx,ii}$ is the transmit power of $eNB_i$, $PL_{ij}$ is the pathloss between $eNB_i$ and $UE_j$, and $BW \times N_0$ is the thermal noise.

### 2.5.4 Error Model

The mmWave module exploits the error model introduced in the ns-3 LTE LENA project, which follows a link abstraction technique for simulating Transport Block (TB) errors in the downlink of an LTE system. In a nutshell, the model described in [105] defines an accurate and lightweight procedure for the computation of the residual errors after PHY layer processing. This is achieved by exploiting:

- Mutual Information-based multi-carrier compression metrics to derive a unique SINR value of the channel, known as *effective SINR*, which is represented as $\gamma_i$ in Eq. (2.3), and

- Link-Level performance curves obtained with a MATLAB bit-level LTE PHY simulator [106], which have been used to match a mathematical approximation of the Block Error Rate (BLER), as reported in Eq. (2.3).

The ultimate goal is to let the receiver derive the error probability of each TB to determine whether the packet can be decoded or not. Because each TB can be composed of multiple Code Blocks (CBs), whose size depends on the channel capacity, the BLER can be formulated as follows:

$$C_{BLER,i}(\gamma_i) = \frac{1}{2}\left[1 - \mathrm{erf}\left(\frac{\gamma_i - b_{C_{SIZE},MCS}}{\sqrt{2}c_{C_{SIZE},MCS}}\right)\right],\tag{2.3}$$

where $\gamma_i$ corresponds to the mean mutual information per coded bit of code block $i$, as explained earlier, and $b_{C_{SIZE},MCS}$ and $c_{C_{SIZE},MCS}$ represent the mean and standard deviation of the Gaussian cumulative distribution, respectively, which have been obtained from the link level performance curves mentioned above. Finally, the TB block error rate is given by:

$$T_{BLER} = 1 - \prod_{i=1}^{C}(1 - C_{BLER,i}(\gamma_i)).\tag{2.4}$$

### 2.5.5 Examples

An example of SINR plots for the three channel models was presented in Fig. 2.4. An example related to the setup of the channel model can be found in the `examples` folder of the mmWave module, in the file `mmwave-3gpp-channel- example.cc`.

Fig. 2.4a shows an example of a rural scenario with an eNB at coordinates $(0,0)$ m and at the height of 35 m, and a UE in position $(100,0)$ m, at the height of 1.5 m and moving at 1 m/s along the y axis, maintaining LOS connectivity. The channel is updated consistently every 100 ms. The top figure shows the SINR when the Beamforming (BF) vector is updated with the long-term covariance matrix method, while in the bottom one it is updated with the beam search method. Notice that the current implementation of the beam search method uses a fixed elevation angle of 90 degrees and sweeps only the horizontal plane. Therefore, the beam search method cannot align with the LOS cluster and the power is reduced by 20 dB. Moreover, after enabling the blockage model, the SINR achieved by the long-term covariance matrix method dropped by 20 dB when the LOS cluster was blocked. However, the beam search method experienced less blockage impact, as it did not align with the LOS cluster. In the other case, without update, the BF vector is computed at $t = 0$ s but never updated, and this causes the SINR to drop as the UE moves. Comparing the blue and black curves, it is possible to observe that for the first 20 s the performance with and without BF update is similar, because of the consistency of the channel and of the low mobility of the UE, but after $t = 20$ s the SINR without update degrades by nearly 30 dB. The last observation is that the long-term covariance matrix method finds the

optimal BF vector whenever the channel is changed, therefore the SINR is very stable. On the other hand, the beam search method shows an SINR drop after 20 s even with update, because when the UE moves both the UE and the eNB are unable to optimally adapt the BF vector and just select one of the available sectors.

Fig. 2.4b plots the average SINR of a ray tracing channel indicating both LOS intervals and NLOS channel states. The ray tracing data contains 5000 samples along a 500 meter route. The SINR has a sudden change when the channel state switches. We note that the SINR curve within LOS is relatively stable, whereas more random variations are introduced for NLOS.

Finally, Fig. 2.4c shows the average SINR trace generated with the NYU channel model [37] in two cases, a walking user blocked by a building (top) or by other humans (bottom). The main difference is that, with buildings, the link capacity drops rapidly and the blocking interval lasts seconds; on the other hand, with humans, the channel deteriorates slowly and the blockage lasts only for a short interval. From the top figure, we can observe that with soft LOS/NLOS transition enabled, the SINR curve changes less suddenly when the channel condition switches. In the bottom graph, three human blockage events, at 1, 4 and 7 seconds, are added on top of the statistical channel.

## 2.6   Physical Layer

In this section, we discuss the key features of the mmWave PHY layer. Specifically, we have implemented a TDD frame and subframe structure, which has similarities to TDD-LTE, but allows for more flexible allocation and placement of control and data channels within the subframe and is suitable for the *variable TTI* MAC scheme described in Sec. 2.7. Moreover, we implemented an error model and HARQ model based on those in LENA, but compatible with our custom mmWave PHY and numerology (for instance, they support larger TB and codeword sizes as well as multi-process stop-and-wait HARQ for both DL and UL).

### 2.6.1   Frame Structure

It is widely contended that 5G mmWave systems will target Time Division Duplex (TDD) operation because it offers improved utilization of wider bandwidths and the opportunity to take advantage of channel reciprocity for channel estimation [32, 42, 61, 107, 108]. In addition, shorter symbol periods and/or slot lengths have been proposed in order to reduce radio link latency [109–111]. The ns–3 mmWave module therefore implements a TDD frame structure which is designed to be configurable and supports short slots in the hope that it will be useful for evaluating different potential designs and numerologies. These parameters, shown in Table 2.1, are accessible through the attributes of the common `MmwavePhyMacCommon` class, which stores all user-defined configuration parameters used by the PHY and MAC classes. Examples related to the setup of the PHY layer parameters can be found in the `mmwave-tdma.cc` and `mmwave-epc-tdma.cc` files.

The frame and subframe structures share some similarities with LTE in that each frame is subdivided into a number of subframes of fixed length [112]. However, in this case, the user is allowed to specify the subframe length in multiples of OFDM symbols[7]. Within each subframe,

---

[7]Though many waveforms are being considered for 5G systems, OFDM is still viewed as a possible candidate. In [107,113], Verizon and the consortium led by Korea Telecom propose a frame structure and OFDM numerology. However, this is still under debate in 3GPP [114]. We naturally chose to adopt OFDM, at least initially, for the mmWave module, which allows us to leverage the existing PHY models derived for OFDM from the LTE LENA module. As soon as the 3GPP NR will be standardized, the protocol stack in our module can be easily adapted to the updated parameters.

| Parameter Name | Default Value | Description |
|---|---|---|
| SubframePerFrame | 10 | Number of subframes in one frame |
| SubframeLength | 100 | Length of one subframe in $\mu s$ |
| SymbolsPerSubframe | 24 | Number of OFDM symbols per slot |
| SymbolLength | 4.16 | Length of one OFDM symbol in $\mu s$ |
| NumSubbands | 72 | Number of subbands |
| SubbandWidth | 13.89 | Width of one subband in MHz |
| SubcarriersPerSubband | 48 | Number of subcarriers in each subband |
| CenterFreq | [6-100] | Possible carrier frequencies in GHz* |
| NumRefScPerSymbol | 864 (25% total) | Reference subcarriers per symbol |
| NumDlCtrlSymbols | 1 | Downlink control symbols per subframe |
| NumUlCtrlSymbols | 1 | Uplink control symbols per subframe |
| GuardPeriod | 4.16 | Guard period for UL-to-DL mode switching in $\mu s$ |
| MacPhyDataLatency | 2 | Subframes between MAC scheduling request and scheduled subframe |
| PhyMacDataLatency | 2 | Subframes between TB reception at PHY and delivery to MAC |
| NumHarqProcesses | 20 | Number of HARQ processes for both DL and UL |

**Table 2.1:** Parameters for configuring the mmWave PHY.
*The NYU channel model [37] supports only 28 and 73 GHz.



**Figure 2.5:** Proposed mmWave frame structure.

a variable number of symbols can be assigned by the MAC scheduler and designated for either control or data channel transmission. The MAC entity therefore has full control over multiplexing of physical channels within the subframe, as discussed in Sec. 2.7. Furthermore, each variable-length time-domain data slot can be allocated by the scheduler to different users for either the uplink or the downlink.

Fig. 2.5 shows an example of frame structure with the numerology taken from our proposed design in [110]. Each frame of length 1 ms is split in time into 10 subframes, each of duration $100\,\mu s$, representing 24 symbols of approximately $4.16\,\mu s$. In this particular scheme, the downlink and uplink control channels are always fixed in the first and last symbol of the subframe, respectively. A switching guard period of one symbol period is introduced each time the direction changes from UL to DL. In the frequency domain, the entire bandwidth of 1 GHz is divided into 72 subbands of width 13.89 MHz, each composed of 48 subcarriers. It is possible to assign

UE data to each of these subbands, as is done with Orthogonal Frequency-Division Multiple Access (OFDMA) in LTE, however only TDMA operation is currently supported for reasons we shall explain shortly.

### 2.6.2 PHY Transmission and Reception

The `MmWaveEnbPhy` and `MmWaveUePhy` classes model the physical layer for the mmWave eNB and the UE, respectively, and encapsulate similar functionalities as the `LtePhy` classes from the LTE module. Broadly, these objects (i) handle the transmission and reception of physical control and data channels (analogous to the Physical Downlonk Control Channel (PDCCH)/Physical Uplink Control Channel (PUCCH) and PDSCH/Physical Uplink Shared Channel (PUSCH) channels of LTE), (ii) simulate the start and the end of frames, subframes and slots, and (iii) deliver received and successfully decoded data and control packets to the MAC layer.

In the `MmWaveEnbPhy` and `MmWaveUePhy` classes, calls to `StartSubFrame()` and `EndSubFrame()` are scheduled at fixed periods, based on the user-specified subframe length, to mark the start and end of each subframe. The timing of variable-TTI slots, controlled by scheduling the `StartSlot()` and `EndSlot()` methods, is dynamically configured by the MAC via the MAC-PHY SAP method `SetSfAllocInfo()`, which enqueues an `SfAllocInfo` allocation element for some future subframe index specified by the MAC. A *subframe indication* to the MAC layer triggers the scheduler at the beginning of each subframe to allocate a future subframe. For the UE PHY, `SfAllocInfo` objects are populated after reception of Downlink Control Information (DCI) messages. At the beginning of each subframe, the current subframe allocation scheme is dequeued, which contains a variable number of `SlotAllocInfo` objects. These, in turn, specify contiguous ranges of OFDM symbol indices occupied by a given slot, along with the designation as either *DL* or *UL* and control (*CTRL*) or data (*DATA*).

The data packets and the control messages generated by the MAC are mapped to a specific subframe and slot index in the *packet burst map* and *control message map*, respectively. Presently, in our custom subframe design, certain control messages which must be decoded by all UEs, such as the DCIs, are always transmitted in fixed PDCCH/PUCCH symbols at the first and last symbol of the subframe, but this static mapping can be easily changed by the user[8]. Other UE-specific control and data packets are recalled at the beginning of each allocated TDMA data slot and are transmitted to the intended device.

To initiate transmission of a data slot, the eNB PHY first calls `AntennaArrayModel::Change-BeamformingVector()` to update the transmit and receive beamforming vectors for both the eNB and the UE. In the case of control slots, no beamforming update is applied since we currently assume an "ideal" control channel. For both DL and UL transmissions, either the `MmWaveSpectrumPhy` method `StartTxDataFrame()` or `StartTxCtrlFrame()` is called to transmit a data or control slot, respectively. The functions of `MmWaveSpectrumPhy`, which is similar to the corresponding LENA class, are as follows. After the reception of data packets, the PHY layer calculates the SINR of the received signal in each subband, taking into account the path loss, MIMO beamforming gains and frequency-selective fading. This triggers the generation of Channel Quality Information (CQI) reports, which are fed back to the base station in either UL data or control slots. The error model instance is also called to probabilistically compute whether a packet should be dropped by the receiver based on the SINR and, in the case of an HARQ retransmission, any soft bits that have been accumulated in the PHY HARQ entity (see

---

[8]As in [109, 110], we assume either FDMA or SDMA-based multiple access in the control regions. However, we do not currently model these modulation schemes nor the specific control channel resource mapping explicitly. We intend for this capability to be available in later versions, which will enable more accurate simulation of the control overhead.

Sec. 2.7.2). Uncorrupted packets are then received by the `MmWavePhy` instance, which forwards them up to the MAC layer SAP.

## 2.7 MAC Layer

TDMA is widely assumed to be the de-facto scheme for mmWave access because of the dependence on analog beamforming, where the transmitter and receiver align their antenna arrays to maximize the gain in a specific direction (rather than with a wide angular spread or omni-directionally, as in conventional systems). Many early designs and prototypes have been TDMA-based [32, 61, 108], with others incorporating SDMA for the control channel only [109]. SDMA or FDMA schemes (as in LTE) are possible with *digital beamforming*, which would allow the base station to transmit or receive in multiple directions at the same time.

Furthermore, one of the foremost considerations driving innovation for the 5G MAC layer is latency. Specifically, the Key Performance Indicator (KPI) of 1 ms over-the-air latency has been proposed as one of the core 5G requirements by such standards bodies as the International Telecommunication Union [2], as well as by recent pre-standardization studies such as those carried out under the METIS 2020 project [115]. However, a well-known drawback of TDMA is that fixed slot lengths or TTIs can result in poor resource utilization and latency, which can become particularly severe in scenarios where many intermittent, small packets must be transmitted to/received from many devices [110].

Based on these considerations, variable TTI-based TDMA frame structures and MAC schemes have been proposed in [20, 109–111, 116]. This approach allows for slot sizes that can vary according to the length of the packet or TB to be transmitted and are well-suited for diverse traffic since they allow bursty or intermittent traffic with small packets as well as high-throughput data like streaming and file transfers to be scheduled efficiently.

The MAC layer implementation can be found in the `MmWaveEnbMac` and `MmWaveUeMac` classes, whose main role is the coordination of procedures such as scheduling and retransmission. Moreover, they interact with the RLC layer to receive periodic reports on the buffer occupancy, i.e., the Buffer Status Reports (BSRs), and with the physical layer classes for the transmission and reception of packets. To carry out their functionalities, the MAC classes interact with several other classes, that we will describe in the following paragraphs.

### 2.7.1 Adaptive Modulation and Coding

The role of the Adaptive Modulation and Coding (AMC) mechanism is to adapt the modulation scheme and the coding applied on top to the channel quality, measured using CQIs. In the simulator, this translates into (i) mapping the CQI into the Modulation and Coding Scheme (MCS), using the error model implemented in the `MmWaveMiErrorModel` and described in Sec. 2.5.4, and (ii) computing the available TB size for a subframe given the MCS. This information is then used by the scheduler to perform radio resource management.

The AMC is implemented in the `MmWaveAmc` class, which uses most of the code of the corresponding LENA module class. Some minor modifications and additional methods were necessary to accommodate the dynamic TDMA MAC scheme and frame structure. For instance, the `GetTbSizeFromMcsSymbols()` and `GetNumSymbolsFromTbsMcs()` methods are used by the scheduler to compute the TB size from the number of symbols for a given MCS value, and vice versa. Also the `CreateCqiFeedbackWbTdma()` method is added to generate wideband CQI reports for variable-TTI slots.

Fig. 2.6 shows the results of the test case provided in `mmwave-amc-test.cc`. This simulation serves to demonstrate the performance of the AMC and CQI feedback mechanisms

**Figure 2.6:** Rate and MCS vs. SINR for a single user under AGWN and fast-fading mmWave channels.

for a single user in the uplink (although a multi-user scenario could easily be configured as well). The default PHY/MAC parameters in Table 2.1 are used along with the default scheduler and default parameters for the statistical path loss, fading and beamforming models (i.e., `MmWavePropagationLossModel` and `MmWaveBeamforming`).

We compute the rate versus the average SINR over a period of 12 seconds (long enough for the small-scale fading to average out). The average PHY-layer rate is computed as the average sum of the sizes of successfully-decoded TBs per second. Every 12 seconds we artificially increase the path loss while keeping the UE position fixed. As the SINR decreases, the MAC will select a lower MCS level to encode the data. The test is performed for the Additive White Gaussian Noise (AGWN) case (i.e., no fading) as well as for small-scale fading. Although the UE position relative to the base station is constant, we can generate time-varying multi-path fading through the `MmWaveBeamforming` class by setting a fixed speed of 1.5 m/s to artificially generate Doppler, which is a standard technique for such an analysis. Also, we assume that the long-term channel parameters do not change for the duration of the simulation.

Fig. 2.6 therefore shows the data rate that it is possible to achieve with a certain SINR and with a certain modulation and coding scheme. If this plot is compared to the one generated from a similar test in Figure 3.1 of the LENA documentation [90], we notice that the AGWN curve from the mmWave test is shifted by approximately 5 dB to the left, indicating that the LENA version is transitioning to a lower MCS at a much higher SINR. This is because the LENA test is using the more conservative average SINR-based CQI mapping, which targets a much smaller TB error probability. In our test, we use the Mutual Information-Based Effective SINR scheme described in Sec. 2.5.4 with a target maximum TB error of 10% in order to maximize the rate for a given SINR [105].

### 2.7.2 Hybrid ARQ Retransmission

Full support for HARQ with soft combining is now included in the mmWave module. HARQ is a technique introduced in [117] and extensively used in LTE networks [20], which enables fast retransmissions with incremental redundancy in order to increase the probability of successful decoding and the efficiency of the transmissions. In LTE, the HARQ mechanism is based on multiple stop and wait retransmission processes, and a maximum of 8 simultaneous HARQ processes can be active at any given time [118]. The HARQ retransmissions have priority with respect to new transmissions, thus the available resources are given first to HARQ processes and

then to the data queued in the RLC buffers. Despite being fundamental in protecting from the losses of packets due to rapid variations in the channel quality, the HARQ mechanism introduces additional latency [20, 404], therefore the optimization of its performance is necessary to enable the target of sub-1-ms latency for ultra-low-latency communications.

The `MmWaveHarqPhy` class along with the functionalities within the different scheduler classes are based heavily on the LENA module code. The scheduler at the eNB uses the information provided by HARQ feedback messages to assign new resources to the HARQ processes that require retransmissions. Each transport block is granted a maximum number of transmission attempts, which is set to 3, as in LTE. However, some novelties are introduced in `MmWaveHarqPhy` in order to account for the more challenging channel conditions of the mmWave scenario. First, multiple HARQ processes per user can be created not only for the downlink but also for the uplink. Second, the number of processes per user is not fixed to 8, but can be configured through the `NumHarqProcesses` attribute in `MmWavePhyMacCommon`. This makes it possible to control (and, if needed, increase) the number of the simultaneous stop and wait retransmission processes and optimize the bandwidth utilization. Third, additional modifications were needed to support larger codeword sizes in both the `MmWaveHarqPhy` and `MmWaveMiErrorModel` classes. Finally, the integration with the flexible TTI physical layer allows a reduction in the latency of the retransmissions, as discussed in [20].

### 2.7.3 Schedulers

We now present the implementations of four scheduler classes for the variable TTI scheme. These differ significantly from the OFDMA-based schedulers available in ns–3 LENA [90] as, instead of allocating Resource Blocks/Resource Block Groups of frequency-domain resources, these TDMA-based schedulers allocate time-domain symbols within a periodic subframe to different users in the DL or UL direction.

Before scheduling new data, Buffer Status Report and CQI messages are first processed. The MCS is computed by the AMC model for each user based on the CQIs for the DL or SINR measurements for the UL data channel. The MCS and the buffer length of each user are used to compute the minimum number of symbols required to schedule the data in the user's RLC buffers. This procedure for estimating the optimal MCS and determining the minimum number of symbols is common to each of the schedulers described in the following.

**Round Robin (RR) Scheduler:** The `MmWaveFlexTtiMacScheduler` class is the default scheduler for the mmWave module. It supports the variable TTI scheme previously described in Sec. 2.6 and assigns OFDM symbols to user flows in *Round-Robin* order. Upon being triggered by a subframe indication, any HARQ retransmissions are automatically scheduled using the available OFDM symbols. While the slot allocated for a retransmission does not need to start at the same symbol index as the previous transmission of the same TB, it does need the same number of contiguous symbols and MCS, since an adaptive HARQ scheme (where the retransmission can be scheduled with a different MCS) has not yet been implemented.

To assign symbols to users, the total number of users with active flows is first calculated. Then the total available data symbols in the subframe are divided evenly among users. If a user requires fewer symbols to transmit its entire buffer, the remaining symbols (i.e., the difference between the available and required slot length) are distributed among the other active users.

One also has the option to set a fixed number of symbols per slot by enabling the *fixed TTI* mode. Although the same general subframe structure is maintained, slots will then be allocated in some multiple of `SymPerSlot` symbols. Setting the `SymPerSlot` attribute of the scheduler class to the number of slots per subframe, for instance, will result in only one UE being scheduled per subframe, which would be highly inefficient in a multi-user cell.

**Proportional Fair (PF) Scheduler:** *Proportional Fair* is another well-known discipline, and is provided by the `MmWaveFlexTtiPfMacScheduler` class. The PF scheduler attempts to prioritize traffic for high-SINR users while maintaining some measure of fairness by ensuring that low-SINR, cell-edge users are also scheduled [119].

**Earliest Deadline First (EDF) Scheduler:** The `MmWaveFlexTtiMaxWeightMacScheduler` class implements an *Earliest Deadline First* policy, which is a priority queue-based policy that weighs flows by their relative deadlines for packet delivery. The deadlines are initially set according to the delay budget of the QoS Class Indicator (QCI) configured by the RRC layer [120,121]. The deadline of the Head-of-Line (HOL) packets of each RLC buffer is then compared, and that with the earliest deadline is scheduled first. Any remaining symbols in the subframe are allocated to the packet with the next smallest relative deadline and so forth until all $N_{sym}$ symbols are assigned. The EDF scheduler is the only deliberately delay-sensitive scheme included in the mmWave module and can be useful for evaluating the latency performance of mmWave links, as in the simulations presented in Sec. 2.10.

**Maximum Rate (MR) Scheduler:** The *Maximum Rate* policy realized in the `MmWaveFlexTtiMrMacScheduler` class schedules only the users with the highest SINRs to maximize cell throughput. Initially, UEs are sorted based on their optimal MCS values. Symbols are distributed in round-robin fashion among UEs at the highest MCS level until the minimum number of symbols required to transmit the entire buffers of these users has been assigned. This is then repeated for UEs at the second highest level, and so forth until all symbols of the subframe are allocated.

The MR scheduler potentially suffers from extremely poor fairness when there are both high- and low-rate users, and some users may not be scheduled at all, thus making it impractical for any real-world multi-user system. However, it may still be useful for testing system capacity and performance.

### 2.7.4 Carrier Aggregation

The modeling of the Carrier Aggregation (CA) feature in the mmWave module for ns-3 has been introduced in [414]. It follows the 3GPP specifications for New Radio (NR) [7], and aligns the PHY and MAC design of the mmWave module to the ns-3 LTE module implementation [69], for which the CA capability was introduced in [122].

As shown in Fig. 2.1, the implementation for the data plane involves the lower layers of the protocol stack (i.e., MAC and PHY), i.e., it is transparent with respect to the functionalities offered by the RLC and Packet Data Convergence Protocol (PDCP) layers. The control functionalities are performed by the RRC layer, which is in charge of sharing the information for the carrier setup between the base station and the UE. In particular, the base station broadcasts information on the primary Carrier Component (CC), and the UE connects to it. Then, when it enters the RRC_CONNECTED state, the base station RRC can instruct the UE to add and/or remove additional carriers with different parameters [7].

In our CA model, and as generally done in the ns-3 mmWave module, we inherit and extend the inter-layer interfaces of the LTE module (i.e., the SAPs) [69] and the classes that implement them, in order to increase the flexibility and account for different channel and propagation conditions for the different carriers, as well as possibly different numerologies, as specified in [7].

**CC Configuration**   Similarly to the LTE implementation [122], the basic class of the CA implementation is the `MmWaveComponentCarrier` class and its `MmWaveEnbComponentCarrier` and `MmWaveUeComponentCarrier` extensions. An instance of this class represents a single carrier, and contains pointers to the associated protocol stack layers and relevant configurations, as

**Figure 2.7:** Information represented by instances of `MmWaveComponentCarrier` (and extensions)

shown in Fig. 2.7. In particular, in our implementation, a `MmWaveComponentCarrier` object contains a reference to a `MmWavePhyMacCommon` object, which is used to specify the numerology, frequency and bandwidth information for the carrier. The `MmWavePhyMacCommon` class was introduced in [66], and prior to the CA implementation, a single `MmWavePhyMacCommon` was created by `MmWaveHelper` during the configuration of the simulation. This object was shared by all the eNB and UE PHY and MAC layer classes, as well as by the channel model classes, to provide access to a set of common parameters. With the CA implementation, instead, an instance of `MmWavePhyMacCommon` is created for each possible carrier, and is associated to the unambiguous identifier of the carrier (i.e., carrier ID stored in the `m_componentCarrierId` private variable of `MmWavePhyMacCommon`). Each of these objects is shared by all classes of the layers at the base station and UE side that are related to the same carrier. The carrier-specific `MmWavePhyMacCommon` instance then defines the carrier frequency (with the attribute `CenterFreq`), the bandwidth (for which it is possible to control the size of the resource blocks and their number) and the frame structure (i.e., the number of symbols per subframe, their duration, and the number of subframes per frame).

**CC Managers** The different `MmWaveComponentCarrier` objects in the UEs and base stations are managed by a single CC manager, i.e., an object that implements respectively the `LteUeComponentCarrierManager` or the `LteEnbComponentCarrierManager` interfaces. The CC manager, together with the `MmWaveUeMac` or `MmWaveEnbMac` classes, models the functionalities of the MAC layer for the mmWave protocol stack. In particular, at the base station side, it receives the BSRs from the RLC layers, and forwards them to the `MmWaveScheduler` instances following different policies according to the particular implementation of the CC manager. The schedulers then allocate the available resources and generate Downlink Control Information for the different carriers. In the current implementation, the scheduling on the different carriers is independent, but we plan to extend it in order to model joint cross-carrier scheduling. The CC manager at the UE side is a simplified version of that at the base station, since it does not need to split the BSRs between the carriers, but limits itself to forwarding them to the base station CC manager. In particular, only the primary CC is used to report the BSRs and the exchange of control information, since it is the only CC in which the Service Radio Bearers are set up.

In the mmWave CA implementation, we provide different implementations of the CC manager at the base station side. As for LTE, there is a `MmWaveNoOpMacComponentCarrierManager` which is used for single-carrier simulations, and a `MmWaveRrMacComponentCarrierManager`, which applies a round robin policy and splits the BSRs equally across the carriers, with the result that they reach a similar throughput. In addition, we include also a bandwidth-aware CC manager. It is likely that different carriers over different frequency bands will use different bandwidths, given that the higher the carrier frequency the larger the bandwidth that can be allocated to mobile network operators[9]. Therefore, a typical use case for CA in the mmWave band would be

---

[9]For example, the International Telecommunication Union is considering the allocation to mobile operators of

32

the aggregation of a CC at relatively low carrier frequency, with a smaller bandwidth, but with better propagation properties (i.e., lower pathloss), and other CCs at much higher frequencies with larger bandwidths. In this case, a round robin CC manager that evenly splits the packets to be transmitted across the different carriers would not yield an optimal performance, given the different data rates that can be supported by the CCs. Therefore, the CC manager implemented in the `MmWaveBaRrMacComponentCarrierManager` class is made aware of the bandwidth available to the different carriers during the simulation setup, and then, when it receives the BSRs from the RLC layer instances in the base station or the UE, it divides the reports according to the bandwidth ratio across the carriers.

Another difference with respect to the LTE implementation is the usage of different channel model objects for the different carriers. In the mmWave module, indeed, the joint modeling of the propagation loss, the small and large scale fading and the beamforming has a fundamental importance for the accuracy of the simulation results. In our previous paper [406] we introduced the implementation of the 3GPP channel model for frequencies above 6 GHz [71], which has features that depend on some carrier-specific parameters, such as the bandwidth and the carrier frequency. Therefore, we decided to use different channel objects for each carrier, and use the `MmWavePhyMacCommon` of the carrier to set up the necessary parameters. Finally, we extended the `MmWaveSpectrumValueHelper` class in order to support the configuration (i.e., bandwidth, numerology and carrier frequency) of the different carriers.

CA configuration in ns-3 mmWave simulations  Thanks to the adoption of a `MmWavePhyMacCommon` object per carrier, the user of the ns-3 mmWave module has a lot of flexibility in configuring the parameters of the simulation. We provide two comprehensive simulation examples in the `mmwave-ca-same-bandwidth.cc` and `mmwave-ca-diff-bandwidth.cc` files in the `examples` folder of the mmWave module. The first step in the simulation configuration is the initialization of a `MmWavePhyMacCommon` per CC. The method `SetAttribute` can be used to set the relevant parameters for the carrier. Then, a map that associates the carrier ID to the `MmWaveComponentCarrier` is created, and passed as a parameter to the `MmWaveHelper` with the method `SetCcPhyParams`. The user then deploys the nodes, installs the relevant `NetDevices`, mobility models and applications as in a non-CA simulation script. It is the `MmWaveHelper`, indeed, that transparently takes care of the initialization of the channel objects and the association to the correct carrier, and of the setup of the mmWave base stations and UEs with the carriers information.

## 2.8   RLC Layer

The RLC layer is inherited directly from the LTE module described in [69], and therefore all the LTE RLC entities are included. Moreover, the RLC AM retransmission entity is modified to be compatible with the mmWave PHY and MAC layers, and Active Queue Management (AQM) for the RLC buffers is introduced as a new optional feature.

### 2.8.1   Modified RLC AM Retransmission

Reordering and retransmission play an important role in RLC AM. Due to the shortened mmWave frame structure, the timers of the RLC entity should also be reduced accordingly, e.g., the `PollRetransmitTimer` is changed to 2 ms from 20 ms. Moreover, the original LTE module does not perform re-segmentation for retransmissions, and the RLC segment waits in

---

bands of approximatively 3 GHz in the 20–30 GHz spectrum, and of 10 GHz in the 60–80 GHz spectrum [123].

the retransmission buffer until the transmission opportunity advertised by the lower layers is big enough. This becomes problematic when the transmission is operated over an intermittent channel, as a sudden channel capacity drop would halt the retransmission entirely. Therefore, we added to the RLC AM layer implementation the capability of performing segmentation also for the retransmission process, in order to support an intermittent mmWave channel. The re-segmentation process deployed in our RLC AM class works as follows: If the number of bytes that can be transmitted in the next opportunity is smaller than the bytes of the segment that should be retransmitted, then the segment will be split into smaller subsegments with a re-segment flag set to be true. The RLC layer at the receiver side will check the flags of the subsegments, and wait until the final one if the flag is set to be true. Finally, the subsegments are assembled to construct the original segment and forwarded to the upper PDCP layer if all subsegments are received correctly. Otherwise, all subsegments are discarded and another retransmission is triggered.

### 2.8.2 Active Queue Management

Active Queue Management techniques are used in the buffers of routers, middleboxes and base stations in order to improve the performance of TCP and avoid the manual tuning of the buffer size. Different strategies have been defined in the literature [124, 125], and several of them are implemented in ns–3 [126–128]. AQM strategies allow the network to avoid congestion at the buffers, because they react early to the increase in the buffer occupancy by dropping some packets before the buffer is full. With respect to the default Drop-tail approach, in which no packet is dropped until the buffer is full, AQM techniques make TCP aware of possible congestion earlier, avoiding the latency increase which is typical of the bufferbloat phenomenon [129].

Some early AQM, such as Random Early Detection (RED) [130, 131], were widely studied in the literature, but failed to find market traction because of the intrinsic complexity of their tuning parameters. Recently, a simpler AQM technique, namely Controlled Delay (CoDel) [132], was proposed to replace RED queues. CoDel adapts to dynamic link rates without parameter configuration, and is able to discriminate "good" and "bad" queues: good queues can quickly empty the buffer, whereas "bad" queues persistently buffer packets. CoDel works by monitoring the minimum queue delay in every 100 ms interval, and only drops packets when the minimum queue delay is more than 5 ms.

In the RLC layer of the LTE module, the default queue management is Drop-tail. In the mmWave module, the RLC layer can use either the default Drop-tail approach or more sophisticated AQM techniques, that can be enabled by setting the `EnableAQM` attribute to true. The default AQM is the CoDel scheme, however it is possible to use any of the queues available in ns–3 by modifying the queue attribute in the `LteRlcAm` class. The evaluation of the AQM scheme will be further discussed in Chapter 6.

## 2.9  Dual Connectivity Extension

The ns–3 mmWave module is also capable of performing simulations with dual-stack UEs connected both to an LTE eNB and to a mmWave gNB. This feature, partially described in [402], was introduced because mmWave 5G networks will likely use multi-connectivity and inter-networking with legacy RATs in order to increase the robustness with respect to mobility and channel dynamics [5, 62, 63, 133, 134, 388]. The source code can be found in the `new-handover` branch of the ns–3 mmWave module repository.

The DC implementation of this simulation module assumes that the core networks of LTE and of mmWave will be integrated, as in one of the options described in [135]. Therefore the LTE and

**Figure 2.8:** Simplified UML diagram of a dual-connected device, an LTE eNB and a MmWave eNB that also support carrier aggregation. We only report the main classes of the DC-CA integration implementation, i.e., the SAP interfaces are omitted.

the mmWave gNBs share the same backhaul network, i.e., they are connected to each other with X2 links and to the MME/PGW nodes with the S1 interface. As to the RAN, the DC solution of this module is an extension of 3GPP's LTE DC proposal [136]. In particular, a single bearer per DC flow is established, with a connection from the core network to the LTE eNB, where the flow is split and forwarded either to the local stack or to the remote mmWave stack. We chose the PDCP layer as the integration layer, since it allows a non-colocated deployment of the base stations and a clean-slate approach in the design of the PHY, MAC and RLC layers [137].

A basic diagram for a DC UE device, an LTE eNB and a mmWave gNB is shown in Fig. 2.8. The core of the DC implementation is the `McUeNetDevice` class, which is a subclass of the ns–3 `NetDevice` and provides an interface between the ns–3 TCP/IP stack and the custom lower layers. The `McUeNetDevice` holds pointers to the custom lower layer stack classes, and has a `Send` method that forwards packets to the TCP/IP stack. This method is linked to a callback on the `DoRecvData` of the `EpcUeNas` class, which as specified by the 3GPP standard acts as a connection between the LTE-like protocol stack and the TCP/IP stack.

The `McUeNetDevice` describes a dual connected UE with a single `EpcUeNas`, but with a dual stack for the lower layers, i.e., there are separate LTE and mmWave PHY and MAC layers. Moreover, there is an instance of the RRC layer for both links. This grants a larger flexibility,

because the functionalities and the implementation of the two layers may differ. Besides, the LTE RRC manages both the LTE connection and the control plane features related to DC, while the mmWave RRC handles only the mmWave link. The usage of a secondary RRC, dedicated to the mmWave link, avoids latency in control commands (i.e., the mmWave gNB does not have to encode and transmit the control Packet Data Unit (PDU)s to the master LTE eNB). The `EpcUeNas` layer has an interface to both RRC entities to exchange information between them.

The LTE RRC manages also the data plane for the DC devices. In particular, for each bearer, a dual connected PDCP layer is initialized and stored in the LTE RRC. The classes describing the DC PDCP layer are `McEnbPdcp` and `McUePdcp`, respectively at the eNB side and at the UE side. They both extend the `LtePdcp` class with a second interface to the RLC. However, while `McUePdcp` simply has to communicate with a local RLC in the UE, the implementation of `McEnbPdcp` requires new interfaces to the class describing the X2 links between base stations (i.e., `EpcX2`). In particular, in downlink the eNB PDCP has to send packets to the X2 link and the mmWave RLC layer has to receive them, and vice versa in uplink.

The DC module can be used to simulate different dual connected modes, i.e., it can support both Fast Switching (FS) and throughput-oriented dual connectivity, according to which RRC and X2 procedures and primitives are implemented. With FS, the UE is in the RRC_CONNECTED state with respect to both eNBs, but only transmits data to one of the two. With the other option, the UE can transmit data simultaneously on both RATs, and different flow control algorithms can be plugged in and tested.

As to the physical layer, the two stacks rely on the mmWave and LTE channel models. Notice that since the two systems operate at different frequencies, modeling the interference between the two RATs is not needed. Each of the two channel models can therefore be configured independently.

In order to use an `McUeNetDevice` as a mobile User Equipment in the simulation, the helper class of the mmWave module was extended with several features, such as (i) the initialization of the objects related to the LTE channel; (ii) the installation and configuration of the LTE eNBs, so that they can be connected to the LTE stack of the `McUeNetDevice`; and (iii) the methods to set up a `McUeNetDevice` and link its layers as shown in Fig 2.8. An example on how to set up a dual-connectivity based simulation is provided in the file `mc-twoenbs.cc`.

**RRC Layer for Dual Connectivity and Mobility.** The RRC layer implementation of the original LTE ns–3 module was extended in order to account for DC-related control procedures. In particular, the multi-connectivity uplink-based measurement framework described in [62] was added with changes to the `MmWaveEnbPhy`, `EpcX2` and `LteEnbRrc` classes. The `MmWaveEnbPhy` instance simulates the reception of uplink reference signals (which are accounted for as overhead in the data bearers resource allocation), computes the SINR for each UE in the scenario[10], and sends this information to the LTE eNB on the X2 link. This also allows the simulation of a delay in the reporting, since the control packets with the SINR values must be transmitted on an ns–3 `PointToPointLink`, which adds a certain latency and has a certain bitrate.

Thanks to this framework, the LTE eNB is able to act as a coordinator for the surrounding mmWave gNBs, and learns which is the best association (in terms of SINR) between UEs and mmWave secondary gNBs. This enables automatic cell selection for mmWave gNBs at the beginning of a simulation, and the control of mobility-related operations. The DC module is indeed capable of simulating FS procedures between mmWave and LTE links and Secondary Cell Handover (SCH) (i.e., handovers between mmWave gNBs that do not involve the MME in the core network) initiated by the central controller in the LTE eNB. It is also possible to use the DC module to simulate X2-based RAT handovers between the LTE and mmWave gNBs, i.e., to

---

[10]The framework assumes that the optimal beam is always chosen, so the actual directional scan procedure described in [62] is not simulated

use standalone UEs based on `McUeNetDevice` that can perform handovers from the LTE to the mmWave gNBs, and vice versa.

Different handover (either inter-RAT or SCH) algorithms can be tested, by implementing them in the `LteEnbRrc` class. In order to make the handover simulation more compliant with the 3GPP specifications, the lossless handover option implemented for ns–3 in [138] was adapted to the DC module in order to forward the RLC buffer content to the target RAT/eNB RLC layer for both the SCH and the FS. Moreover, in order to model the additional latency given by the interaction with the MME for inter-RAT handovers for standalone UEs, the link between the eNBs and the MME is modeled in this module as a `PointToPointLink`, while in the original ns–3 LTE module it is an ideal connection.

## 2.10   Use Cases

In this section, we illustrate various examples of scenarios[11] that can be simulated to show the utility of the module for the analysis of novel mmWave protocols and for testing higher-layer network protocols over 5G mmWave networks. After each particular use case example, we also provide to the interested reader some references to recent papers that report additional results obtained using the ns–3 mmWave simulation module.

### 2.10.1   Simulation Setup Walk-through

In order to proficiently use the mmWave ns–3 module, a basic knowledge of ns–3 is required. We therefore advise the interested users to first study the extensive documentation referenced in Sec. 2.3. Moreover, we provide some basic ns–3 scripts in the `examples` folder of the mmWave module, that can be a basis for the design of any simulation script that uses the mmWave module. In the following paragraphs, we will describe the basic structure of a typical example in simple steps.

The first step is to configure all the attributes needed in a simulation. A complete list of attributes related to the mmWave module can be found in the `mmWaveAttributesList` file in the module repository. The second step involves the setup of the `MmWaveHelper` object, which provides methods to create the entities involved in the simulation (e.g., the channel-related objects and the `MmWavePhyMacCommon` object), to install the mmWave stack over ns–3 nodes (for both UEs and eNBs), to perform the initial attachment of a UE to the closest eNB and to enable or disable the generation of simulation traces. Moreover, if the scenario of interest is an end-to-end scenario, the core network and the internet must be set up as well. The first is created by the `MmWavePointToPointEpcHelper`, which also provides a pointer to the Packet Gateway (PGW) node. This is then usually connected to a remote host, and the internet stack (i.e., the TCP/IP protocol suite) is added to the UEs and to the remote host.

In the third step, the positions and velocities of the eNBs and UEs are specified using one or more `MobilityHelper` objects and different mobility models. Moreover, buildings and obstacles can be added to the scenario using the ns–3 buildings module and the `Buildings` and `BuildingsHelper` objects. The fourth step requires the setup of applications in the UEs and in the remote host (if an end-to-end scenario is considered), in order to simulate downlink and uplink traffic. ns–3 provides a wide range of different applications, and helpers that take care of their setup. They can run on either UDP or TCP sockets, and several TCP congestion control

---

[11]The simulations in this section are all configured with the basic PHY and MAC parameters in Table 2.1, with other notable parameters given in the sequel.

(a) Empirical CDF of rate

(b) Empirical CDF of latency

**Figure 2.9:** Distributions of PHY-layer throughput and IP-layer latency for 70 UEs, 10 Mbps/UE arrival rate



(a) Empirical CDF of rate

(b) Empirical CDF of latency

**Figure 2.10:** Distributions of PHY-layer throughput and IP-layer latency for 7 UEs, 100 Mbps/UE arrival rate

versions are available. Finally, the simulation can be run using the `Simulator` object of ns–3, and traces are generated.

### 2.10.2 Multi-User Scheduling Simulation

In this experiment, the throughput and latency of users of a mmWave cell with 1 GHz of bandwidth are simulated for variable TTI and each of the scheduling policies described in Sec. 2.7.[12] We shall see how the choice of the scheduler has a significant impact on the subframe utilization and latency of the multi-user cell. In these scenarios, UEs have similar distances from the eNB but are assigned the constant speed of 25 m/s (typical of vehicular users), which results in a lower achievable rate, on average, as well as increased packet errors compared to walking users due to the more rapid variation in the channel.

The simulation is again run over 10 drops for each of two scenarios and using default parameters from Table 2.1. In the first scenario, 70 UEs are simulated with each UE generating IP-layer traffic at an average arrival rate of 10 Mbps. In the second scenario, 7 UEs are simulated with a 100 Mbps arrival rate per UE.

These specific combinations of users and rates are deliberately chosen because they illustrate the cut-off point at which the system is no longer able to service most users at the requested rate, leading to backlogged queues and increased latency. That is, we wish to analyze the performance

---

[12]The multi-user scheduling experiment can be reproduced by running the `mmwave-epc-tdma.cc` example simulation.

| | | PHY-layer throughput [Mbps] | | | |
|---|---|---|---|---|---|
| | Policy | Cell | Mean UE | Mean 5% Worst UE | Max UE |
| 70 UE – 10 Mbps | RR | 1815.92 | 25.94 | 1.11 | 48.80 |
| | PF | 1494.61 | 21.35 | 3.26 | 43.94 |
| | MR | 2273.18 | 32.47 | 0.00 | 151.36 |
| | EDF | 925.80 | 13.23 | 7.31 | 31.02 |
| 7 UE – 100 Mbps | RR | 715.26 | 102.18 | 49.18 | 134.28 |
| | PF | 758.32 | 108.33 | 52.32 | 158.16 |
| | MR | 766.26 | 109.47 | 47.26 | 158.22 |
| | EDF | 647.98 | 92.57 | 63.89 | 121.76 |

**Table 2.2:** DL PHY throughput for RR, PF, MR and EDF scheduling policies.

| | | IP-layer latency [ms] | | |
|---|---|---|---|---|
| | Policy | Mean UE | Mean 5% Worst UE | Max UE |
| 70 UE – 10 Mbps | RR | 7.47 | 69.35 | 118.62 |
| | PF | 2.83 | 34.48 | 106.54 |
| | MR | 0.65 | 1.89 | 3.07 |
| | EDF | 1.63 | 7.91 | 30.65 |
| 7 UE – 100 Mbps | RR | 0.67 | 2.01 | 2.37 |
| | PF | 0.55 | 0.68 | 0.77 |
| | MR | 0.56 | 0.78 | 1.09 |
| | EDF | 0.69 | 1.41 | 1.44 |

**Table 2.3:** IP-layer latency for RR, PF, MR and EDF scheduling policies.

at the knee in the curve of the delay taken as a function of the system utilization. In the variable TTI system, this bottleneck effect has the following potential causes: (i) the number of users that must be serviced exceeds the number of available slots (ultimately limited by the number of time-domain symbols), independently of the total throughput requested by the users; (ii) the number of users that are connected to eNB is smaller than the number of available slots, but the total throughput they request exceeds the available resources in the given time period; or (iii) a combination of the previous cases.

These effects are demonstrated in Figs. 2.9 and 2.10 for the 70 UE/10 Mbps and 7 UE/100 Mbps arrival rate scenarios, respectively. The mean, maximum and cell-edge (i.e., 5% worst-case) user PHY rates and IP-to-IP layer latencies are also provided in Tables 2.2 and 2.3 along with the utilization and Jain's Fairness Index in Table 2.4.

For the 70 UE case, Fig. 2.9a shows the distribution of the mean rate experienced by each UE over the simulation duration. It can be seen that the MR and RR policies exhibit the greatest disparity between users scheduled with high and low rates.

It can also be observed that the PHY rate significantly exceeds the 10 Mbps arrival rate for some users, which leads to the poor utilization for these two policies, as shown in Table 2.4. The reason why the utilization (defined as the ratio of the received IP-layer rate to the allocated PHY-layer rate for each terminal) suffers in these cases is that the UEs with higher achievable rates are heavily favored by the MR and RR schedulers. As these users are typically scheduled at a higher MCS level, even a single 4.16 μs-long time-domain symbol has the capacity to transmit kilobytes of data, which cannot be fully taken advantage of given the low 10 Mbps arrival rate. Insufficient data is buffered at the MAC layer to utilize the full slot and useless padding bits must be added. This effect is felt less by users under the PF and EDF policies, which are inherently

|          | Policy | Fairness | Utilization |
|----------|--------|----------|-------------|
| 70 UE – 10 Mbps | RR | 0.71 | 0.53 |
|          | PF | 0.76 | 0.73 |
|          | MR | 0.28 | 0.39 |
|          | EDF | 0.96 | 0.87 |
| 7 UE – 100 Mbps | RR | 0.95 | 0.74 |
|          | PF | 0.90 | 0.77 |
|          | MR | 0.91 | 0.76 |
|          | EDF | 0.99 | 0.84 |

**Table 2.4:** Fairness index and utilization (received IP-layer rate/allocated PHY rate) for RR, PF, MR and EDF scheduling policies.

more fair and allow more resources to be scheduled for lower-MCS users.

The ensuing effect of these trends on latency is shown in Fig. 2.9b. Here latency is measured as the time between the arrival time of packets at the PDCP layer of the eNB stack and the time they are delivered to the IP layer at the UE. Naturally, the MR scheduler offers the best delay performance because only 40% of the users, i.e., those with the highest rates, are ever scheduled (unscheduled users with zero rate are not included in the figure). The RR policy offers the highest worst-case delays but is able to achieve mean latencies of less than 1 ms for more than 60% of the users. Of all the policies besides MR, Earliest-Deadline First offers the best worst-case delays, as it attempts to balance the delay of all users by scheduling them exclusively based on their relative deadlines (not taking into account achievable rates). The EDF scheduler is able to achieve a mean UE latency of 1.6 ms, which, as we will see from the experiment in the next section, drops below 1 ms for 60 or fewer users (with the same arrival rate).

Finally, it can be observed from Fig. 2.10b that, despite having the same total packet arrival rate of 700 Mbps as in the 70 UE case, latencies are much lower overall in the 7 UE/100 Mbps per UE case. This can be clearly explained by the higher utilization factor in Table 2.4. In this scenario, the number of available slots for scheduling different users is no longer the bottleneck. Though we still see that a significant number of users are scheduled at rates that exceed their 100 Mbps arrival rates, the utilization is notably better than in the 10 Mbps case for all scheduling policies. Thus, the channel capacity itself is better utilized, allowing most users to be scheduled at the requested rate, thereby avoiding additional queue wait time and delay.

Additional results on the impact of the scheduling on throughput and latency can be found in [65].

### 2.10.3   Latency Evaluation for Variable and Fixed TTI Schemes

The results and the discussion introduced in this section are taken from our previous article [20], where the interested reader will find a more comprehensive treatment of techniques to achieve low latency in mmWave 5G cellular systems. While the qualitative benefits of variable TTI over fixed TTI may seem self-evident, in this section we quantify the performance gains for a multi-user TDMA mmWave system with 1 GHz of bandwidth. We also demonstrate that, with the low-latency scheduling loop enabled by the proposed frame structure, LTE-style Hybrid ARQ can still be employed for enhanced link-layer reliability without excessively exceeding the delay constraints. We model the subframe formats shown in Fig. 2.5 for two subframe periods: the default 100 µs subframe, equivalent to 24 OFDM symbols, and a 66.67 µs subframe, equivalent to 16 OFDM symbols. The symbol length of 4.16 $\mu s$ is based on the numerology in [32]. Each subframe has one fixed DL-CTRL and one UL-CTRL symbol, with the remaining symbols used for DL or UL data slots. For fixed TTI mode, the entire subframe is allocated to a single user,

**(a)** 10 Mbps per UE arrival rate (100-bytes packets)  **(b)** 100 Mbps per UE arrival rate (1200-bytes packets)

**Figure 2.11:** Latency and Deadline Miss Ratio as a function of the downlink IP-layer arrival rate for fixed and variable TTI radio frame structures.

whereas for variable TTI mode, the scheduler may allocate any number of data symbols within the subframe to match the throughput required by each user.

We also note that UEs are again modeled as moving at 25 m/s, typical of vehicular speeds, which causes fast channel variations and frequent packet errors from small-scale fading (it is observed that between 0.5% and 3% of the transport blocks are lost and require retransmission).

We consider a simple traffic model with Poisson arrivals, where each UE sends an average of 12.5K packets per second (100-byte packets resulting in an average rate of 10 Mbps), as well as a separate, higher-throughput case where each UE sends an average of 83K packets per second (1200-byte packets resulting in an average rate of 100 Mbps). Scheduling is performed based on the EDF policy where the scheduler attempts to deliver each IP packet within 1 ms from its arrival at the PDCP layer and packets are assigned a priority based on how close they are to the deadline. Priority is therefore always given to HARQ retransmissions. We simulate the performance for between 10 and 100 UEs for a 10 Mbps (per UE) arrival rate and between 1 and 10 UEs for the 100 Mbps case, equivalent to a total IP-layer arrival rate of between 100 and 1000 Mbps in both cases.

In Fig. 2.11, the downlink radio link latency is averaged among the best 95% of the users (i.e., the 5% of UEs experiencing the highest latency are not considered). The Deadline Miss Ratio (DMR), which represents the fraction of packets delivered after the 1 ms deadline, is also

| Description | Value |
|---|---|
| Subframe length in µs | 100 or 66.67 |
| OFDM symbols per slot | 24 or 16 |
| HARQ processes (DL and UL) | 20 DL/20 UL |
| Number of UEs | Case 1: $\{10, 20, 30, 40, 50, 60, 70, 80, 100\}$ |
| Number of UEs | Case 2: $\{1, 2, 3, 4, 5, 6, 7, 8, 10\}$ |
| Traffic model | Case 1: Poisson, $\lambda = 12.5K$ pck/s, 100 B packets |
| Traffic model | Case 2: Poisson, $\lambda = 83K$ pck/s, 1200 B packets |

**Table 2.5:** Additional parameters for variable and fixed TTI latency experiment.

given for the top 95th percentile UEs. We see that, for a 10 Mbps arrival rate (Fig. 2.11a), variable TTI is able to achieve sub-ms average latency and a DMR of about 10% with over 60 users (corresponding to a 600 Mbps total packet arrival rate) and consistently outperforms fixed TTI. Fixed TTI, despite the relatively short subframe compared to LTE, exceeds 1 ms average latency and has a DMR of over 60% even for the 20 UE case and of more than 90% for 40 or more users. This result shows that variable TTI will be essential for reliable, low-latency service, particularly when considering use cases with many lower-rate devices, such as Machine-Type Devices (MTDs) [6].

For the higher-load (100 Mbps arrival rate per UE) case in Fig. 2.11b, we expect the deviation between the variable and fixed TTI schemes to be less pronounced, as the bottleneck is now the multi-user channel capacity and not the minimum slot size. However, we do find a reduction in radio link latency of around 500 µs, or 30%, for the variable TTI scheme in some cases.

We also find that, for a smaller number of users, the shorter 66.67 µs subframe offers some improvement over the longer 100 µs subframe thanks to the decreased turnaround time. In particular, the DMR is consistently less for the 100 Mbps/UE case for both variable and fixed TTI. However, this trend reverses with more users due to the lower ratio of control to data symbols in the 100 µs subframe case. We note that the control overhead could be somewhat mitigated by multiplexing data in the DL-CTRL region.

We also note that, in real-world implementations, there may be some additional delay related to beam tracking (i.e., for computing and applying the optimal TX/RX beamforming vectors), although the performance limitations of adaptive beamforming transceivers and channel tracking techniques in future implementations are still unknown. We assume that this delay can be neglected in our analysis because data is constantly being transmitted to each UE and channel state feedback is being transmitted by the UEs to the eNB in each subframe period (which is well within the coherence time observed in many studies), thus ensuring that the channel state information is always up-to-date at the eNB.

The performance of a mmWave cellular network with respect to the end-to-end or the RAN latency has been studied in several papers with the ns-3 mmWave module. In [20, 65] the authors discuss architectural and protocol solutions for low latency networks based on mmWave, while [57,404,411,413] propose and evaluate latency-reduction techniques for transport protocols and video streaming.

## 2.11   Integration with DCE and Examples

DCE was introduced in [89] as a powerful tool that combines the flexibility of a network simulator such as ns–3 with the robustness of the TCP/IP stack of the Linux kernel and the authenticity of real applications. There are several benefits in using this tool. First, the Linux kernel implements protocols which are not yet available for ns–3, or which are in an early development phase and present some limitations. An example is MPTCP, the multipath extension of TCP which makes it possible to transmit data on multiple subflows (i.e., a mobile user could simultaneously transmit on a Wi-Fi subflow and a cellular subflow) [139]. At the time of writing, it was implemented for ns–3 by different projects [140, 141], but none of them is completely compliant with the MPTCP specification, and they are not integrated in the main ns–3 release and validated. With DCE, instead, it is possible to use the MPTCP code developed and tested by the same MPTCP protocol designers [142]. Second, the Linux kernel TCP/IP stack is the most widely used in real production environments and datacenters, besides being the basis for the Android mobile operating system. Therefore, it is a very well tested codebase, with very few bugs. Moreover, its usage in network simulations provides a higher level of realism. Finally, with DCE it is possible

**Figure 2.12:** Random realization of the simulation scenario. The grey rectangle represents a building.



**(a)** BALIA



**(b)** Uncoupled Cubic

**Figure 2.13:** Throughput with MPTCP.

to use real POSIX socket-based applications. For example, the well known iPerf tool [143] can be used to measure the maximum achievable datarate in the network. It is also possible to simulate a website, with an http daemon in the server and wget as a client. Besides, standard ns–3 applications (`OnOffApplication`, `BulkSendApplication`) can be used with the Linux TCP/IP stack thanks to DCE Cradle.

In order to integrate DCE with the ns–3 mmWave module, it is necessary to patch the `KernelFdSocketFactory` class so that it recognizes the `MmWaveUeNetDevice`. The patch can be found in the `utils` folder of the ns–3 mmWave module repository. Then, replace the standard ns–3 folder with our mmWave module. Notice that, if MPTCP is used as the transport protocol, the DC extension must be used with the patch provided in the `utils` folder.

**MPTCP on mmWave links:** The latest Linux kernel implementation of MPTCP compatible with DCE can be found in the `net-next-nuse` library of the LibOS project [144]. The standard DCE distribution already provides MPTCP examples, which can be promptly extended in order to account for mmWave and LTE subflows, as long as they operate on links with different carrier frequencies (i.e., it is possible to simulate an MPTCP connection on a 2.1 GHz LTE link and a 28 GHz mmWave link). It is possible to simulate different state of the art congestion control algorithms for MPTCP, either coupled or uncoupled, as shown in [389, 404].

An example is in the file `dce-example-mptcp-mmwave`, which creates the scenario shown in Fig. 2.12. The application used is iPerf, and the mobile device creates two uplink subflows to a

remote server, the first on mmWave, and the second on LTE. The UE moves along the y-axis and switches from a LOS to a NLOS condition, and then returns to LOS. Fig. 2.13 shows the output of a simulation, with the TCP throughput for two different congestion control algorithms for MPTCP, together with the per-subflow Radio Access Network throughput. In particular, Fig. 2.13a shows the performance of Balanced Link Adaptation Algorithm (BALIA) [145], which is a coupled congestion control algorithm, i.e., it tries to adapt the congestion window of each MPTCP subflow according to the congestion experienced on all links.

In Fig. 2.13b, instead, the congestion controls on the LTE and on the mmWave subflows are uncoupled, i.e., each subflow is independent, and TCP Cubic is used. The first observation is that the LTE subflow has, as expected, a much smaller throughput compared to the mmWave subflow, and thus the total throughput measured with iPerf is similar to that of the mmWave connection. The second is that the uncoupled solution manages to reach a more stable throughput in NLOS conditions, compared to the coupled solution, as was observed in [389, 404], showing that the current coupled congestion control algorithms are not well suited for a deployment over these kinds of links.

## 2.12 Potential uses and future extensions

This simulation tool was developed and extensively used to evaluate the solutions proposed throughout this thesis. Moreover, given the full-stack nature of the simulation framework, the 5G mmWave research community has also started to leverage this tool to bring and test innovation at every layer [81, 146–151]. Each model can be easily extended while maintaining full backward compatibility. The fundamental components are in the form of functions, classes and interfaces, which can be implemented to design novel algorithms, procedures, and, more in general, architectures. For example, the *scheduling* and *allocation* strategies proposed in [152–154] can be readily integrated and tested in our framework with some simple tweaks. Similarly, due to the importance of coping with *mobility* and frequent *handovers*, innovative approaches like the one proposed in [155], which exploits *caching*, can take advantage of the modular structure of the ns–3 framework to test flexible and reprogrammable logics. Additionally, as previously mentioned, several papers already fully exploit the capabilities of this simulator to capture the performance of *TCP* in mmWave networks, and propose some novel approaches to mitigate the limiting effects of congestion control procedures with intermittent multi-Gbps mmWave links [57, 64, 147, 409, 411].

As part of our future work, we aim at expanding the code to include additional components such as:

- 3GPP-inspired *signaling/beamtracking* procedures [156] to better accommodate novel techniques like those proposed in [154, 157];

- novel applications such as *virtual & augmented reality*, to ultimately test key 5G metrics as done in [151], where the authors leverage our mmWave module to run a performance analysis of traditional video delivery over mmWaves, and in [413], where ns-3 mmWave is used to assess the performance of network coding and multi connectivity for reliable video streaming over mmWave;

- *vehicular* channel and traffic models to test and capture the end-to-end performance of mmWave communications for high-mobility scenarios [158–160];

- *public safety* scenarios [422], including aerial communications and robotics, where the propagation environments and the performance requirements differ from those of traditional cellular networks, as detailed in [391, 421].

- improvement of ns-3 *mobility models*. The realistic modeling of the mobility patterns of user terminals is fundamental for a proper evaluation of the network performance. The ns-3 `mobility` module provides several classes implementing different mobility models, either random, e.g., `RandomWalk2D`, `GaussMarkov`, `RandomWaypoint`, or deterministic, e.g., `ConstantVelocity`, `Waypoint`. While the deterministic models can be used only when the mobility patterns of the moving terminals are completely known, the random models provides the possibility to account for the user mobility in a statistical way. However, in some cases they may be too general to properly consider the complex temporal and spatial correlations which characterize the human motion [161], and may not be suitable to simulate emerging use cases, such as vehicular and aerial deployments [162]. Moreover, the device rotation could be even more problematic than mobility alone in high frequency communications. A basic implementation of the device rotation could simply follow the direction in which the user is moving. To further increase the realism, a choice of statistical models for device rotation could also be added. To integrate the device rotation with the beamforming and channel modeling, a complete framework should be introduced, thus separating the global frame of reference from local ones.

- improvement of ns-3 *building model*. Buildings and obstacles can severely impact the quality of mobile communications. As already mentioned, ns-3 features a module that implements building in the scenarios, and is used by different propagation models (e.g., `Hybrid-BuildingsPropagationLossModel` for LTE and the `MmWave3gppBuildingPropagation-LossModel`) to tune the propagation condition to the presence of a NLOS condition or an urban canyon. Nonetheless, the objects of the `Buildings` cannot be moved during the simulation runtime, and, if a mobility model makes a user enter and exit a building, its indoor/outdoor status is not properly handled.[13] With respect to the integration with mobility models, they would benefit from a tighter integration with the buildings, as we propose in [422].

Moreover, we plan to address the challenge of *scalability*. 5G networks will likely comprise a large number of nodes, with high mobility, and thus channel states must be updated frequently. In addition, the use of low-latency applications requires that packet timelines must be scheduled at very fast interaction times. In general, there is a trade-off between the accuracy and the complexity of the simulations. For example, the high level of detail in the scheduling of synchronization signals for the cellular stacks may yield an excessive overhead in terms of simulation complexity. Nonetheless, this could be avoided when the users are not using the network resources (e.g., in the case of machine-to-machine communications). Therefore, it makes sense to investigate and implement mechanisms that allow the simulation to reduce the number of scheduled events and consequently the simulation runtime, e.g., in case of bursty and sporadic traffic. Even though it should be noted that this approach is not always possible (e.g., in scenarios with non-bursty and heavy traffic), finding scalable approaches where the simulation complexity is minimized according to the simulated scenario is a worthwhile research direction.

## 2.13 Conclusions

In this chapter, we have presented the current status of the ns–3 framework for simulation of mmWave cellular systems. The code, which is publicly available at GitHub[14] and in the official

---

[13]https://www.nsnam.org/bugzilla/show_bug.cgi?id=3018
[14]https://github.com/nyuwireless-unipd/ns3-mmwave

ns-3 App Store[15], is highly modular and customizable to allow researchers to test novel 5G protocols. We have shown some performance trends based on the mmWave channel models available. A detailed explanation of our configurable physical and MAC layers is provided, along with a corroborating set of simulation results for varying configurations. Implementations of advanced 5G architectural features, such as dual connectivity, are also available, and we have reported different representative results. We have also shown that the module can be interfaced with the higher-layer protocols and core network models from the ns–3 LTE module to enable full-stack simulations of end-to-end connectivity, along with the simulation of real applications through the implementation of direct code execution. The module is demonstrated through several example simulations showing the performance of our custom mmWave stack as well as custom congestion control algorithms, specifically designed for efficient utilization of the mmWave channel.

---

[15]`https://apps.nsnam.org/app/mmwave/`

# Part II

# The Architecture: System Level Design of 5G mmWave Networks

# 3

# Multi Connectivity

## 3.1  Introduction

Mobility management is a key challenge for cellular networks at mmWave frequencies. Indeed, as mentioned in Sec. 1.2, the UE mobility, combined with small movements of obstacles and reflectors, or even changes in the orientation of a handset relative to the body or a hand, can cause the channel to rapidly appear or disappear. Moreover, the high propagation loss calls for a small cell and heterogenous deployment. Therefore, the number of handovers and beam switch events in mmWave systems is expected to be high, and thus there is a need for an efficient mobility management framework that allows the end users to quickly update the serving beam pair or switch to another base station [42].

One of the main tools to improve the robustness of mmWave systems is *multi-connectivity* [61]: each mobile device (UE or user equipment in 3GPP terminology) maintains connections to multiple cells, possibly including both 5G mmWave cells and/or conventional 4G cells. In the event that one link is blocked, the UE can find alternate routes to preserve the connection. In cellular systems, this robustness is called *macro-diversity* and is particularly vital for mmWave systems [61].

How to implement multi-connectivity in the network layer for mmWave systems remains largely an open problem. Current 3GPP cellular systems offer multiple mechanisms for fast switching of paths between different cells including conventional handover, multi-connectivity and carrier aggregation – these methods are summarized below. However, mmWave systems present unique challenges:

- Most importantly, the dynamics of mmWave channels imply that the links to any one cell can deteriorate rapidly, necessitating much faster link detection and re-routing [42].

- Due to the high isotropic pathloss, mmWave signals are transmitted in narrow beams, typically formed with high-dimensional phased arrays. In any link, channel quality must be continuously scanned across multiple possible directions which can dramatically increase the time it takes to detect that the link has failed and a path switch is necessary [62].

- One of the main goals of 5G is to achieve ultra-low latency [6] (possibly $< 1$ ms). Thus, service unavailability during path switches must be kept to a minimum.

### 3.1.1 Contributions

To address these challenges, in this chapter, which is based on $[388, 402]^1$, we expand the main results and findings of our previous works [62, 98, 163], to provide the first global end-to-end comprehensive evaluation of a mobility management framework for handover and path switching of mmWave systems under realistic dynamic scenarios, and assess how a dual-connectivity[2] (DC) approach can enable faster, more robust and better performing mobility management schemes.

In particular, in [62] we proposed a novel multi-connectivity uplink measurement framework that, with the joint effort of the legacy LTE frequencies, enables fair and robust cell selection, in addition to efficient and periodic tracking of the user, suitable for several control-plane cellular applications (i.e., we showed that periodic measurement reports can be used to trigger handovers or adapt the beams of the user and its serving cell, to grant good average throughput and deal with the channel dynamics experienced at mmWave frequencies). In [98], we evaluated the tracking performance of a user's signal quality considering real experiments in common blockage scenarios, combined with outdoor statistical models. Finally, in [402], we discussed two possible ways to integrate 5G and LTE networks to improve the reliability of next-generation mobile users, and described a preliminary ns-3 simulation framework to evaluate the performance of both.

By extending our previous contributions, in this chapter we also propose:

- The use of a DC scheme to enable the base stations to efficiently track the UE channel quality along multiple links and spatial directions within those links. In addition, to allow fast detection of link failures, we demonstrate that the uplink control signaling enables the network to track the angular directions of communication to the UE on all possible links simultaneously, so that, when a path switch is necessitated, no directional search needs to be performed (this approach greatly saves switch time, since directional scanning dominates the delay in establishing a new link [164, 165]).

- The use of a local coordinator that manages the traffic between the cells. The coordinator performs both control plane tasks of path switching and data plane tasks as a traffic anchor, at the PDCP layer. In conventional cellular systems, these control and data plane functions are performed in the MME and Service Gateway (SGW), which are often far from the cells. In contrast, the local coordinator is placed in close proximity to the cells, significantly reducing the path switch time.

- The design of faster network handover procedures (namely *fast switching* and *secondary cell handover*) that, by exploiting our DC framework, improve the mobility management in mmWave networks, with respect to the standard standalone Hard Handover (HH) scheme. These procedures are controlled by the LTE RRC layer and, since the UE is connected to the LTE and the mmWave eNBs, it is possible to perform quick fallback to LTE with the fast switching command.

- A dynamic Time-to-Trigger (TTT) adaptation to enhance the switch decision timing in highly uncertain link states.

Moreover, we evaluate the proposed switching and handover protocols by extending the evaluation methodology we have described in Chapter 2. The ns-3-based framework we implemented

---

[1]Part of this chapter is based on joint work with Marco Giordani.

[2]Although many of the ideas and techniques discussed in this study apply to more general multi-connectivity scenarios, for concreteness in the following we will specifically refer to *dual* connectivity, in which a UE is simultaneously connected to one 5G mmWave base station and one legacy LTE eNB.

for this work makes it possible to use detailed measurement-based channel models that can account for both the spatial characteristics of the channel and the channel dynamics arising from blocking and other large-scale events, which is important for a detailed and realistic assessment. We believe that this is the first exhaustive contribution which provides a global evaluation of the performance of a dual-connectivity architecture with respect to a traditional standalone HH scheme in terms of handover and mobility management specifically tailored to a dynamic mmWave scenario. In particular, we simulated the user's motion in a typical urban environment. Separately, actual local blockage dynamics were measured and superimposed on the statistical channel model, to obtain a realistic spatial dynamic channel model. We believe that this is the first work in which such detailed mmWave dynamic models have been used in studying handover.

Our study reveals several important findings on the interaction of transport layer mechanisms, buffering, and its interaction with physical-layer link tracking and handover delays. We also demonstrate that the proposed dual connectivity framework offers significant performance improvements in the handover management of an end-to-end network with mmWave access links, including (i) reduced packet loss, (ii) reduced control signaling, (iii) reduced latency, and (iv) higher throughput stability. Moreover, we show that a dynamic TTT approach should be preferred for handover management, since it can deliver non-negligible improvements in specific mobility scenarios in which state-of-the-art methods fail.

### 3.1.2 Related Work

Dual connectivity to different types of cells (e.g., macro and pico cells) has been proposed in Release 12 of Long Term Evolution-Advanced (LTE-A) [136] and in [166]. However, these systems were designed for conventional sub-6 GHz frequencies, and the directionality and variability of the channels typical of mmWave frequencies were not addressed. Some other previous works, such as [167], consider only the bands under 6 GHz for the control channel of 5G networks, to provide robustness against blockage and a wider coverage range, but this solution could not provide the high capacities that can be obtained when exploiting mmWave frequencies. The potential of combining legacy and mmWave technologies in outdoor scenarios has also been investigated in [134], highlighting the significant benefits that a mmWave network achieves with flexible, dynamic support from LTE technologies. Articles [63,133] propose a multi-connectivity framework as a solution for mobility-related link failures and throughput degradation of cell-edge users, enabling increased reliability with different levels of mobility.

Although the literature on handover in more traditional sub-6 GHz heterogeneous networks is quite mature, papers on handover management for mmWave 5G cellular are very recent, and research in this field has just started. The survey in [168] presents multiple vertical handover decision algorithms that are essential for heterogeneous wireless networks, while article [169] investigates the management of the handover process between macro, femto and pico cells, proposing a theoretical model to characterize the performance of a mobile user in heterogeneous scenarios as a function of various handover parameters. However, these works are focused on low frequency legacy cellular systems. When dealing with mmWaves, frequent handover, even for fixed UEs, is a potential drawback that needs to be addressed. In [170], the handover rate in 5G systems is investigated and in [171] a scheme for handover management in high-speed railway is proposed by employing the received signal quality from measurement reports. In [172, 173] the impact of user mobility in multi-tier heterogeneous networks is analyzed and a framework is proposed to solve the dynamic admission and mobile association problem in a wireless system with mobility. Finally, the authors of [174] present an architecture for mobility, handover and routing management.

**Figure 3.1:** LTE-5G tight integration architecture.

## 3.2 Framework Description For Dual Connectivity

We propose a dual connectivity architecture, introduced here for the control and user planes as an extension of 3GPP's LTE DC proposal [136] to the needs of mmWave communications. In the proposed solution, the UE is simultaneously connected to both LTE and mmWave eNBs. The LTE cell is a backup for the user plane: since the UE is already connected, when the signal quality of the mmWave link degrades, there is no need to perform a complete handover; a single RRC control message from the LTE eNB to the UE is enough. Moreover, for the control plane, this scheme enables a coordinated measurement collection as described in [62, 163]. Fig. 3.1 shows a block diagram of the proposed architecture, as presented in [402]. For each DC device there is a single connection point to the Core Network (CN), through the S1 interface that links the LTE eNB to the CN: the mmWave eNB does not exchange control messages with the MME. The two eNBs are connected via an X2 link, which may be a wired or wireless backhaul. Each LTE eNB coordinates a cluster of mmWave eNBs which are located under its coverage. Notice that the coordinator may also be placed in a new node in the core network, or can be based on NFV logic.

In the following paragraphs, we will present in detail how the DC framework enables (i) channel monitoring over time, (ii) a PDCP layer integration across different radio access networks, and (iii) faster network handover procedures.

### 3.2.1 Control Plane For Measurement Collection

Monitoring the channel quality is an essential component of any modern cellular system, since it is the basis for enabling and controlling many network tasks including rate prediction, adaptive modulation and coding, path selection and also handover. In this work, we follow the multi-cell measurement reporting system proposed in [62, 163], where each UE directionally broadcasts a SRS in a time-varying direction that continuously sweeps the angular space. Each potential serving cell scans all its angular directions and monitors the strength of the received SRS,

| RT (mmWave eNB$_j$) | |
| --- | --- |
| UE$_1$ | SINR$_{1,j}$ |
| UE$_2$ | SINR$_{2,j}$ |
| ... | ... |
| UE$_N$ | SINR$_{N,1}$ |

| Complete Report Table (CRT) | | | |
| --- | --- | --- | --- |
| UE | mmWave eNB$_1$ | ... | mmWave eNB$_M$ |
| UE$_1$ | SINR$_{1,1}$ | ... | SINR$_{1,M}$ |
| UE$_2$ | SINR$_{2,1}$ | ... | SINR$_{2,M}$ |
| ... | ... | ... | ... |
| UE$_N$ | SINR$_{N,1}$ | ... | SINR$_{N,M}$ |

**Table 3.1:** An example of RT (left) and CRT (right), referred to $N$ users and $M$ available mmWave eNBs in the network. We suppose that the UE can send the sounding signals through $N_{\mathrm{UE}}$ angular directions and each mmWave eNB can receive them through $N_{\mathrm{eNB}}$ angular directions. Each pair is the maximum SINR measured in the best direction between the eNB and the UE.

building a report table (RT) based on the channel quality of each receiving direction, to better capture the dynamics of the channel[3]. A centralized coordinator (which may reside in the LTE eNB) obtains complete directional knowledge from all the RTs sent by the potential cells in the network to make the optimal serving cell selection and scheduling decisions. In particular, due to the knowledge gathered on the signal quality in each angular direction for each eNB-UE pair, the coordinator is able to match the beams of the transmitter and of the receiver to provide maximum performance.

In this work, we assume that nodes select one of a finite number of directions for measuring the signal quality, and we let $N_{\mathrm{eNB}}$ and $N_{\mathrm{UE}}$ be the number of directions at each eNB and UE, respectively. Supposing that $M$ cells are deployed within the coverage of the coordinator, the procedure works as follows.

### First Phase – Uplink Measurements

Each UE directionally broadcasts uplink sounding reference signals in dedicated slots, steering through directions $d_1, \ldots, d_{N_{\mathrm{UE}}}$, one at a time, to cover the whole angular space. The SRSs are scrambled by locally unique identifiers (e.g., C-RNTI) that are known to the mmWave eNBs and can be used for channel estimation. If analog beamforming is used, each mmWave eNB scans through directions $D_1, \ldots, D_{N_{\mathrm{eNB}}}$ one at a time or, if digital beamforming is applied, collects measurements from all of them at once. Each mmWave eNB fills a RT, as in Table 3.1 left, whose entries represent the highest SINR between UE$_i$, $i = 1, \ldots, N$, transmitting through its best direction $d_{\mathrm{UE,opt}} \in \{d_1, \ldots, d_{N_{\mathrm{UE}}}\}$, and eNB$_j$, $j = 1, \ldots, M$, receiving through its best possible direction $D_{\mathrm{eNB,opt}} \in \{D_1, \ldots, D_{N_{\mathrm{eNB}}}\}$:

$$\mathrm{SINR}_{i,j} = \max_{\substack{d_{\mathrm{UE}} = d_1, \ldots, d_{N_{\mathrm{UE}}} \\ D_{\mathrm{eNB}} = D_1, \ldots, D_{N_{\mathrm{eNB}}}}} \mathrm{SINR}_{i,j}(d_{\mathrm{UE}}, D_{\mathrm{eNB}}) \tag{3.1}$$

### Second Phase – Coordinator Collection

Once the RT of each mmWave eNB has been filled for each UE, each mmWave cell sends this information, through the X2 link, to the coordinator[4] which, in turn, builds a complete report table (CRT), as depicted in Table 3.1 right. When accessing the CRT, the optimal mmWave eNB

---

[3]Unlike in traditional LTE systems, the proposed framework is based on the channel quality of uplink (UL) rather than downlink (DL) signals. This eliminates the need for the UE to send measurement reports back to the network and thereby removes a possible point of failure in the control signaling path.

[4]The complexity of this framework resides in the central coordinator, which has to aggregate the RT from the $M$ mmWave eNBs that are under its control and perform for each of the $N$ UEs a search operations among $M$ entries. As the number of the mmWave eNBs $M$ increases, the search space increases linearly.

| BF Architecture | | Delay $D$ |
|---|---|---|
| mmWave eNB Side | UE Side | |
| Analog | Analog | 25.6 ms |
| Hybrid | Analog | 25.6/$L$ ms |
| Digital | Analog | 1.6 ms |

**Table 3.2:** Delay $D$ for each mmWave eNB to fill each RT. A comparison among different BF architectures (analog, hybrid and fully digital) is reported. We assume $T_{\text{sig}} = 10\ \mu$s, $T_{\text{per}} = 200\ \mu$s (to maintain an overhead $\phi_{\text{ov}} = 5\%$), $N_{\text{UE}} = 8$ and $N_{\text{eNB}} = 16$.

(with its optimal direction $D_{\text{eNB,opt}}$) is selected for each UE (with optimal direction $d_{\text{UE,opt}}$), considering the absolute maximum SINR in each CRT's row. The criterion with which the best mmWave eNB is chosen will be described in Sec. 3.2.3.

Third Phase – Network Decision

The coordinator reports to the UE, on a legacy LTE connection, which mmWave eNB yields the best performance, together with the optimal direction $d_{\text{UE,opt}}$ in which the UE should steer its beam, to reach the candidate serving mmWave eNB in the optimal way. The choice of using the LTE control link is motivated by the fact that the UE may not be able to receive from the optimal mmWave link if not properly configured and aligned. Moreover, since path switches and handover events in the mmWave regime are commonly due to link failures, the control link to the serving mmWave cell may not be available. Finally, the coordinator also notifies the designated mmWave eNB, through the X2 link, about the optimal direction $D_{\text{eNB,opt}}$ in which to steer the beam, for serving each UE. We highlight that the procedure described in this section allows to optimally adapt the beam even when a handover is not strictly required. In particular, if the user's optimal mmWave eNB is the same as the current one, but a new steering direction pair $(d_{\text{UE,opt}}, D_{\text{eNB,opt}})$ is able to provide a higher SINR to the user, a beam switch is prompted, to realign with the eNB and guarantee better communication performance.

According to [175], we assume that the SRSs are transmitted periodically once every $T_{\text{per}} = 200\ \mu$s seconds, for a duration of $T_{\text{sig}} = 10\ \mu$s seconds (which is deemed sufficient to allow proper channel estimation at the receiver), to maintain a constant overhead $\phi_{\text{ov}} = T_{\text{sig}}/T_{\text{per}} = 5\%$. The switching time for beam switching is in the scale of nanoseconds, and so it can be neglected [176]. The scanning for the SRSs for each UE-eNB direction and the filling of each RT require $N_{\text{eNB}}N_{\text{UE}}/L$ scans, where $L$ is the number of directions in which the receiver can look at any one time. Since there is one scanning opportunity every $T_{\text{per}}$ seconds, the total delay is

$$D = \frac{N_{\text{eNB}}N_{\text{UE}}T_{\text{per}}}{L}. \tag{3.2}$$

The value of $L$ depends on the beamforming (BF) capabilities. In the uplink-based design, $L = 1$ if the eNB receiver has analog BF and $L = N_{\text{eNB}}$ if it has a fully digital transceiver. According to Eq. (3.2), the value of $D$ is independent of the number of users and of the MAC layer scheduling. Since each UE sends its sounding reference signals at the same time and the mmWave eNBs synchronously receive those messages through exhaustive search schemes, the proposed framework scales well with the network density.

Table 3.2 reports the delay $D$ for different configurations of a system with $N_{\text{UE}} = 8$ and $N_{\text{eNB}} = 16$ directions required to collect each instance of the CRT at the LTE eNB side by implementing the framework described above. For example, by implementing a hybrid BF with

$L = 2$ RF chains, the eNB can simultaneously receive through $L = 2$ directions at the same time [43] so the overall delay is $D = 12.8$ ms.

From the protocol stack point of view, unlike in [136], both Radio Access Technologies (RATs) have a complete RRC layer in the eNBs and in the UE. This allows a larger flexibility, since the design of the mmWave RRC layer can be decoupled from that of the LTE stack. Moreover, the LTE RRC is used for the management of the LTE connection but also to send and receive commands related to DC, while the mmWave RRC is used to manage only the mmWave link and the reporting of measurements to the coordinator. The choice of using a dedicated RRC link for the secondary eNB is motivated by the desire to reduce the latency of control commands, since it avoids the encoding and transmission of the control PDUs of the secondary cell to the master cell. The mmWave signaling radio bearers are used only when a connection to LTE is already established, and this can offer a ready backup in case the mmWave link suffers an outage.

### 3.2.2 User Plane (PDCP Layer Integration)

In a DC architecture, the layer at which the LTE and the mmWave protocol stacks merge is called *integration layer*. In this study we propose the PDCP layer as the integration layer. In fact, it allows a non co-located[5] deployment, since synchronization among the lower layers is not required, and it does not impose any constraint on the design of mmWave PHY to RLC layers, so that a clean slate approach can be used to address mmWave specific issues and reach 5G performance requirements.

For each bearer, a PDCP layer instance is created in the LTE eNB and interfaced with the X2 link that connects to the remote eNB. Local and remote RLC layer instances are created in the LTE and mmWave eNB, respectively. The packets are routed from the S-GW to the LTE eNB, and once in the PDCP layer they are forwarded either to the local LTE stack or to the remote mmWave RLC. If there exists at least one mmWave eNB not in outage and the UE is connected to it, then the mmWave RAT is chosen, i.e., the LTE connection is used only when no mmWave eNB is available. This choice is motivated by the fact that the theoretical capacity of the mmWave link is greater than that of the LTE link [33], and that the LTE eNB will typically serve more users than the mmWave eNBs; however, when the mmWave eNBs are in outage (as it may happen in a mmWave context) and would therefore provide zero throughput to their users, an LTE link may be a valid fallback alternative to increase the robustness of the connection. In addition, integration at the PDCP layer ensures ordered delivery of packets to the upper layers, which is useful in handover circumstances.

### 3.2.3 Dual Connectivity-aided Network Procedures

The DC framework allows to design network procedures that are faster than the standard standalone hard handover (HH), thus improving the mobility management in mmWave networks. The standalone HH architecture will be the baseline for the performance evaluation of Sec. 3.4: the UE is connected to either the LTE or the mmWave RAT and, in order to switch from one to the other, it has to perform a complete handover, or, if the mmWave connectivity is lost, an initial access to LTE from scratch. Besides, in order to perform a handover between mmWave eNBs, the UE has to interact with the MME in the core network, introducing additional delays.

---

[5]MmWave eNBs will be deployed more densely than already installed LTE eNBs, therefore it would be costly to have only co-located cells. Moreover a high density of LTE eNBs would decrease the effectiveness of the coverage layer. Finally, the PDCP layer can also be deployed in the core network, in a new node (*coordinator*), which can be a gateway for a cluster of LTE eNBs and the mmWave eNBs under their coverage, or can be deployed in a macro LTE cell.

**(a)** Switch from LTE RAT to mmWave RAT.



**(b)** Switch from mmWave RAT to LTE RAT.

**Figure 3.2:** Proposed RAT switch procedures.



**Figure 3.3:** Secondary cell Handover procedure (SCH).

The DC architecture, instead, allows to perform fast switching between the LTE and mmWave RATs and SCH across mmWave eNBs.

The fast switching procedure is used when all the mmWave eNBs for a certain UE are in outage. Since the handling of the state of the user plane for both the mmWave and the LTE RATs is carried out by the LTE RRC, it is possible to correctly modify the state of the PDCP layer and perform a switch from the mmWave to the LTE RAT. The proposed switch procedure, shown in Fig. 3.2, simply requires an RRC message (RRC Connection Switch command) to the UE, sent on the LTE link, and a notification to the mmWave eNB via X2 if the switch is from mmWave to LTE, in order to forward the content of the RLC buffers to the LTE eNB.

The DC solution therefore allows to have an uninterrupted connection to the LTE anchor point. However, it is possible to switch from a secondary mmWave eNB to a different mmWave eNB with a procedure which is faster than a standard intra RAT handover, since it does not involve the interaction with the core network. The Secondary Cell Handover procedure is shown in Fig. 3.3. The Random Access (RA) procedure [119] is aided by the measurement collection framework described in Sec. 3.2.1, which allows to identify the best beam to be used by the UE and avoids the need for the UE to perform an initial beam search. Moreover, if the UE is capable of maintaining timing control with multiple mmWave eNBs, the RA procedure in the target mmWave eNB can be skipped.

We also propose an algorithm for SCH, based on the SINR measurements reported by the

mmWave eNBs to the coordinator and on a threshold in time (TTT). When a mmWave eNB has a better SINR than the current one (and neither of the two is in outage), the LTE coordinator checks for TTT seconds if the condition still holds, and eventually triggers the SCH. Notice that, if during the TTT the SINR of a third cell becomes better than that of the target cell by less than 3 dB, the handover remains scheduled for the original target eNB, while, if the original cell SINR becomes the highest, then the SCH is canceled. The TTT is computed in two different ways. With the *fixed* TTT option it always has the same value[6] (i.e., $f_{TTT} = 150$ ms), while for the *dynamic* TTT case we introduce a dependency on the difference $\Delta$ between the SINRs of the best and of the current cell:

$$f_{TTT}(\Delta) = TTT_{max} - \frac{\Delta - \Delta_{min}}{\Delta_{max} - \Delta_{min}}(TTT_{max} - TTT_{min}) \qquad (3.3)$$

so that the actual TTT value is smaller when the difference in SINR between the current eNB and the target is higher. The parameters that were used in the performance evaluation carried out in this study are $TTT_{max} = 150$ ms, $TTT_{min} = 25$ ms, $\Delta_{min} = 3$ dB, $\Delta_{max} = 8$ dB.

Finally, if at a given time all the mmWave eNBs are in outage, then the UE is instructed to switch to the LTE eNB. If instead only the current mmWave eNB is in outage, the UE immediately performs a handover to the best available mmWave eNB, without waiting for a TTT.

## 3.3    Performance Evaluation Framework

In order to assess the performance of the proposed DC architecture with respect to the traditional standalone hard handover (HH) baseline we use ns–3-based system level simulations, based on the DC framework described in Chapter 2, with the DC extension discussed in Sec. 2.9. This approach has the advantage of including many more details than would be allowed by an analytical model (which, for such a complex system, would have to introduce many simplifying assumptions), and makes it possible to evaluate the system performance accounting for realistic (measurement-based) channel behaviors and detailed (standard-like) protocol stack implementations.

The source code of the DC framework is publicly available[7], as well as the ns–3 script (`mc-example-udp.cc`) used for the simulation scenario considered in this chapter.

### 3.3.1    Semi-Statistical Channel Model

The channel model we specifically implemented for this chapter is based on recent real-world measurements at 28 GHz in New York City, to provide a realistic assessment of mmWave micro and picocellular networks in a dense urban deployment [37, 96, 177, 178]. Unfortunately, most of the studies have been performed in stationary locations with minimal local blockage, making it difficult to estimate the rapid channel dynamics that affect a realistic mmWave scenario. Dynamic models such as [179] do not yet account for the spatial characteristics of the channel.

Measuring a wideband spatial channel model with dynamics is not possible with our current experimental equipment, as such measurements would require that the transmitting and receiving directions be swept rapidly during the local blocking event. Since our available platform relies on horn antennas mounted on mechanically rotating gimbals, such rapid sweeping is not possible.

In this work, we follow the alternate approximate *semi-statistical* method proposed in [98] to generate realistic dynamic models for link evaluation:

---

[6]This approach recalls the standard HO for LTE networks.
[7]https://github.com/nyuwireless/ns3-mmwave/

(i) We first randomly generate the statistical parameters of the mmWave channel, according to [37] and [96], which would reflect the characteristics of a stationary ground-level mobile with no local obstacles.

(ii) Since there are no statistical models for the blocking dynamics, local blocking events are measured experimentally and modulated on top of the static parameters, in case an obstacle is physically deployed through the path that links the UE to one of the mmWave eNBs[8].

Handover decisions described in Sec. 3.2.3 are based on the SINR values saved in the CRT, built at the coordinator's side. Specifically, the SINR between a mmWave $eNB_j$ and a test UE can be computed in the following way:

$$\text{SINR}_{j,\text{UE}} = \frac{\frac{P_{\text{TX}}}{PL_{j,\text{UE}}} G_{j,\text{UE}}}{\sum_{k \neq j} \frac{P_{\text{TX}}}{PL_{k,\text{UE}}} G_{k,\text{UE}} + W_{\text{tot}} \times N_0} \tag{3.4}$$

where $G_{i,\text{UE}}$ and $PL_{i,\text{UE}}$ are the beamforming gain and the pathloss obtained between $eNB_i$ and the UE, respectively, $P_{\text{TX}}$ is the transmit power and $W_{\text{tot}} \times N_0$ is the thermal noise power.

In the following, we describe in detail how the real experiments in common blockage scenarios are combined with the outdoor statistical model for ns–3, to get a realistic expression for the SINR samples which takes into account the dynamics experienced in a mmWave channel.

**MmWave Statistical Channel Model**

The parameters of the mmWave channel that are used to generate the time-varying channel matrix $\mathbf{H}$ include: (i) spatial clusters, described by central azimuth and elevation angles; (ii) fractions of power; (iii) angular beamspreads; and (iv) small-scale fading, which models every small movement (e.g., a slight variation of the handset orientation) and is massively affected by the Doppler shift and the real-time position (AoA, AoD) of the UE, which may change very rapidly, especially in dense and high-mobility scenarios (for this reason, we chose to adapt the channel's small scale fading parameters as frequently as possible, that is once every time slot of 125 $\mu s$).

These parameters are defined and explained in [37, 96], while a complete description of the channel model can be found in [163]. Notice that, following the approach of [390], the large scale fading parameters of the $\mathbf{H}$ matrix are updated every 100 ms, to simulate a sudden change of the link quality.

The pathloss is defined as $PL(d)[dB] = \alpha + \beta 10 \log_{10}(d)$, where $d$ is the distance between the receiver and the transmitter and the values of the parameters $\alpha$ and $\beta$ are given in [37]. In case an obstacle is obstructing the path that links the UE and a specific mmWave eNB in the network, a NLOS pathloss state is emulated by superimposing the experimentally measured blockage traces to the statistical realization of the channel, as explained in Sec. 3.3.1.

When just relying on the statistical characterization of the mmWave channel, the SINR expression obtained by applying Equation (3.4) assumes a baseline LOS pathloss where no local obstacles affect the propagation of the signal. In the next paragraph, a channel sounding system

---

[8] An important simplification is that we assume that the local blockage equally attenuates all paths, which may not always be realistic. For example, a hand may block only paths in a limited number of directions. However, in any fixed direction, most of the power is contributed only by paths within a relatively narrow beamwidth and thus the approximation that the paths are attenuated together may be reasonable.

is presented for measuring the dynamics of the blockage.

**Measurement of Local Blockage**

The key challenge in measuring the dynamics of local blockage is that we need relatively fast measurements. To perform these fast measurements, we used a high-bandwidth baseband processor, built on a PXI (a rugged PC-based platform for measurement and automation systems) from National Instruments, which engineers a real-world mmWave link. A detailed description of the experimental testbed can be found in [98] and [180].

Using this system, the experiments were then conducted by placing moving obstacles (e.g., a person walking or running) between the transmitter and the receiver, and continuously collecting Power Delay Profile (PDP) samples during each blocking event[9]. The experiments show that obstacles can cause up to 35-40 dB of attenuation with respect to the LOS baseline SINR values, and this local blocking attenuation factor is thus used to modulate the time-varying channel response from the statistical channel model.

**Final Semi-Statistical SINR Trace**

Once a statistical instance of $\text{SINR}_{j,\text{UE}}$ is obtained from Equation (3.4), a raw estimate of the real SINR at the UE is derived by superimposing the local blocking dynamics measured experimentally, when an obstacle is physically present in the path between that UE and $\text{eNB}_j$. In particular, we denote by $\Gamma_{\text{stat},j}$ the maximum static SINR between the $\text{eNB}_j$ and the UE receiver, when assuming that no local obstacles are present. Then the maximum wideband SINR when also considering a dynamic model for the link evaluation (that is the value inserted in the $j$-th column of the CRT, at a specific time instant) is obtained as:

$$\Gamma_j = \begin{cases} \Gamma_{\text{stat},j} & \text{if no obstacles are in the path between UE and } \text{eNB}_j \text{ (LOS} \\ & \text{condition)} \\ \delta + \Gamma_{\text{stat},j} & \text{if an obstacle is in the path between UE and } \text{eNB}_j \text{ (NLOS} \\ & \text{condition)} \end{cases} \tag{3.5}$$

where $\delta$ is a scaling factor that accounts for the SINR drop measured experimentally in various blocking scenarios and collected using the instrumentation described in the previous paragraph.

This final semi-statistical SINR trace is composed of samples of $\Gamma_j$ generated every 125 $\mu s$ (from both the statistical trace and the experimental measurements). Finally, according to Sec. 3.2, the HO decisions are made once the coordinator has built a CRT, that is every $D$ seconds. Thus, the original SINR trace has been downsampled, keeping just one sample every $D$ seconds.

### 3.3.2   SINR Filtering

The mmWave eNBs estimate the wideband SINR $\Gamma_j$ from the sounding reference signals that are transmitted by the UE and are collected by each mmWave eNB, to build a CRT at the

---

[9]PDPs were measured at a rate of one PDP every 32 $\mu s$ but, since we found that the dynamics of the channel varied considerably slower than this rate, we decimated the results by a factor of almost four, recording one PDP every about 125 $\mu s$, that matches the slot duration of the ns-3 framework.

**Figure 3.4:** SINR evolution, with respect to a specific mmWave eNB in the network, whose samples are collected every $D = 1.6$ ms, according to the measurement framework described in Sec. 3.2.1. Each sample is obtained by following the semi-statistical channel model proposed in [98] and explained in this section. The red line is referred to the true SINR trace $\Gamma$, while the black line is referred to its estimate $\bar{\Gamma}$, after noise and a first-order filter are applied to the true SINR $\Gamma$.

coordinator's side[10]. However, the *raw estimate* of the SINR $\hat{\Gamma}_j$, that is what is really measured in a realistic communication system, may deviate from $\Gamma_j$ due to noise (whose effect can be very significant when considering low SINR regimes). To reduce the noise, $\hat{\Gamma}_j$ is filtered, producing a time-averaged SINR trace $\bar{\Gamma}_j$. According to [98], a simple first-order filter can properly restore the desired SINR stream and perform reliable channel estimation even without designing more complex and expensive adaptive nonlinear filters. Therefore, $\bar{\Gamma}_j$ is obtained as

$$\bar{\Gamma}_i = (1 - \eta)\bar{\Gamma}_{i-1} + \eta\hat{\Gamma}_i, \tag{3.6}$$

for some constant $\eta \in (0, 1)$ chosen in order to minimize the estimation error $e_i = |\Gamma_i - \bar{\Gamma}_i|^2$.

As an example, the SINR trace in Fig. 3.4 (whose samples are collected every $D = 1.6$ ms) is obtained by following the semi-statistical channel model proposed in [98] and explained in this section. For time $t < 19.4$, the UE is in a NLOS pathloss condition with respect to its eNB, therefore a scaling factor $\delta$ measured experimentally is applied to the statistical trace to account for the dynamics of the local blockage. For time $t > 19.4$, the UE enters a LOS state until the end of the simulation. The SINR collapses and spikes within the trace (i.e., at times $t = 18.1$ or $t = 18.9$) are mainly caused by the update of the large scale fading parameters of the statistical mmWave channel, while the rapid fluctuations of the SINR are due to the adaptation of the small scale fading parameters of **H** (and mainly to the Doppler effect experienced by the moving user). Finally, the red and black lines are referred to the true measured SINR trace $\Gamma$ and its estimate $\bar{\Gamma}$ (after the noise and a first-order filter are applied), respectively. We observe that, for low SINR regimes, $\bar{\Gamma}$ presents a noisy trend but appears still similar to the original trace while, when considering good SINR regimes (e.g., when the UE is in LOS), the estimated trace almost overlaps with its measured original version.

### 3.3.3 Simulation Parameters

The reference scenario (for which one example of random realization is presented in Fig. 3.5) is a typical urban grid having area $200 \times 115$ meters, where 4 non-overlapping buildings of random size and height are deployed, in order to randomize the channel dynamics (in terms on LOS-NLOS transitions) for the moving user. Three mmWave eNBs are located at coordinates $eNB_2 = (0; 50)$, $eNB_3 = (200; 50)$ and $eNB_4 = (100; 110)$, at a height of 10 meters. The LTE $eNB_1$ is co-located with $eNB_4$. We consider a single UE that is at coordinates $(50; -5)$ at the

---

[10]The estimation of the channel is relatively straightforward in 3GPP LTE [119, 181] and is based on the cell reference signal (CRS) that is continuously and omnidirectionally sent from each eNB. However, a CRS will likely not be available in mmWave systems, since downlink transmissions at mmWave frequencies will be directional and specific to the UE [98].

**Figure 3.5:** Random realization of the simulation scenario. The grey rectangles are $4$ randomly deployed non-overlapping buildings.

beginning of the simulation. It then moves along the x-axis at speed $v$ m/s, until it arrives in position $(150; -5)$. The simulation duration $T_{\mathrm{sim}}$ therefore depends on the UE speed $v$ and is given by $T_{\mathrm{sim}} = \frac{l_{\mathrm{path}}}{v} = 20$ s, where $l_{\mathrm{path}} = 100$ m is the length of the path of the UE during the simulation and the default value of the mobile speed has been taken to be $v = 5$ m/s.

Our results are derived through a Monte Carlo approach, where multiple independent simulations are repeated, to get different statistical quantities of interest. In each experiment: (i) we randomly deploy the obstacles; (ii) we apply the measurement framework described in Sec. 3.2.1 to collect one CRT every $D$ seconds; and (iii) we eventually employ one of the HO algorithms presented in Sec. 3.2.3.

The goal of these simulations is to assess the difference in performance between a system using dual connectivity, with fast switching and SCH, and another where hard handover (HH) is used, for different values of $D$, i.e., when varying the periodicity of the CRT generation at the LTE eNB side. Indeed the comparison between these two configurations can be affected by several parameters, which are based on realistic system design considerations and are summarized in Table 3.3 [402]. On the other hand, the performance of the two options does not depend on the interference, since its impact is similar in both schemes. The value of the delay to the MME node ($D_{\mathrm{MME}}$) is chosen in order to model both the propagation delay to a node which is usually centralized and far from the access network, and the processing delays of the MME server. We also model the additional latency $D_{\mathrm{X2}}$ introduced by the X2 connections between each pair of eNBs, which has an impact on (i) the forwarding of PDCP PDUs from the LTE eNB to the mmWave ones; (ii) the exchange of control messages for the measurement reporting framework and (iii) the network procedures which require coordination among eNBs. Thus, the latency $D_{\mathrm{X2}}$ may delay the detection at the LTE eNB coordinator of an outage with respect to the current mmWave link. In order to avoid performance degradation, the value of $D_{\mathrm{X2}}$ should be smaller than 2.5 ms, as recommended by [182].

We consider an SINR threshold $\Gamma_{out} = -5$ dB, assuming that, if $\bar{\Gamma}_j(t) < \Gamma_{out}$, no control signals are collected by eNB$_j$ at time $t$ when the UE is transmitting its SRSs. Reducing $\Gamma_{out}$ allows the user to be potentially found by more suitable mmWave cells, at the cost of designing more complex (and expensive) receiving schemes, able to detect the intended signal in more noisy channels. eNBs are equipped with a Uniform Planar Array (UPA) of $8 \times 8$ elements, which allow them to steer beams in $N_{\mathrm{eNB}} = 16$ directions, whereas UEs have a UPA of $4 \times 4$ antennas, steering beams through $N_{\mathrm{UE}} = 8$ angular directions.

The behavior of the UDP transport protocol (whose interarrival packet generation time is $T_{\mathrm{UDP}}$) is tested, to check whether our proposed dual connectivity framework offers good resilience in mobility scenarios. Only downlink traffic is considered.

61

| Parameter | Value | Description |
|---|---|---|
| mmWave $W_{\text{tot}}$ | 1 GHz | Bandwidth of mmWave eNBs |
| mmWave $f_{\text{c}}$ | 28 GHz | mmWave carrier frequency |
| mmWave $P_{\text{TX}}$ | 30 dBm | mmWave transmission power |
| LTE $W_{\text{tot}}$ | 20 MHz | Bandwidth of the LTE eNB |
| LTE $f_{\text{c}}$ | 2.1 GHz | LTE carrier frequency |
| LTE DL $P_{\text{TX}}$ | 30 dBm | LTE DL transmission power |
| LTE UL $P_{\text{TX}}$ | 25 dBm | LTE UL transmission power |
| NF | 5 dB | Noise figure |
| $\Gamma_{out}$ | $-5$ dB | Minimum SINR threshold |
| eNB antenna | $8 \times 8$ | eNB UPA MIMO array size |
| UE antenna | $4 \times 4$ | UE UPA MIMO array size |
| $N_{\text{eNB}}$ | 16 | eNB scanning directions |
| $N_{\text{UE}}$ | 8 | UE scanning directions |
| $T_{\text{sig}}$ | 10 $\mu s$ | SRS duration |
| $\phi_{\text{ov}}$ | 5% | Overhead |
| $T_{\text{per}}$ | 200 $\mu s$ | Period between SRSs |
| $v$ | 5 m/s | UE speed |
| $B_{\text{RLC}}$ | 10 MB | RLC buffer size |
| $D_{\text{X2}}$ | 1 ms | One-way delay on X2 links |
| $D_{\text{MME}}$ | 10 ms | One-way MME delay |
| $T_{\text{UDP}}$ | $\{20, 80\}\,\mu s$ | UDP packet interarrival time |
| $s_{\text{UDP}}$ | 1024 byte | UDP payload size |
| $D$ | $\{1.6, 12.8, 25.6\}$ ms | CRT intergeneration delay |

**Table 3.3:** Simulation parameters.

## 3.4 Results And Discussion

In this section, we present some results that have been derived for the scenario presented in Sec. 3.3.3. Different configurations have been compared in terms of packet loss, latency, PDCP throughput, RRC and X2 traffic in order to:

i) compare DC with fast switching and SCH versus the traditional standalone hard handover architectures;

ii) compare the performance of the dynamic and the fixed TTT HO algorithms;

iii) validate our proposed measurement reporting system varying the CRT intergeneration periodicity $D$ and the UDP interarrival packet time $T_{\text{UDP}}$.

### 3.4.1 Packet Loss and Handover

In Fig. 3.6(a) we plot the *average number of handover* (or switch) events. As expected, we notice that this number is much higher when considering the DC configuration. The reason is that, since the DC-aided fast switching and SCH procedures are faster than the traditional standalone hard handover, the UE has more chances to change its current cell and adapt to the channel dynamics in a more responsive way. Moreover, when increasing the delay $D$, i.e., when reducing the CRT generation periodicity, the number of handovers reduces, since the UE may have fewer opportunities to update its serving cell, for the same simulation duration. Finally, we see that a dynamic HO procedure requires, on average, a larger number of handover events, to account for the situations in which TTT< 150 ms, when the UE may change its serving cell earlier than it would have done if a fixed TTT algorithm had been applied.

**(a)** Number of handover events during $T_{\text{sim}}$ seconds.

**(b)** UDP packet loss ratio.

**Figure 3.6:** Average number of handover events and packet loss ratio, for different values of the delay $D$, for a fixed and dynamic TTT HO algorithm. Narrow bars refer to a hard handover configuration, while wide colored bars refer to a dual connectivity implementation. The RLC buffer size is $B = 10$ MB and the interarrival packet time is $T_{\text{UDP}} = 20\ \mu s$.

Another element to consider in this performance analysis is the *packet loss ratio* $R_{\text{loss}}$, plotted in Fig. 3.6(b)[11], and defined as the ratio between lost and sent packets, averaged over the $N$ different iterations for each set of parameters. Since the UDP source constantly injects packets into the system, with interarrival time $T_{\text{UDP}}$, it can be computed as $R_{\text{loss}} = 1 - rT_{\text{UDP}}/T_{\text{sim}}$ where $r$ is the total number of received packets and $T_{\text{sim}}$ is the duration of each simulation. We first notice that, with the use of the DC solution, fewer packets are lost. In fact, there are mainly two elements that contribute to the losses: (i) some UDP packets, which are segmented in the RLC retransmission buffer, cannot be reassembled at the PDCP layer and are therefore lost; (ii) during handover, the target eNB RLC transmission buffer receives both the packets sent by the UDP application with interpacket interval $T_{\text{UDP}}$ and the packets that were in the source eNB RLC buffer. If the latter is full, then the target eNB buffer may overflow and discard packets.

Both these phenomena are stressed by the fact that the standalone HH procedure takes more time than both the DC-aided fast switching and SCH procedures. Moreover, during a complete outage event, with the HH solution, until the UE has completed the Non Contention Based Random Access procedure with the LTE eNB, packets cannot be sent to the UE and must be buffered at the RLC layer. This worsens the overflow behavior of the RLC buffer. Instead, with fast switching, the UE does not need to perform random access, since it is already connected and, as soon as packets get to the buffer of the LTE eNB, they are immediately transmitted to the UE.

Fig. 3.6(b) also shows that the packet loss ratio increases when $D$ increases since, if handover or switch events are triggered less frequently, the RLC buffer occupancy increases, and so does the probability of overflow.

Finally, almost no differences are registered when considering a dynamic or fixed TTT HO algorithm, nor when increasing the CRT delay from $D = 12.8$ ms to $D = 25.6$ ms (this aspect will be explained in more detail later).

**(a)** Latency, for $T_{\mathrm{UDP}} = 20\ \mu s$.

**(b)** Latency, for $T_{\mathrm{UDP}} = 80\ \mu s$.

**Figure 3.7:** Average latency, for different values of the delay $D$ and the UDP packet interarrival time $T_{\mathrm{UDP}}$, for a fixed and dynamic TTT HO algorithm. Narrow bars refer to a hard handover configuration, while wide colored bars refer to a dual connectivity implementation. The RLC buffer size is $B_{\mathrm{RLC}} = 10$ MB.

### 3.4.2 Latency

The latency is measured for each packet, from the time it leaves the PDCP layer of the LTE eNB to when it is successfully received at the PDCP layer of the UE. Therefore, it is the latency of only the correctly received packets, and it accounts also for the forwarding latency $D_{\mathrm{X2}}$ on the X2 link. Moreover, this metric captures the queuing time in the RLC buffers, and the additional latency that occurs when a switch or handover happens, before the packet is forwarded to the target eNB or RAT.

Fig. 3.7 shows that the DC framework outperforms the standalone hard handover: in fact, as we pointed out in Sec. 3.4.1, handovers (which dominate the HH configuration) take more time than the fast switching and SCH procedures, and therefore with DC the UE experiences a reduced latency and no service interruptions. This result is even more remarkable when realizing that, from Fig. 3.6, the absolute number of handover (or switch) events is higher when using DC: despite this consideration, the overall latency is still higher for a system where hard handover is implemented[12].

Furthermore, the latency increases as $D$ increases. In fact, when reducing the intergeneration time of the CRT, the UE is attached to a suboptimal mmWave eNB (or to the LTE eNB) for a longer period of time: this increases the buffer occupancy, thus requiring a stronger effort (and longer time) for forwarding many more packets to the new candidate cell, once the handover (or switch) is triggered. Finally, there are no remarkable differences between $D = 12.8$ and $D = 25.6$ ms.

According to Fig. 3.7(b), the latency gap between the HH and DC configurations is much more impressive when considering $T_{\mathrm{UDP}} = 80\ \mu s$. In fact, with this setup, the RLC buffer is empty most of the time and, when a handover (or a switch) is triggered, very few UDP packets need to be forwarded to the destination mmWave or LTE eNB, thus limiting the impact of latency.

We finally recall that, as already introduced in Sec. 3.2.3, the *handover interruption time* (HIT,

---

[11]The presented figure has been obtained when setting $T_{\mathrm{UDP}} = 20\ \mu s$. We have also tested the configuration $T_{\mathrm{UDP}} = 80\ \mu s$, but we saw that, across the different realizations of the simulation, $R_{\mathrm{loss}}$ was zero, due to the fact that the UDP traffic injected in the system was sufficiently well handled by the buffer, with no overflow.

[12]The latency gap is even more remarkable when considering a dynamic TTT HO algorithm. In fact, although the UE experiences, on average, almost 15% more handovers than in the fixed TTT configuration, the overall latency of the two configurations shown in Fig. 3.7 is comparable, due to the fact that with dynamic TTT some SCHs are more timely.

**(a)** PDCP throughput in Mbps, for $T_{\mathrm{UDP}} = 20\ \mu s$.

**(b)** PDCP throughput in Mbps, for $T_{\mathrm{UDP}} = 80\ \mu s$.

**Figure 3.8:** Average PDCP throughput in Mbps, for different values of the delay $D$ and the UDP packet interarrival time $T_{\mathrm{UDP}}$, for fixed and dynamic TTT HO algorithm. Narrow bars refer to a hard handover configuration, while wide colored bars refer to a dual connectivity implementation. The RLC buffer size is $B_{\mathrm{RLC}} = 10$ MB.

i.e., the time in which the user's connectivity is interrupted during the handover operations) takes different values, according to the implemented handover scheme (either DC or HH). When considering a switch to LTE, the HIT is negligible if a DC approach is used, since the UE is already connected to both the LTE and the mmWave RATs. There may be an additional forwarding latency for the switch from mmWave to LTE, which however is already accounted for in Fig. 3.7. On the other hand, when referring to the baseline HH architecture, the UE has to perform a complete handover to switch from one RAT to the other, thus introducing a significant additional delay. When considering the handover between mmWave eNBs, instead, the HIT is comparable for both the DC and the HH schemes. However, in the first case, the procedure does not involve any interaction with the core network and the UE is informed about the new mmWave eNB to handover to and the best angular direction to set through an LTE message (while, when choosing the HH configuration, the handover completion is postponed since the UE has to exhaustively scan again the angular space and perform a complete initial beam search to receive a connection-feedback message from the new serving mmWave eNB). In general, the DC approach is thus preferred in terms of reduced interruption time too.

### 3.4.3   PDCP Throughput

The throughput over time at the PDCP layer is measured by sampling the logs of received PDCP PDUs every $T_s = 5$ ms and summing the received packet sizes to obtain the total number of bytes received $B(t)$. Then the throughput $S(t)$ is computed in bit/s as $S(t) = B(t) \times 8/T_s$. In order to get the mean throughput $S_{\mathrm{PDCP}}$ for a simulation, these samples are averaged over the total simulation time $T_{\mathrm{sim}}$, and finally over the $N$ simulations, to obtain the parameter $\mathbb{E}[S_{\mathrm{PDCP}}]$. Notice that the PDCP throughput (which is mainly a measure of the rate that the radio access network can offer, given a certain application rate), is mostly made up of the transmission of new incoming packets, but it may also account for the retransmissions of already transmitted ones.

In Fig. 3.8, it can be observed that the throughput achievable with the dual connectivity solution is slightly higher than with hard handover. The reason is that, when relying on the LTE eNB for dealing with outage events, the UE experiences a non-zero throughput, in contrast to the hard handover configuration which cannot properly react to a situation where no mmWave eNBs are within reach. Moreover, the difference in throughput increases as the application rate increases, in accordance with the results on packet loss described in the previous section.

**(a)** Variance/Mean ratio, for $T_{\mathrm{UDP}} = 20 \; \mu s$.

**(b)** Variance/Mean ratio, for $T_{\mathrm{UDP}} = 80 \; \mu s$.

**Figure 3.9:** Average ratio $R_{\mathrm{var}}$, for different values of the delay $D$ and the UDP packet interarrival time $T_{\mathrm{UDP}}$, for a fixed and dynamic TTT HO algorithm. Narrow bars refer to a hard handover configuration, while wide colored bars refer to a dual connectivity implementation. The RLC buffer size is $B_{\mathrm{RLC}} = 10$ MB.

As expected, the PDCP throughput decreases as $D$ increases, since the CRT are generated less frequently and the beam pair between the UE and its serving mmWave eNB is monitored less intensively. This means that, when the channel conditions change (e.g., due to the user motion, to a pathloss condition modification or to the small and large scale fading parameters update), the communication quality is not immediately recovered and the throughput is affected by portions of time where suboptimal network settings are chosen.

Moreover, as pointed out in Sec. 3.4.2, we cannot see notable differences between the fixed and dynamic TTT HO procedures and between the $D = 12.8$ and the $D = 25.6$ ms CRT delays. Also a lower UDP rate, according to Fig. 3.8(b), presents comparable PDCP throughput gains with respect to the HH option.

Finally, it is interesting to notice that, when the system implements a DC architecture for handover management, the traditional trade-off between latency and throughput no longer holds. In fact, despite the increased number of handover and switch events shown in Fig. 3.6(a), with respect to the baseline HH configuration, the UE experiences both a reduced latency and an increased PDCP throughput, thus enhancing the overall network quality of service.

### 3.4.4   Variance Ratio

In order to compare the variance of the rate experienced in time by a user, according to the different HO algorithms implemented (DC or HH, for fixed and dynamic TTT), we used the ratio

$$R_{\mathrm{var}} = \frac{\sigma_{S_{\mathrm{PDCP}}}}{\mathbb{E}[S_{\mathrm{PDCP}}]}, \tag{3.7}$$

where $\mathbb{E}[S_{\mathrm{PDCP}}]$ is the mean value of the PDCP throughput measured for each HO configuration and $\sigma_{S_{\mathrm{PDCP}}}$ is its standard deviation, obtained over $N$ repetitions. High values of $R_{\mathrm{var}}$ reflect remarkable channel instability, thus the rate would be affected by local variations and periodic degradations.

Let $R_{\mathrm{var,DC}}$ and $R_{\mathrm{var,HH}}$ be the variance ratios of Equation (3.7) for the fast switching with dual connectivity and hard handover configurations, respectively. From Fig. 3.9, we observe that $R_{\mathrm{var,HH}}$ is higher than $R_{\mathrm{var,DC}}$, for each value of the delay $D$, the HO metric and the UDP packet interarrival time $T_{\mathrm{UDP}}$, making it clear that the LTE eNB employed in a DC configuration can stabilize the rate, which is not subject to significant variations. In fact, in the portion of time

in which the UE would experience zero gain if a hard handover architecture were implemented (due to an outage event), the rate would suffer a noticeable discrepancy with respect to the LOS values, thus increasing the rate variance throughout the simulation. This is not the case for the DC configuration, in which the UE can always be supported by the LTE eNB, even when a blockage event affects the scenario. This result is fundamental for real-time applications, which require a long-term stable throughput to support high data rates and a consistently acceptable Quality of Experience for the users.

Furthermore, it can be seen that $R_{\text{var}}$ increases when the CRT are collected more intensively. In fact, even though reducing $D$ ensures better monitoring of the UE's motion and faster reaction to the channel variations (i.e., LOS/NLOS transitions or periodic modification of the small and large scale fading parameters of $\mathbf{H}$), the user is affected by a higher number of handover and switch events, as depicted in Fig. 3.6(a): in this way, the serving cell will be adapted regularly during the simulation, thereby causing large and periodic variation of the experienced throughput. For the same reason, $R_{\text{var}}$ is higher when applying a dynamic TTT HO algorithm, since the handovers and switches outnumber those of a fixed TTT configuration.

Finally, to compare the DC and the HH architectures, we can consider the ratio $R_{\text{DC/HH}} = R_{\text{var,DC}}/R_{\text{var,HH}}$. It assumes values lower than 1, reflecting the lower variance of a DC configuration, with respect to the baseline HH option. We can therefore affirm that (i) $R_{\text{DC/HH}} < 1$ for every parameter combination and (ii) although the dynamic TTT HO approach shows an absolute higher variance than the fixed TTT one, the hard handover baseline suffers much more because of the aggressiveness of the dynamic TTT configuration than the DC architecture, and therefore $R_{\text{DC/HH,dyn}} < R_{\text{DC/HH,fixed}}$.

### 3.4.5 RRC Traffic

The RRC traffic is an indication of how many control operations are done by the UE-mmWave eNB pairs. Moreover, it is dependent also on the RRC PDU size[13].

Fig. 3.10 shows the RRC traffic for different values of the delay $D$. Notice that the RRC traffic is independent of the buffer size $B$, since even 10 MB are enough to buffer the RRC PDUs, and of the UDP packet interarrival time $T_{\text{UDP}}$. It can be seen that fast switching causes an RRC traffic which is lower than for hard handover. The reason for this behavior is that, when implementing a DC solution, part of the control channel occupancy is due to the switches between the mmWave eNB and the LTE eNB, which use smaller control PDUs than standalone handover events with the HH architecture. A lower RRC traffic is better, since it allows to allocate more resources to data transmission and, given the same amount of control overhead, it allows to scale to a larger number of users [402].

The RRC traffic is then higher for the dynamic TTT HO configuration due to the corresponding higher number of required handovers and switches shown in Fig. 3.6(a).

Finally, we highlight that the RRC traffic measured for a CRT intergeneration periodicity $D = 1.6$ ms is lower than for $D \in \{12.8, 25.6\}$ ms, despite its higher number of required handovers and switches. The reason is that, when the CRT are very frequent, the UE is more intensively monitored, and can thus react more promptly when an outage or a channel update occurs. In this way, retransmissions of control PDUs are less probable and thus fewer messages need to be exchanged at the RRC layer.

---

[13]For example, a switch message contains 1 byte for each of the bearers that should be switched, while an RRC connection reconfiguration message (which triggers the handover) carries several data structures, for a minimum of 59 bytes for a single bearer reconfiguration.

**Figure 3.10:** Average amount of traffic at the RRC layer in bit/s, for different values of the delay $D$, for a fixed and dynamic TTT HO algorithm. Narrow bars refer to a hard handover configuration, while wide colored bars refer to a dual connectivity implementation. The RLC buffer size is $B_{\mathrm{RLC}} = 10$ MB.



**(a)** $\mathbb{E}[S_{\mathrm{X2}}]/\mathbb{E}[S_{\mathrm{PDCP}}]$, for $T_{\mathrm{UDP}} = 20 \ \mu s$.



**(b)** $\mathbb{E}[S_{\mathrm{X2}}]/\mathbb{E}[S_{\mathrm{PDCP}}]$, for $T_{\mathrm{UDP}} = 80 \ \mu s$.

**Figure 3.11:** Average ratio of X2 and PDCP throughput, for different values of the delay $D$ and of the UDP packet interarrival time $T_{\mathrm{UDP}}$, for a fixed and dynamic TTT HO algorithm. Narrow bars refer to a hard handover configuration, while wide colored bars refer to a dual connectivity implementation. The RLC buffer size is $B_{\mathrm{RLC}} = 10$ MB.

### 3.4.6 X2 Traffic

One drawback of the DC architecture is that it needs to forward PDCP PDUs from the LTE eNB to the mmWave eNB, besides forwarding the content of RLC buffers during switching and SCH events. On the other hand, the HH option only needs the second kind of forwarding during handovers. Therefore, the load on the X2 links connecting the different eNBs is lower for the HH solution, as can be seen in Fig. 3.11, which shows the ratio between the average $\mathbb{E}[S_{\mathrm{X2}}]$ of the sum of the throughput $S_{\mathrm{X2}}$ in the six X2 links of the scenario and the average PDCP throughput $\mathbb{E}[S_{\mathrm{PDCP}}]$. It can be seen that for the DC architecture the ratio is close to 1, therefore the X2 links for such configuration must be dimensioned according to the target PDCP throughput for each mmWave eNB. For both architectures the ratio is higher for the lower UDP interarrival time, since there are more packets buffered at the RLC layer that must be forwarded, and also for lower delay $D$, since there are more handover events. However, as we will discuss in more detail in Sec. 3.4.7, the forwarding cost (in terms of inbound traffic to the mmWave eNB) of the DC architecture is similar to that of HH.

**(a)** Evolution of PDCP throughput for HH.



**(b)** Evolution of PDCP throughput for DC.

**Figure 3.12:** Evolution, for a specific simulation of duration $T_{\text{sim}} = 20$ seconds, of the PDCP throughput and of the UE's instantaneous mmWave eNB association. We compare both the hard handover (above) and the dual connectivity (below) configurations, for the fixed TTT HO algorithm and a delay $D = 1.6$ ms. The RLC buffer size is $B_{\text{RLC}} = 10$ MB. The green line represents the current cell over time, where cells from 2 to 4 are mmWave eNBs and cell 1 is the LTE eNB.

### 3.4.7 Final Comments

**Dual Connectivity vs. Hard Handover**: It can be seen that, in general, a multi-connectivity architecture performs better than the hard handover configuration. The main benefit is the short time it takes to change radio access network and its enhancements are shown in terms of mainly: (i) *latency*, which is reduced up to 50% because the fast switching and SCH procedures are in general much faster than traditional handovers (although the number of SCH or switching events may be higher with DC), as observed in Fig. 3.7 and Fig. 3.6(a); (ii) *packet loss*, which is reduced since PDUs are less frequently buffered, thus reducing the overflow probability, as shown in Fig 3.6(b). This is shown by the lower PDCP throughput of Fig. 3.12(a), referred to the HH configuration, with respect to that of the DC architecture of Fig. 3.12(b); (iii) *control signaling* related to the user plane which, despite an increase of the RRC traffic for the LTE eNB, is smaller with the DC solution (this allows the LTE eNB to handle the load of more UEs). This is supported by the results shown in Fig 3.10; (iv) *throughput variance*, where smaller rate variations are registered, with a reduction of $R_{\text{var}}$ of up to 40%, as observed in Fig. 3.9. As an example, Fig. 3.12(a) shows periodic wide fluctuations of the throughput (which sometimes is even zero, when outages occur), while it settles on steady values when DC is applied, as in Fig. 3.12(b).

We also showed that, when the system implements the DC configuration, despite the increased number of handovers and switches, the UE can *jointly* achieve both a reduced latency and an increased PDCP throughput, enhancing its overall quality of service. We have also examined the main cost of the DC architecture, showing in Sec. 3.4.6 that the X2 traffic for the DC option is higher than for the HH configuration because of the forwarding of packets from the LTE eNB to the mmWave ones. However, we must recall that, with the HH solution, the mmWave eNBs receive the packets from the core network through the S1 link, which is not used for the mmWave eNBs in the DC configuration. Therefore, when considering the overall inbound traffic to the mmWave eNBs on both the X2 and the S1 links, the costs of the two architectures may be

**(a)** Corner case. The grey rectangles are buildings.

**(b)** Latency.

**Figure 3.13:** Average latency, for different values of the delay $D$ and for $T_{\text{UDP}} = 20\ \mu s$, comparing a fixed and dynamic TTT HO algorithm. The colored bars refer to a dual connectivity implementation for HO management. The RLC buffer size is $B = 10$ MB and a *corner scenario* is implemented, for a user moving at speed $v$.

equivalent. Given these considerations, we argue that the use of multi-connectivity for mobility management is to be preferred to the traditional hard handover approach.

**UDP interarrival time**: We observed that the general behaviors are similar for most metrics. However, the latency is much lower when $T_{\text{UDP}} = 80\ \mu s$, since RLC buffers are empty most of the time and fewer packets need to be forwarded during the switching and handover events. This justifies the wider gap between DC and HH architectures, with respect to the $T_{\text{UDP}} = 20\ \mu s$ case.

**CRT intergeneration delay and beamforming architecture**: We noticed remarkable differences between $D = 1.6$ and $D = 25.6$ ms (validating the choice of designing a *digital* BF architecture, more complex but more efficient in terms of both latency and throughput) but almost no distinction between the $D = 12.8$ and $D = 25.6$ ms configurations: we conclude that a *hybrid* BF system at the mmWave eNB side is not to be preferred to an *analog* one, since the complexity is increased while the overall performance is almost equivalent.

**Fixed vs. Dynamic TTT**: We showed that the second approach never results in a performance degradation for any of the analyzed metrics. Moreover, we showed that it may also deliver tangible improvements in some specific scenarios where the traditional methods fail, such as the one shown in Fig. 3.13. In this *corner scenario*, the UE turns left at a T-junction and loses LOS with respect to both mmWave eNBs at the bottom. However, the mmWave eNB on top of the scenario is now in LOS, thus the handover should be triggered as quickly as possible. From the result in Fig. 3.13(b), we observe that in this case a dynamic and more aggressive approach is able to massively reduce latency compared to the fixed configuration, since a reduced TTT may be vital in this specific scheme, in which the user experiences a degraded rate until the handover to the LOS mmWave eNB is completed. We indeed state that, since the dynamic TTT algorithm never underperforms the fixed TTT approach but is able to greatly improve the performance in specific scenarios, it should be preferred for handover management.

## 3.5 Conclusion And Future Work

A limitation for the deployment of mmWave 5G systems is the rapidly changing dynamic channel caused by user mobility. The UE may be suddenly in outage with respect to all the mmWave eNBs, and a classic standalone architecture with traditional handovers cannot react quickly

enough. In this study we proposed a dual connectivity framework that, with the aid of a macro LTE eNB, can collect measurements and track the channel dynamics and perform fast switching to fall back to LTE and SCH for a fast handover among the mmWave eNBs. We showed, with an extensive simulation campaign, that the proposed framework is able to improve the performance of an end-to-end network with mmWave access links with respect to several metrics, including latency, throughput (in terms of both average and stability), radio control signaling and packet loss. Moreover, we presented and studied the performance of a dynamic TTT algorithm for SCH, showing that in some specific cases it may gain significantly with respect to a standard fixed TTT handover algorithm.

# 4

# Beam Management in 5G mmWave Networks

## 4.1 Introduction

As mentioned in Chapter 1, the mmWave spectrum is considered as an enabler of the 5G performance requirements in micro and picocellular networks [39, 42]. These frequencies offer much more bandwidth than current cellular systems in the congested bands below 6 GHz, and initial capacity estimates have suggested that networks operating at mmWaves can offer orders of magnitude higher bit-rates than 4G systems [33]. Nonetheless, the higher carrier frequency makes the propagation conditions harsher than at the lower frequencies traditionally used for wireless services, especially in terms of robustness [32]. Signals propagating in the mmWave band suffer from increased pathloss and severe channel intermittency, and are blocked by many common materials such as brick or mortar [40], and even the changing position of the body relative to the mobile device can lead to rapid drops in signal strength.

To deal with these impairments, next-generation cellular networks such as 3GPP NR [7] must provide a set of mechanisms by which UEs and mmWave gNB stations establish highly directional transmission links, typically using high-dimensional phased arrays, to benefit from the resulting beamforming gain and sustain an acceptable communication quality. Directional links, however, require fine alignment of the transmitter and receiver beams, achieved through a set of operations known as *beam management*. They are fundamental to perform a variety of control tasks including (i) Initial Access (IA) [175, 183, 416] for idle users, which allows a mobile UE to establish a physical link connection with a gNB, and (ii) beam tracking, for connected users, which enable beam adaptation schemes, or handover, path selection and radio link failure recovery procedures [184, 388]. In current LTE systems, these control procedures are performed using omnidirectional signals, and beamforming or other directional transmissions can only be performed after a physical link is established, for data plane transmissions. On the other hand, in the mmWave bands, it may be essential to exploit the antenna gains even during initial access and, in general, for control operations. Omnidirectional control signaling at such high frequencies, indeed, may generate a mismatch between the relatively short range at which a cell can be detected or the control signals can be received (control-plane range), and the much longer range at which a user could send and receive data when using beamforming (data-plane range). However, directionality can significantly delay the access procedures and make the performance more sensitive to the beam alignment. These are particularly important issues in 5G networks, in particular when considering high-mobility environments and blockage, and motivate the need

to extend current LTE control procedures with innovative mmWave-aware beam management algorithms and methods.

### 4.1.1 Contributions

This chapter, which was published also in [394, 395, 416, 426][1], is a tutorial on the design and dimensioning of beam management frameworks for mmWave cellular networks. In particular, we consider the parameters of interest for 3GPP NR networks, which will support carrier frequencies up to 52.6 GHz [114]. We also report an analysis of beam management techniques, including initial access and tracking strategies, for cellular networks operating at mmWaves under realistic NR settings and channel configurations, and describe how to optimally design fast, accurate and robust control-plane management schemes through measurement reports in different scenarios. Finally, we introduce a context-based beam management scheme for UAV networks. More specifically, in this chapter we:

- Provide an overview of the most effective measurement collection frameworks for 5G systems operating at mmWaves. We focus on DL and UL frameworks, according to whether the reference signals are sent from the gNBs to the UEs or vice versa, respectively, and on NSA and SA architectures, according to whether the control plane is managed with the support of an LTE overlay or not, respectively. A DL configuration is in line with the 3GPP specifications for NR and reduces the energy consumption at the UE side, but it may be lead to a worse beam management performance than in the UL. Moreover, when considering stable and dense scenarios which are marginally affected by the variability of the mmWave channel, an SA architecture is preferable for the design of fast IA procedures, while an NSA scheme may be preferable for reducing the impact of the overhead on the system performance and enable more robust and stable communication capabilities.

- Simulate the performance of the presented measurement frameworks in terms of *signal detection accuracy*, using a realistic mmWave channel model based on real-world measurements conducted in a dense, urban scenario in which environmental obstructions (i.e., urban buildings) can occlude the path between the transmitter and the receiver. The tutorial shows that accurate beam management operations can be guaranteed when configuring narrow beams for the transmissions, small subcarrier spacings, denser network deployments and by adopting *frequency diversity* schemes.

- Analyze the *reactiveness* (i.e., how quickly a mobile user gets access to the network and how quickly the framework is able to detect an updated channel condition), and the *overhead* (i.e., how many time and frequency resources should be allocated for the measurement operations). In general, fast initial access and tracking schemes are ensured by allocating a large number of time/frequency resources to the users in the system, at the expense of an increased overhead, and by using advanced beamforming capabilities (e.g., digital or hybrid beamforming), which allow the transceiver to sweep multiple directions at any given time.

- Illustrate some of the complex and interesting trade-offs to be considered when designing solutions for next-generation cellular networks by examining a wide set of parameters based on 3GPP NR considerations and agreements (e.g., the frame structure and other relevant physical-layer aspects).

---

[1]Part of this chapter is based on joint work with Marco Giordani.

- Experimentally evaluate the performance of beam management schemes in the context of UAV networks at mmWaves. In particular, we review the benefits of side-information-aided beam management and present a GPS-aided beam tracking algorithm for UAV-based aerial cells. We prototype the proposed algorithm on a mmWave aerial link using a DJI M600 Pro and 60 GHz radios and prove its effectiveness in reducing the average link establishment latency by 66% with respect to state-of-the-art non-aided schemes.

In general, the results prove that the optimal design choices for implementing efficient and fast initial access and reactive tracking of the mobile user strictly depend on the specific environment in which the users are deployed, and must account for several specific features such as the base stations density, the antenna geometry, the beamforming configuration and the level of integration and harmonization of different technologies.

### 4.1.2 Organization

The sections of this chapter are organized as follows. Sec. 4.2 reports the related work on beam management at mmWave frequencies. Sec. 4.3 provide basic information on the 3GPP Release 15 frame structure for NR, and presents the candidate DL and UL measurement signals that can be collected by the NR nodes for the beam management operations. Sec. 4.4 describes the beam management frameworks whose performance will be analyzed, simulated and compared in the remainder of the chapter. Sec. 4.5 defines the parameters that affect the performance of beam management in NR. Sec. 4.6 reports a performance evaluation and some considerations on the trade-offs and on which are the best configurations for beam management frameworks. Additional considerations and comprehensive remarks, aiming at providing guidelines for selecting the optimal IA and tracking configuration settings as a function of the system parameters, are stated in Sec. 4.7. Finally, Sec. 4.8 presents the experimental evaluation of beam management schemes for UAVs, using 60 GHz radios on a drone. Finally, Sec. 4.9 concludes the chapter.

## 4.2 Related Work

Measurement reporting is quite straightforward in LTE [181]: the DL channel quality is estimated from an omnidirectional signal called the Cell Reference Signal (CRS), which is regularly monitored by each UE in connected state to create a wideband channel estimate that can be used both for demodulating downlink transmissions and for estimating the channel quality [98]. However, when considering mmWave networks, in addition to the rapid variations of the channel, CRS-based estimation is challenging due to the directional nature of the communication, thus requiring the network and the UE to constantly monitor the direction of transmission of each potential link. Tracking changing directions can decrease the rate at which the network can adapt, and can be a major obstacle in providing robust and ubiquitous service in the face of variable link quality. In addition, the UE and the gNB may only be able to listen to one direction at a time, thus making it hard to receive the control signaling necessary to switch paths.

To overcome these limitations, several approaches in the literature, as summarized in Table 4.1, have proposed directional-based schemes to enable efficient control procedures for both the idle and the connected mobile terminals, as surveyed in the following paragraphs.

Papers on IA[2] and tracking in 5G mmWave cellular systems are very recent. Most literature refers to challenges that have been analyzed in the past at lower frequencies in ad hoc wireless

---

[2]We refer to works [175, 183, 187] for a detailed taxonomy of recent IA strategies.

**Table 4.1:** Relevant literature on measurement reporting, initial access and beam management strategies for mmWave networks.

| Topic | Relevant References |
|---|---|
| IEEE 802.11ad [185] | [102, 176, 186]. Not suitable for long-range, dynamic and outdoor scenarios. |
| Initial Access [175, 183, 187] | [164, 165, 188] Exhaustive search. [189–191] More advanced searching schemes. [192–195] Context-aware initial access. [196–198] Performance comparison. |
| Beam Management [388] | [157, 199, 200] Mobility-aware strategies. [62, 63, 155, 163, 201, 402] Multi-connectivity solutions. |

network scenarios or, more recently, referred to the 60 GHz IEEE 802.11ad WLAN and WPAN scenarios (e.g., [102, 185, 186]). However, most of the proposed solutions are unsuitable for next-generation cellular network requirements and present many limitations (e.g., they are appropriate for short-range, static and indoor scenarios, which do not match well the requirements of 5G systems). Therefore, new specifically designed solutions for cellular networks need to be found.

In [164, 188], the authors propose an exhaustive method that performs directional communication over mmWave frequencies by periodically transmitting synchronization signals to scan the angular space. The result of this approach is that the growth of the number of antenna elements at either the transmitter or the receiver provides a large performance gain compared to the case of an omnidirectional antenna. However, this solution leads to a long duration of the IA with respect to LTE, and poorly reactive tracking. Similarly, in [165], measurement reporting design options are compared, considering different scanning and signaling procedures, to evaluate access delay and system overhead. The channel structure and multiple access issues are also considered. The analysis demonstrates significant benefits of low-resolution fully digital architectures in comparison to single stream analog beamforming. Additionally, more sophisticated discovery techniques (e.g., [189, 190]) alleviate the exhaustive search delay through the implementation of a multi-phase hierarchical procedure based on the access signals being initially sent in few directions over wide beams, which are iteratively refined until the communication is sufficiently directional. In [191] a low-complexity beam selection method by low-cost analog beamforming is derived by exploiting a certain sparsity of mmWave channels. It is shown that beam selection can be carried out without explicit channel estimation, using the notion of compressive sensing.

The issue of designing efficient beam management solutions for mmWave networks is addressed in [199], in which the author designs a mobility-aware user association strategy to overcome the limitations of the conventional power-based association schemes in a mobile 5G scenario. Other relevant papers on this topic include [157], in which the authors propose smart beam tracking strategies for fast mmWave link establishment and maintenance under node mobility. In [200], the authors proposed the use of an extended Kalman filter to enable a static base station, equipped with a digital beamformer, to effectively track a mobile node equipped with an analog beamformer after initial channel acquisition, with the goal of reducing the alignment error and guarantee a more durable connectivity. Recently, robust IA and tracking schemes have been designed by leveraging out-of-band information to estimate the mmWave channel. In [62, 163, 388, 402] an approach where 5G cells operating at mmWaves (offering much higher rates) and traditional 4G cells below 6 GHz (providing much more robust operation) are employed in parallel have been proved to enable fast and resilient tracking operations. In [63], a framework which integrates both LTE and 5G interfaces is proposed as a solution for mobility-related link failures and throughput degradation of cell-edge users, relying on coordinated transmissions from cooperating cells are coordinated for both data and control signals. In [155], a novel approach for

analyzing and managing mobility in joint sub-6GHz–mmWave networks is proposed by leveraging on device caching along with the capabilities of dual-mode base stations to minimize handover failures, reduce inter-frequency measurement, reduce energy consumption, and provide seamless mobility in emerging dense heterogeneous networks. Moreover, the authors in [201] illustrate how to exploit spatial congruence between signals in different frequency bands and extract mmWave channel parameters from side information obtained in another band. Despite some advantages, the use of out-of-band information for the 5G control plane management poses new challenges that remain unsolved and which deserve further investigation.

Context information can also be exploited to improve the cell discovery procedure and minimize the delay [192, 193], while capturing the effects of position inaccuracy in the presence of obstacles. In the scheme proposed in [194], booster cells (operating at mmWave) are deployed under the coverage of an anchor cell (operating at LTE frequencies). The anchor base station gets control over IA informing the booster cell about user locations, in order to enable mmWave gNB to directly steer towards the user position. Finally, in [195], the authors studied how the performance of analog beamforming degrades in the presence of angular errors in the available Context Information during the initial access or tracking procedures, according to the status of the UE (connected or non-connected, respectively). With respect to the study we discuss in Sec. 4.8, however, the context-based solutions that can be found in the state of the art are limited to traditional cellular network scenarios (i.e., fixed base station and users with low mobility), and have not been implemented in an experimental prototype.

The performance of the association techniques also depends on the beamforming architecture implemented in the transceivers. Preliminary works aiming at finding the optimal beamforming strategy refer to WLAN scenarios. For example, the algorithm proposed in [176] takes into account the spatial distribution of nodes to allocate the beamwidth of each antenna pattern in an adaptive fashion and satisfy the required link budget criterion. Since the proposed algorithm minimizes the collisions, it also minimizes the average time required to transmit a data packet from the source to the destination through a specific direction. In 5G scenarios, papers [164, 188, 189] give some insights on trade-offs among different beamforming architectures in terms of users' communication quality. In this context, articles [196, 197] evaluate the mmWave cellular network performance while accounting for the beam training, association overhead and beamforming architecture. More recently the authors in [198], based on current 5G NR slot design considerations, compare the performance of several IA schemes in terms of coverage and search delays, and for different antenna array settings. The results show that, although employing wide beams, initial beam training with full pilot reuse is nearly as good as perfect beam alignment. However, they lack considerations on the latest 3GPP specifications for NR. Finally, paper [202] provides an overview of the main features of NR with respect to initial access and multi-beam operations, and article [203] reports the details on the collection of channel state information in NR. However, the aforementioned papers only present a high level overview, and do not include a comprehensive performance evaluation of NR beam management frameworks at mmWave frequencies.

The above discussion makes it apparent how next-generation mmWave cellular networks should support a mechanism by which the users and the infrastructure can quickly determine the best directions to establish the mmWave links, an operation which may increase the latency and the overhead of the communication and have a substantial impact on the overall network performance. In the remainder of this chapter we will provide guidelines to characterize the optimal beam management strategies as a function of a variety of realistic system parameters.

**Table 4.2:** Reference signals for beam management operations, for users in idle and connected states, in downlink or uplink.

|  | Initial Access (Idle UE) | Tracking (Connected UE) |
|---|---|---|
| Downlink | SS blocks (carrying the PSS, the SSS, and the PBCH). See references [18, 114, 204–210]. | CSI-RSs and SS blocks. See references [18, 114, 204–206, 211–216]. |
| Uplink | 3GPP does not use uplink signals for initial access, but the usage of SRSs has been proposed in [62, 163, 388] | SRSs. See references [18, 114, 204, 217, 218]. |

## 4.3 Frame Structure and Signals for 3GPP NR at mmWave Frequencies

Given that NR will support communication at mmWave frequencies, it is necessary to account for beamforming and directionality in the design of its PHY and MAC layers. The NR specifications will thus include a set of parameters for the frame structure dedicated to high carrier frequencies, as well as synchronization and reference signals that enable beam management procedures [114]. In this regard, in Sec. 4.3.1 and Sec. 4.3.2 we introduce the 3GPP frame structure and measurement signals proposed for NR, respectively, which will provide the necessary background for the remainder of this tutorial.

### 4.3.1 NR Frame Structure

The 3GPP technical specification in [18] and the report in [114] provide the specifications for the PHY layer. Both Frequency Division Duplexing (FDD) and TDD will be supported.

The *waveform* is OFDM with a cyclic prefix. Different numerologies[3] will be used, in order to address the different use cases of 5G [5]. The frame structure follows a time and frequency grid similar to that of LTE, with a higher number of configurable parameters. The subcarrier spacing is $15 \times 2^n$ kHz, $n \in \mathbb{Z}, n \leq 4$. In Release 15, there will be at most 3300 subcarriers, for a maximum bandwidth of 400 MHz. A *frame* lasts 10 ms, with 10 subframes of 1 ms. It will be possible to multiplex different numerologies for a given carrier frequency, and the whole communication must be aligned on a subframe basis. A *slot* is composed of 14 OFDM symbols. There are multiple slots in a subframe, and their number is given by the numerology used, since the symbol duration is inversely proportional to the subcarrier spacing [219]. *Mini-slots* are also supported: they can be as small as 2 OFDM symbol and have variable length, and can be positioned asynchronously with respect to the beginning of the slot (so that low-latency data can be sent without waiting for the whole slot duration).

### 4.3.2 NR Measurements for Beam Management

Regular beam management operations are based on the control messages which are periodically exchanged between the transmitter and the receiver nodes. In the following paragraphs we will review the most relevant DL and UL measurement signals supported by 3GPP NR for beam management purposes, as summarized in Table 4.2.

**Downlink Measurements: SS Blocks.**

In the most recent versions of the 3GPP specifications [18], the concept of SS block and burst emerged for periodic synchronization signal transmission from the gNBs. An SS block is a group of 4 OFDM symbols [18, Sec. 7.4.3] in time and 240 subcarriers in frequency (i.e., 20 resource

---

[3]The term numerology refers to a set of parameters for the waveform, such as subcarrier spacing and cyclic prefix duration for OFDM [17].

blocks), as shown in Fig. 4.1. It carries the PSS, the SSS and the PBCH. The DeModulation Reference Signal (DMRS) associated with the PBCH can be used to estimate the Reference Signal Received Power (RSRP) of the SS block. In a slot of 14 symbols, there are two possible locations for SS blocks: symbols 2-5 and 8-11.

The SS blocks are grouped into the first 5 ms of an SS burst [18, 208], which can have different periodicities $T_{\mathrm{SS}} \in \{5, 10, 20, 40, 80, 160\}$ ms [204]. When accessing the network for the first time, the UE should assume a periodicity $T_{\mathrm{SS}} = 20$ ms [205].

The maximum number $L$ of SS blocks in a burst is frequency-dependent [205, 208], and above 6 GHz there could be up to 64 blocks per burst. When considering frequencies for which beam operations are required [220], each SS block can be mapped to a certain angular direction. To reduce the impact of SS transmissions, SS can be sent through wide beams, while data transmission for the active UE is usually performed through narrow beams, to increase the gain produced by beamforming [210].

**Downlink Measurements: CSI-RS.**

It has been agreed that CSI-RSs can be used for Radio Resource Management (RRM) measurements for mobility management purposes in connected mode [7]. As in LTE, it shall be possible to configure multiple CSI-RS to the same SS burst, in such a way that the UE can first obtain synchronization with a given cell using the SS bursts, and then use that as a reference to search for CSI-RS resources [204, 211]. Therefore, the CSI-RS measurement window configuration should contain at least the periodicity and time/frequency offsets relative to the associated SS burst. Fig. 4.2 shows the two options we consider for the time offset of the CSI-RS transmissions. The first option, shown in Fig. 4.2a, allows the transmission of the first CSI-RS $T_{\mathrm{CSI}}$ ms after the end of an SS burst. The second one, shown in Fig. 4.2b, has an additional parameter, i.e., an offset in time $O_{\mathrm{CSI}}$, which represents the time interval between the end of the SS burst and the first CSI-RS. The CSI-RSs, which may not necessarily be broadcast through all the available frequency resources [212, 221], may span $N =1$, 2 or 4 OFDM symbols [18, 222]. For periodic CSI-RS transmissions, the supported periodicities are $T_{\mathrm{CSI,slot}} \in \{5, 10, 20, 40, 80, 160, 320, 640\}$ slots [18], thus the actual periodicity in time depends on the slot duration.

As we assessed in the previous sections of this work, when considering directional communications, the best directions for the beams of the transceiver need to be periodically identified (e.g., through beam search operations), in order to maintain the alignment between the communicating nodes. For this purpose, SS- and CSI-based measurement results can be jointly used to reflect the different coverage which can be achieved through different beamforming architectures [206, 214]. As far as CSI signals are concerned, the communication quality can be derived by averaging the signal quality from the $N_{\mathrm{CSI,RX}}$ best beams among all the available ones, where



**Figure 4.1:** SS block structure [7].

**(a)** Option 1: the first CSI-RS is sent $T_{\text{CSI}}$ ms after an SS burst. **(b)** Option 2: the first CSI-RC is sent $O_{\text{CSI}}$ ms after an SS burst.

**Figure 4.2:** Examples of CSI-RS measurement window and periodicity configurations. SS blocks are sent every $T_{\text{SS}}$ ms, and they embed time and frequency offsets indicating the time and frequency allocation of CSI-RS signals within the frame structure.

the value of $N_{\text{CSI,RX}}$ can be configured to 1 or more than 1 [204, 211][4]. Nevertheless, to avoid the high overhead associated with wide spatial domain coverage with a huge number of very narrow beams, on which CSI-RSs are transmitted, it is reasonable to consider transmitting only subsets of those beams, based on the locations of the active UEs. This is also important for UE power consumption considerations [204, 216]. For example, the measurement results based on SS blocks (and referred to a subset of transmitting directions) can be used to narrow down the CSI-RS resource sets based on which a UE performs measurements for beam management, thereby increasing the energy efficiency.

**Uplink Measurements: SRS** The SRSs are used to monitor the uplink channel quality, and are transmitted by the UE and received by the gNBs. According to [204, 217], their transmission is scheduled by the gNB to which the UE is attached, which also signals to the UE the resource and direction to use for the transmission of the SRS. The UE may be configured with multiple SRSs for beam management. Each resource may be periodic (i.e., configured at the slot level), semi-persistent (also at the slot level, but it can be activated or deactivated with messages from the gNB) and a-periodic (the SRS transmission is triggered by the gNB) [204, 218]. The SRSs can span 1 to 4 OFDM symbols, and a portion of the entire bandwidth available at the UE [18, 217].

## 4.4 Beam Management Frameworks for 5G Cellular Systems

In this section, we present three measurement frameworks for both initial access and tracking purposes, whose performance will be investigated and compared in Sec. 4.6.

As we introduced in the above sections of this chapter, the NR specifications include a set of basic beam-related procedures [114] for the control of multiple beams at frequencies above 6 GHz and the related terminologies, which are based on the reference signals described in Sec. 4.3. The different operations are categorized under the term *beam management*, which is composed of four different operations:

- *Beam sweeping*, i.e., covering a spatial area with a set of beams transmitted and received according to pre-specified intervals and directions.

---

[4]In [204] it is specified that, for the derivation of the quality of a cell, the UEs should consider an absolute threshold, and average the beams with quality above the threshold, up to $N_{\text{CSI,RX}}$ beams. If there are no beams above threshold, then the best one (regardless of its absolute quality) should be selected for the cell quality derivation.

- *Beam measurement*, i.e., the evaluation of the quality of the received signal at the gNB or at the UE. Different metrics could be used [206]. We consider here the Signal to Noise Ratio (SNR), which is the average of the received power on synchronization signals divided by the noise power.

- *Beam determination*, i.e., the selection of the suitable beam or beams either at the gNB or at the UE, according to the measurements obtained with the beam measurement procedure.

- *Beam reporting*, i.e., the procedure used by the UE to send beam quality and beam decision information to the RAN.

These procedures are periodically repeated to update the optimal transmitter and receiver beam pair over time.

We consider a *NSA* or a *standalone (SA)* architecture. Non-standalone is a deployment configuration in which a NR gNB uses an LTE cell as support for the control plane management [223] and mobile terminals exploit *multi-connectivity* to maintain multiple possible connections (e.g., 4G and 5G overlays) to different cells so that drops in one link can be overcome by switching data paths [29, 62, 63, 163, 388, 402]. Mobiles in a NSA deployment can benefit from both the high bit-rates that can be provided by the mmWave links and the more robust, but lower- rate, legacy channels, thereby opening up new ways of solving capacity issues, as well as new ways of providing good mobile network performance and robustness. Conversely, with the standalone option, there is no LTE control plane, therefore the integration between LTE and NR is not supported [394].

The measurement frameworks can be also based on a *downlink* or an *uplink* beam management architecture. In the first case, the gNBs transmit synchronization and reference signals (i.e., SS blocks and CSI-RSs) which are collected by the surrounding UEs, while in the second case the measurements are based on SRSs forwarded by the mobile terminal instead. Notice that the increasing heterogeneity in cellular networks is dramatically changing our traditional notion of a communication cell [6], making the role of the uplink important [224] and calling for the design of innovative UL-driven solutions for both the data and the control planes.

In the following, we will describe in detail the three considered measurement schemes[5]. Table 4.3 provides a summary of the main features of each framework.

### 4.4.1   Standalone-Downlink (SA-DL) Scheme

The SA-DL configuration scheme is shown in Fig. 4.3. No support from the LTE overlay is provided in this configuration. The beam management procedure is composed of the following phases:

(i) *Beam sweeping.* The measurement process is carried out with an exhaustive search, i.e., both users and base stations have a predefined codebook of directions (each identified by a beamforming vector) that cover the whole angular space and are used sequentially to transmit/receive synchronization and reference signals [188].

(ii) *Beam measurements.* The mmWave-based measurements for IA are based on the SS blocks. The tracking is done using both the measurements collected with the SS bursts and the CSI-RSs. These last elements cover a set of directions which may or may not cover the entire set of available directions according to the users' needs, as explained in Sec. 4.3. No support from the LTE overlay is provided in this configuration.

---

[5]Notice that we do not consider the SA-UL configuration for both IA and tracking applications. In fact, we believe that uplink-based architectures will likely necessitate the support of the LTE overlay for the management of the control plane and the implementation of efficient measurement operations.

**Table 4.3:** Comparison of the beam management frameworks.

|  | SA-DL | NSA-DL | NSA-UL |
|---|---|---|---|
| Multi-RAT connectivity | Not available | LTE overlay available for robust control operations and quick data fallback [62, 63, 402]. | |
| Reference signal transmission | Downlink | Downlink | Uplink |
| Network coordination | Not available | Possibility of using a centralized controller [388]. | |

| Beam management phase | SA-DL | NSA-DL | NSA-UL |
|---|---|---|---|
| Beam sweep | Exhaustive search based on SS blocks [188]. | | Based on SRS [163]. |
| Beam measurement | UE-side | UE-side | gNB-side |
| Beam determination | The UE selects the optimal communication direction. | | Each gNB sends information on the received beams to a central controller, which selects the best beam pair [62]. |
| Beam reporting | Exhaustive search at the gNB side [225]. | The UE signals the best beam pair using LTE, a RACH opportunity in that direction is then scheduled. | The gNB signals the best beam pair using LTE, a RACH opportunity in that direction is then scheduled. |



**Figure 4.3:** Signals and messages exchanged during the SA-DL beam management procedure (with the beam reporting step of the IA). Notice that the duration of the three phases is not in scale, since it depends on the actual configuration of the network parameters.

(iii) *Beam determination.* The mobile terminal selects the beam through which it experienced the maximum SNR, if above a predefined threshold. The corresponding sector will be chosen for the subsequent transmissions and receptions and benefit from the resulting antenna gain.

(iv) *Beam reporting.* For IA, as proposed by 3GPP, after beam determination the mobile terminal has to wait for the gNB to schedule the RACH opportunity towards the best direction that the UE just determined, for performing random access and implicitly informing the selected serving infrastructure of the optimal direction (or set of directions) through which it has to steer its beam, in order to be properly aligned. It has been agreed that for each SS block the gNB will specify one or more RACH opportunities with a certain time and frequency offset and direction, so that the UE knows when to transmit the RACH preamble [7, 225]. This may require an additional complete directional scan of the gNB, thus further increasing the time it takes to access the network. For the tracking in connected mode, the UE can provide feedback using the mmWave control channel it has already established, unless there is a link failure and no directions can be recovered using CSI-RS. In this case the UE must repeat the IA procedure or try to recover the link using the SS bursts while the user experiences a service unavailability.

### 4.4.2 Non-Standalone-Downlink (NSA-DL) Scheme

The sub-6-GHz overlay can be used with different levels of integration. As shown in Fig. 4.4, the first three procedures are as in the SA-DL scheme. However, non-standalone enables an improvement in the beam reporting phase. Thanks to the control-plane integration with the overlay, the LTE connection can be used to report the optimal set of directions to the gNBs, so that the UE does not need to wait for an additional beam sweep from the gNB to perform the beam reporting or the IA procedures. Thanks to this signaling, a random access opportunity can therefore be immediately scheduled for that direction with the full beamforming gain. Moreover, the LTE link can be also used to immediately report a link failure, and allow a quick data-plane fallback to the sub-6-GHz connection, while the UE recovers the mmWave link.

### 4.4.3 Non-Standalone-Uplink (NSA-UL) Scheme

Unlike in traditional LTE schemes, this framework (first proposed in [62] and then used in [388]) is based on the channel quality of the UL rather than that of the DL signals and, with the joint support of a central coordinator (i.e., an LTE eNB operating at sub-6 GHz frequencies), it enables efficient measurement operations. In this framework, a user searches for synchronization signals from conventional 4G cells. This detection is fast since it can be performed omnidirectionally and there is no need for directional scanning. Under the assumption that the 5G mmWave eNBs are roughly time synchronized to the 4G cell, and the round trip propagation times are not large, an uplink transmission from the UE will be roughly time aligned at any closeby mmWave cell[6] [163]. The NSA-UL procedure[7] is shown Fig. 4.5 with a detailed breakout of the messages exchanged by the different parties. In detail, it is composed of:

(i-ii) *Beam sweeping and beam measurements.* Each UE directionally broadcasts SRSs in the mmWave bands in time-varying directions that continuously sweep the angular space. Each

---

[6]For example, if the cell radius is 150 m (a typical mmWave cell), the round trip delay is only 1 $\mu$s.

[7]Unlike the conventional DL-based measurement configuration, the uplink scheme has not been considered by 3GPP. Nevertheless, we will freely adapt the same NR frame structure proposed for the downlink case to the NSA-UL scheme, using for the uplink SRSs the resources that would be allocated to SS blocks in a downlink framework.

**Figure 4.4:** Signals and messages exchanged during the NSA-DL beam management procedure (with the beam reporting step of the IA). Notice that the duration of the three phases is not in scale, since it depends on the actual configuration of the network parameters.



**Figure 4.5:** Signals and messages exchanged during the NSA-UL beam management procedure (with the beam reporting step of the IA). Notice that the duration of the three phases is not in scale, since it depends on the actual configuration of the network parameters.

potential serving gNB scans all its angular directions as well, monitoring the strength of the received SRSs and building a *report table* based on the channel quality of each receiving direction, to capture the dynamics of the channel.

(iii) *Beam determination.* Once the report table of each mmWave gNB has been filled for each UE, each mmWave cell sends this information to the LTE eNB which, due to the knowledge gathered on the signal quality in each angular direction for each gNB-UE pair, obtains complete directional knowledge over the cell it controls. Hence, it is able to match the beams of the transmitters and the receivers to provide maximum performance.

(iv) *Beam reporting.* The coordinator reports to the UE, on a legacy LTE connection, which

| Parameter | $\Delta_f$ | $D$ | $N_{\mathrm{SS}}$ | $T_{\mathrm{SS}}$ | CSI | $N_{\mathrm{CSI,RX}}$ | $K_{\mathrm{BF}}$ | $M, N_\theta, N_\phi$ | $N_{\mathrm{user}}$ | $\lambda_b$ |
|---|---|---|---|---|---|---|---|---|---|---|
| Accuracy | ✓ | ✓ | x | x | ✓* | x | x | ✓ | x | ✓ |
| Reactiveness | ✓ | x | ✓ | ✓ | ✓* | ✓* | ✓ | ✓ | ✓ | x |
| Overhead | ✓ | ✓ | ✓ | ✓ | ✓ | x | ✓ | ✓ | ✓ | x |

gNB yields the best performance, together with the optimal direction in which the UE should steer its beam, to reach the candidate serving cell in the optimal way. The choice of using the LTE control link during the tracking is motivated by the fact that the UE may not be able to receive from the optimal mmWave link if not properly aligned, thereby removing a possible point of failure in the control signaling path. Moreover, since path switches and cell additions in the mmWave regime are common due to link failures, the control link to the serving mmWave cell may not be available either. Finally, the coordinator notifies the designated gNB, through a backhaul high-capacity link, about the optimal direction in which to steer the beam for serving each UE.

## 4.5   Performance Metrics and 3GPP Frameworks Parameters

In this section we define the metrics that will be used to compare and characterize the performance of the different beam management frameworks. Moreover, we will list the relevant parameters that affect the performance of the frameworks in 3GPP NR.

### 4.5.1   Performance Metrics

The performance of the different architectures and beam management procedures for IA and tracking will be assessed using three different metrics. The *detection accuracy* is measured in terms of probability of misdetection $P_{\mathrm{MD}}$, defined as the probability that the UE is not detected by the base station (i.e., the Signal-to-Noise-Ratio (SNR) is below a threshold $\Gamma$) in an uplink scenario, or, vice versa, the base station is not detected by the UE in a downlink scenario. The *reactiveness* differs according to the purpose of the measurement framework. For non-connected users, i.e., for IA, it is represented by the average time to find the best beam pair. For connected users, i.e., for tracking, it is the time required to receive the first CSI-RS after an SS burst, and thus react to channel variations or mobility in order to eventually switch beams, or declare a Radio Link Failure (RLF). Moreover, we also consider the time it takes to react to the RLF. Finally, the *overhead* is the amount of time and frequency resources allocated to the framework with respect to the total amount of available resources, taking into account both the IA (i.e., SS blocks or SRSs and the RACH) and the tracking (i.e., CSI-RSs).

### 4.5.2   3GPP Framework Parameters

In this section, we list the parameters that affect the performance of the measurement architectures, as summarized in Table 4.4. Moreover, we provide insights on the impact of each parameter on the different metrics.

**Frame Structure** − As depicted in Fig. 4.6, we consider the frame structure of 3GPP NR, with different subcarrier spacings $\Delta_f$. Given that in [18] the only subcarrier spacings considered for IA at frequencies above 6 GHz are $\Delta_f = 120$ and 240 kHz, i.e., $15 \times 2^n$ kHz, with $n \in [3, 4]$,

**Figure 4.6:** SS block structure. For configurations (a) and (b), each blue rectangle is an SS block (with 4 OFDM symbols) of duration 17.84 $\mu$s (i.e., $\Delta_f = 240$ kHz) and bandwidth $B_{SS} = 57.6$ MHz. For configurations (c) and (d) (for which $\Delta_f = 120$ kHz), instead, the blocks last 35.68 $\mu$s and have bandwidth $B_{SS} = 28.8$ MHz. Cases (a) and (c) implement a *frequency repetition* scheme (with $N_{rep} = 5$ and 11, respectively) while, for cases (b) and (d), a *data* solution (i.e., $N_{rep} = 1$) is preferred.

we will only consider these cases. The slot duration in ms is given by [219]

$$T_{\text{slot}} = \frac{1}{2^n}, \tag{4.1}$$

while the duration of a symbol in $\mu$s is [219]

$$T_{\text{symb}} = \frac{71.35}{2^n}. \tag{4.2}$$

Therefore, for $n = 3$ and 4 the slot duration is $125\,\mu$s or $62.5\,\mu$s, respectively. Moreover, according to the 3GPP specifications [18], the maximum number of subcarriers allocated to the SS blocks is 240, thus the bandwidth reserved for the SS blocks would be respectively 28.8 and 57.6 MHz. As mentioned in Sec. 4.3, we consider a maximum channel bandwidth $B = 400$ MHz per carrier [114].

**Frequency Diversity** – It is possible to configure the system to exploit frequency diversity, $D$. Given that 240 subcarriers are allocated in frequency to an SS, the remaining bandwidth in the symbols which contain an SS block is $B - 240\Delta_f$. Therefore, it is possible to adopt two different strategies: (i) *data* (as represented in Figs. 4.6(b) and (d)), i.e., the remaining bandwidth $B - 240\Delta_f$ is used for data transmission towards users , or (ii) *repetition* (as displayed in Figs. 4.6(a) and (c)), i.e., the information in the first 240 subcarriers is repeated in the remaining subcarriers to increase the robustness against noise and enhance the detection capabilities. The number of repetitions is therefore $N_{rep} = 1$ if frequency diversity is not used (i.e., $D = 0$, and a single chunk of the available bandwidth is used for the SS block), and $N_{rep} = 11$ or $N_{rep} = 5$ when repetition is used (i.e., $D = 1$) with $\Delta_f = 120$ kHz or $\Delta_f = 240$ kHz, respectively. There is a guard interval in frequency among the different repetitions of the SS blocks, to provide a good trade-off between frequency diversity and coherent combining [164]. Notice that 3GPP does not provide specifications for the repetition scheme.

**SS Block Configuration** – We consider different configurations of the SS blocks and bursts. The maximum number $N_{SS}$ of SS blocks in a burst for our frame structure and carrier frequencies is $L = 64$. We assume that, if $N_{SS} < L$, the SS blocks will be transmitted in the first $N_{SS}$ opportunities. The actual maximum duration of an SS burst is $D_{\max,SS} = 2.5$ ms for $\Delta_f = 240$ kHz and $D_{\max,SS} = 5$ ms for $\Delta_f = 120$ kHz. We will also investigate all the possible values

86

**Figure 4.7:** Relationship between beamwidth and antenna array size.

**Table 4.5:** Relationship between $M$, $\theta$ and $N_\theta$, for the azimuth case. Each gNB sector sweeps through $\Delta_{\theta,\mathrm{gNB}} = 120°$, while the UE scans over $\Delta_{\theta,\mathrm{UE}} = 360°$. In our evaluation, we consider a single antenna array at the UE modeled as a uniform rectangular array with isotropic antenna elements, following the approach of the literature [55]. Real handheld devices will be equipped with multiple patch antennas able to cover the whole angular space.

| $M$ | $\theta$ [deg] | $N_\theta$ gNB | $N_\theta$ UE |
|-----|------|------|------|
| 4   | 60   | 2    | 6    |
| 16  | 26   | 5    | 14   |
| 64  | 13   | 10   | 28   |

for the SS burst periodicity $T_{\mathrm{ss}}$, as defined in [204, 209], i.e., $T_{\mathrm{SS}} \in \{5, 10, 20, 40, 80, 160\}$ ms.

**CSI-RS Configuration** − As for the tracking, there are different options for the configuration of the CSI-RS structure. These options include (i) the number $N_{\mathrm{CSI}}$ of CSI-RS per SS burst period, (ii) the CSI-RS periodicity $T_{\mathrm{CSI,slot}} \in \{5, 10, 20, 40, 80, 160, 320, 640\}$ slots, and (iii) the offset $O_{\mathrm{CSI}}$ with respect to the end of an SS burst. In the analysis in Sec. 4.6 we will also refer to $T_{\mathrm{CSI}} = T_{\mathrm{CSI,slot}}T_{\mathrm{slot}}$, which represents the absolute CSI-RS periodicity in ms. These settings will be specified by the system information carried by the SS blocks of each burst. Other CSI-related parameters are the number of symbols of each CSI-RS transmission, i.e., $N_{\mathrm{symb,CSI}} \in \{1, 2, 4\}$, and the portion of bandwidth $\rho B$ allocated to the CSI-RSs. Moreover, the user will listen to $N_{\mathrm{CSI,RX}}$ CSI-RSs through an equivalent number of directions, when in connected state. We will consider $N_{\mathrm{CSI,RX}} \in \{1, 4\}$.

**Array Geometry** − As shown in Fig. 4.7 and Table 4.5, another fundamental parameter is the array geometry, i.e., the number of antenna elements $M$ at the gNB and UE and the number of directions that need to be covered, both in azimuth $N_\theta$ and in elevation $N_\phi$. In general, the antenna elements can be deployed as uniform linear or planar arrays, i.e., ULA and UPA respectively, and can be arranged as either rectangular or square arrays. Among the possible antenna designs, the most suitable approach is the use of UPAs, since they can enable 3D beamforming by adapting the beam in both azimuth and elevation planes [226]. In the simulations, the spacing of the elements is set to $\lambda/2$, where $\lambda$ is the wavelength, since this pattern was shown to offer excellent system capacity in small-cell urban deployments, as well as easy packageability (e.g., at 28 GHz, a $4 \times 4$ array has a size of roughly 1.5 cm $\times$ 1.5 cm) [164]. At the gNB we consider a single sector in a three sector site, i.e., the azimuth $\theta$ varies from $-60$ to 60 degrees, for a total of $\Delta_\theta = 120$ degrees. The elevation $\phi$ varies between $-30$ and 30 degrees, for a total of $\Delta_\phi = 60$ degrees, and also includes a fixed mechanical tilt of the array pointing towards the ground. There exists a strong correlation among beamwidth, number of antenna elements and BF gain. The more antenna elements in the system, the narrower the beams, the higher the gain that can be achieved by beamforming, and the more precise and directional the

transmission. Thus, given the array geometry, we compute the beamwidth $\Delta_{\text{beam}}$ at 3 dB of the main lobe of the beamforming vector, and then $N_\theta = \Delta_\theta / \Delta_{\text{beam}}$ and $N_\phi = \Delta_\phi / \Delta_{\text{beam}}$.

**Beamforming Architecture** –  Different beamforming architectures, i.e., analog, hybrid or digital, can be used both at the UE and at the gNB. *Analog beamforming* shapes the beam through a single RF chain for all the antenna elements, therefore the processing is performed in the analog domain and it is possible to transmit/receive in only one direction at any given time. This model saves power by using only a single pair of ADCs, but has a little flexibility since the transceiver can only beamform in one direction. *Hybrid beamforming* uses $K_{\text{BF}}$ RF chains (with $K_{\text{BF}} \leq M$), thus is equivalent to $K_{\text{BF}}$ parallel analog beams and enables the transceiver to transmit/receive in $K_{\text{BF}}$ directions simultaneously. Nevertheless, when hybrid beamforming is used for transmission, the power available at each transmitting beam is the total node power constraint divided by $K_{\text{BF}}$, thus potentially reducing the received power. The papers [227, 228] survey the main architectures for practical implementations of hybrid beamforming. *Digital beamforming* requires a separate RF chain and data converters for each antenna element and therefore allows the processing of the received signals in the digital domain, potentially enabling the transceiver to direct beams at infinitely many directions. Indeed, the availability of a sample for each antenna allows the transceiver to apply arbitrary weights to the received signals, and perform a more powerful and flexible processing than that in the analog domain. As in the hybrid case, the use of digital beamforming to transmit multiple beams simultaneously leads to a reduced transmit power being available to each (i.e., the total power constraint applies to the sum of all beams, not to each of them individually). Moreover, the digital transceiver can process at most $M$ simultaneous and orthogonal beams without any inter-beam interference (i.e., through a zero-forcing beamforming structure [229]). For this reason, we limit the number of parallel beams that can be generated to $M$. Furthermore, as previously mentioned, we implement a digital beamforming scheme only at the receiver side to avoid higher energy consumption in tranmsission. For the sake of completeness, we also consider an omnidirectional strategy at the UE i.e., without any beamforming gain but allowing the reception through the whole angular space at any given time.

Although currently available network deployments implementing beamforming capabilities use relatively small numbers of antennas, in the last few years new investigation towards the implementation of practical beamforming solutions for mmWave systems using a significantly larger amount of antenna elements in the array has been conducted both in the scientific community and in industry, e.g., in [230–232]. An overview of recent results and practical low-power architectures can be found in [233]. Nevertheless, significant research is still needed to uncover and solve the technical challenges remaining between the promises of beamforming schemes and their implementation in commercial systems.

**Network Deployment** –  Finally, the last parameters are the number of users $N_{\text{user}} \in \{5, 10, 20\}$ per sector of the gNBs and the density of base stations $\lambda_b$, expressed in gNB/km$^2$.

**Table 4.6:** Main simulation parameters.

| Parameter | Value | Description |
|---|---|---|
| $B$ | 400 MHz | Total bandwidth of each mmWave gNB |
| $f_c$ | 28 GHz | mmWave carrier frequency |
| $P_{\text{TX}}$ | 30 dBm | Transmission power |
| $\Gamma$ | $-5$ dB | SNR threshold |

**Table 4.7:** Notation.

| Symbol | Meaning |
|---|---|
| $\Delta_f$ | Subcarrier spacing |
| $T_{\text{slot}}$ | Duration of a slot |
| $T_{\text{symb}}$ | Duration of a symbol |
| $B$ | Bandwidth |
| $D$ | Usage of frequency diversity |
| $N_{rep}$ | Number of repetitions in frequency of an SS block |
| $P_{\text{MD}}$ | Probability of misdetection |
| $\Gamma$ | SNR threshold for the misdetection |
| $\lambda_b$ | gNB density |
| $N_{\text{SS}}$ | Number of SS blocks per burst |
| $L$ | Maximum number of SS blocks per burst |
| $D_{\text{max,SS}}$ | Maximum duration of an SS burst |
| $T_{\text{SS}}$ | SS burst periodicity |
| $S_D$ | Number of SS blocks for a complete sweep |
| $T_{\text{IA}}$ | Time required to perform IA |
| $T_{last}$ | Time to transmit the SS blocks in the last (or only) burst |
| $T_{\text{BR}}$ | Time to perform beam reporting during IA |
| $N_{\text{CSI}}$ | Number of CSI-RSs per SS burst periodicity |
| $T_{\text{CSI}}$ | CSI-RS periodicity |
| $T_{\text{CSI,slot}}$ | CSI-RS periodicity in slot |
| $O_{\text{CSI}}$ | Time offset between the end of the SS burst and the first CSI-RS |
| $N_{\text{symb,CSI}}$ | Number of OFDM symbols for a CSI-RS |
| $\rho$ | Portion of bandwidth $B$ for CSI-RSs |
| $N_{\text{CSI,RX}}$ | Number of directions that a UE monitors |
| $Z_{\text{CSI}}$ | Number of CSI-RSs to be transmitted |
| $T_{tot,\text{CSI}}$ | Time available for the CSI-RS transmission between two SS bursts |
| $N_{\text{CSI}}$ | Number of CSI-RS that can be transmitted between two bursts |
| $T_{\text{tr}}$ | Average time needed to receive the first CSI-RS |
| $N_{\text{CSI},\perp}$ | Number of orthogonal CSI-RSs between two SS bursts |
| $N_{\text{max,neigh}}$ | Number of neighbors that can be supported with orthogonal CSI-RSs |
| $T_{\text{RLF}}$ | RLF recovery delay |
| $M$ | Number of antenna elements at the transceiver |
| $\theta$ | Azimuth angle |
| $\phi$ | Elevation angle |
| $\Delta_\theta$ | Angular range for the azimuth |
| $\Delta_\phi$ | Angular range for the elevation |
| $N_\theta$ | Number of directions to cover in azimuth |
| $N_\phi$ | Number of directions to cover in elevation |
| $\Delta_{\text{beam}}$ | Beamwidth at 3 dB |
| $K_{\text{BF}}$ | Number of beams that the transceiver handles simultaneously |
| $N_{\text{user}}$ | Number of users |
| $R_{\text{SS}}$ | Time and frequency resources occupied by SS blocks |
| $\Omega_{\text{5ms}}$ | SS blocks overhead in 5 ms |
| $\Omega_{T_{\text{SS}}}$ | SS blocks overhead in $T_{\text{SS}}$ |
| $\Omega_{\text{CSI}}$ | CSI-RS overhead in $T_{\text{SS}}$ |
| $\Omega_{tot}$ | Total overhead in $T_{\text{SS}}$ |
| $\mathcal{U}[a,b]$ | Uniform random variable in the interval $[a,b]$ |

**Figure 4.8:** CDF of the SNR, for different antenna configurations. $\Delta_f = 120$ kHz, $N_{rep} = 0$. The red dashed line represents the SNR threshold $\Gamma = -5$ dB that has been considered throughout this work.

### 4.5.3 Channel Model

The simulations for the detection accuracy performance evaluation are based on realistic system design configurations. Our results are derived through a Monte Carlo approach, where multiple independent simulations are repeated, to get different statistical quantities of interest. The channel model is based on recent real-world measurements at 28 GHz in New York City, to provide a realistic assessment of mmWave micro and picocellular networks in a dense urban deployment. A complete description of the channel parameters can be found in [37], while the main simulation parameters for this paper are reported in Table 4.6.

## 4.6 Results and Discussion

In this section, we present some simulation results aiming at (i) evaluating the performance of the presented initial access schemes in terms of detection accuracy (i.e., probability of misdetection), as reported in Sec. 4.6.1; (ii) describing the analysis and the results related to the performance of the measurement frameworks for the reactiveness and the overhead, respectively in Sec. 4.6.2-4.6.3 and Sec. 4.6.4. Table 4.7 reports the notation used in this section.

### 4.6.1 Detection Accuracy Results

**Array size and gNB density** – Fig. 4.8 shows the Cumulative Distribution Function (CDF) of the SNR between the mobile terminal and the gNB it is associated to, for different antenna configurations and considering two density values. Notice that the curves are not smooth because of the progressive transitions of the SNR among the different path loss regimes, i.e., LOS, NLOS and outage. We see that better detection accuracy performance can be achieved when densifying the network and when using larger arrays. In the first case, the endpoints are progressively closer, thus ensuring better signal quality and, in general, stronger received power. In the second case, narrower beams can be steered thus guaranteeing higher gains produced by beamforming. We also notice that, for good SNR regimes, the $M_{\mathrm{gNB}} = 4, M_{\mathrm{UE}} = 4$ and $M_{\mathrm{gNB}} = 64, M_{\mathrm{UE}} = 4$ configurations present good enough SNR values: in these regions, the channel conditions are sufficiently good to ensure satisfactory signal quality (and, consequently, acceptable misdetection) even when considering small antenna factors. Finally, the red line represents the SNR threshold $\Gamma = -5$ dB that we will consider in this work.

90

**Figure 4.9:** $P_{\mathrm{MD}}$ as a function of $\lambda_b$, for different antenna configurations.



**Figure 4.10:** $P_{\mathrm{MD}}$ as a function of $\lambda_b$, for different subcarrier spacings $\Delta_f$ and repetition strategies and for different antenna configurations. $M_{\mathrm{gNB}} = 4, M_{\mathrm{UE}} = 4, \Gamma = -5$ dB.

Similar considerations can be deduced from Fig. 4.9, which illustrates how the misdetection probability monotonically decreases when the gNB density $\lambda_b$ progressively increases or when the transceiver is equipped with a larger number of antenna elements, since more focused beams can be generated in this case. Moreover, we notice that the beamforming strategy in which the UE transmits or receives omnidirectionally, although guaranteeing fast access operations, does not ensure accurate IA performance and leads to degraded detection capabilities. More specifically, the gap with a fully directional architecture (e.g., $M_{\mathrm{gNB}} = 64, M_{\mathrm{UE}} = 16$) is quite remarkable for very dense scenarios, and increases as the gNB density increases. For example, the configuration with 16 antennas (i.e., $M_{\mathrm{UE}} = 16$) and that with a single omnidirectional antenna at the UE reach the same $P_{\mathrm{MD}}$, but at different values of gNB density $\lambda_b$, respectively 30 and 35 gNB/km$^2$: the omnidirectional configuration requires a higher density (i.e., 5 gNB/km$^2$ more) to compensate for the smaller beamforming gain.

**Subcarrier spacing and frequency diversity** – Fig. 4.10 reports the misdetection probability related to $\lambda_b$, for different subcarrier spacings $\Delta_f$ and repetition strategies $D$. First, we see that, if no repetitions are used (i.e., $D = 0$), lower detection accuracy performance is associated with the $\Delta_f = 240$ kHz configuration, due to the resulting larger impact of the thermal noise and the consequent SNR degradation. Furthermore, the detection efficiency can be enhanced by repeating the SS block information embedded in the first 240 subcarriers in the remaining

91

subcarriers (i.e., $D = 1$), to increase the robustness of the communication and mitigate the effect of the noise in the detection process. In fact, if a frequency diversity approach is preferred, the UE (in the DL measurement technique) or the gNB (in the UL measurement technique) has $N_{rep} > 1$ attempts to properly collect the synchronization signals exchanged during the beam sweeping phase, compared to the single opportunity the nodes would have had if they had not implemented any repetition strategy. We also observe that the $\Delta_f = 120$ kHz with no frequency diversity configuration and the $\Delta_f = 240$ kHz scheme with $N_{rep} = 5$ produce the same detection accuracy results, thus showing how the effect of increasing the subcarrier spacing and the number of repetitions of the SS block information in multiple frequency subbands is similar in terms of misdetection capabilities. Finally, we observe that the impact of the frequency diversity $D$ and the subcarrier spacing $\Delta_f$ is less significant when increasing the array factor, as can be seen from the reduced gap between the curves plotted in Fig. 4.10 for the $M_{\mathrm{gNB}} = 4, M_{\mathrm{UE}} = 4$ and $M_{\mathrm{gNB}} = 64, M_{\mathrm{UE}} = 4$ configurations. The reason is that, when considering larger arrays, even the configuration with $\Delta_f = 240$ kHz and no repetitions has an average SNR which is high enough to reach small misdetection probability values.

### 4.6.2 Reactiveness Results for IA

**Analysis** – For initial access, reactiveness is defined as the delay required to perform a full iterative search in all the possible combinations of the directions. The gNB and the UE need to scan respectively $N_{\theta,\mathrm{gNB}} N_{\phi,\mathrm{gNB}}$ and $N_{\theta,\mathrm{UE}} N_{\phi,\mathrm{UE}}$ directions to cover the whole horizontal and vertical space. Moreover, they can transmit or receive respectively $K_{\mathrm{BF,gNB}}$ and $K_{\mathrm{BF,UE}}$ beams simultaneously. Notice that, as mentioned in Sec. 4.5.2, for digital and omnidirectional architectures $K_{\mathrm{BF}} = \min\{N_\theta N_\phi, M\}$, for hybrid $K_{\mathrm{BF}} = \min\{N_\theta N_\phi, M\}/\nu$, where $\nu$ is a factor that limits the number of directions in which it is possible to transmit or receive at the same time, and for analog $K_{\mathrm{BF}} = 1$.

Then the total number of SS blocks needed is[8]

$$S_D = \left\lceil \frac{N_{\theta,\mathrm{gNB}} N_{\phi,\mathrm{gNB}}}{K_{\mathrm{BF,gNB}}} \right\rceil \left\lceil \frac{N_{\theta,\mathrm{UE}} N_{\phi,\mathrm{UE}}}{K_{\mathrm{BF,UE}}} \right\rceil. \tag{4.3}$$

Given that there are $N_{\mathrm{SS}}$ blocks in a burst, the total delay from the beginning of an SS burst transmission in a gNB to the completion of the sweep in all the possible directions is

$$T_{\mathrm{IA}} = T_{\mathrm{SS}} \left( \left\lceil \frac{S_D}{N_{SS}} \right\rceil - 1 \right) + T_{last}, \tag{4.4}$$

where $T_{last}$ is the time required to transmit the remaining SS blocks in the last burst (notice that there may be just one burst, thus the first term in Eq. (4.4) would be 0). This term depends on the subcarrier spacing and on the number of remaining SS blocks which is given by

$$N_{\mathrm{SS,left}} = S_D - N_{\mathrm{SS}} \left( \left\lceil \frac{S_D}{N_{\mathrm{SS}}} \right\rceil - 1 \right). \tag{4.5}$$

---

[8]We recall that hybrid or digital architectures consume more power than analog ones, if the same number of bits in the ADCs is used, and thus are more likely to be implemented only at the receiver side. Nevertheless, some ADC configurations enable energy efficient digital beamforming (e.g., 3 bits ADC [44]), with a power consumption comparable to that of an analog implementation.

**(a)** gNB Analog, UE Analog

**(b)** gNB Analog, UE Hybrid (DL-based configuration)

**(c)** gNB Analog, UE Digital (DL-based configuration)

**(d)** gNB Digital, UE Analog (UL-based configuration)

**Figure 4.11:** $T_{\text{IA}}$ as a function of $N_{\text{SS}}$ with $T_{\text{SS}} = 20$ ms.

Then, $T_{last}$ is

$$T_{last} = \begin{cases} \frac{N_{\text{SS,left}}}{2} T_{slot} - 2T_{symb} & \text{if } N_{\text{SS,left}} \bmod 2 = 0 \\ \left\lfloor \frac{N_{\text{SS,left}}}{2} \right\rfloor T_{slot} + 6T_{symb} & \text{otherwise,} \end{cases} \tag{4.6}$$

The two different options account for an even or odd remaining number of SS blocks. In the first case, the SS blocks are sent in $N_{\text{SS,left}}/2$ slots, with total duration $N_{\text{SS,left}}/2T_{slot}$, but the last one is actually received in the 12*th* symbol of the last slot, i.e., 2 symbols before the end of that slot, given the positions of the SS blocks in each slot described in [18, 208]. If instead $N_{\text{SS,left}}$ is odd, six symbols of slot $\lfloor N_{\text{SS,left}}/2 \rfloor + 1$ are also used.

A selection of results is presented in the next paragraphs.

**Number of SS blocks per burst and beamforming technology** – In Fig. 4.11 we consider first the impact of the number of SS blocks in a burst, with a fixed SS burst periodicity $T_{\text{SS}} = 20$ ms and for different beamforming strategies and antenna configurations. In particular in Fig. 4.11a, in which both the UE and the gNB use analog beamforming, the initial access delay heavily depends on the number of antennas at the transceivers since all the available directions must be scanned one by one. It may take from 0.6 s (with $N_{\text{SS}} = 64$) to 5.2 s (with $N_{\text{SS}} = 8$) to transmit and receive all the possible beams, which makes the scheme infeasible for practical usage. A reduction in the sweeping time can be achieved either by using an omnidirectional antenna at the UE or by decreasing the number of antennas both at the UE and at the gNB. In this case, the only configurations that manage to complete a scan in a single SS burst are those with 4 antennas at both sides and $N_{\text{SS}} \geq 16$, or that with $M_{\text{gNB}} = 64$, an omnidirectional UE and $N_{\text{SS}} = 64$.

**Figure 4.12:** $T_{IA}$ as a function of $T_{SS}$ for the downlink configuration with analog gNB and hybrid UE. $N_{\mathrm{SS}} = 64$



**Figure 4.13:** $T_{\mathrm{IA}}$ for different antenna configurations and subcarrier spacing $\Delta_f$, with gNB Analog, UE Analog.

Another option is the usage of hybrid or digital beamforming at the UE in a downlink-based scheme, or at the eNB in an uplink-based one. Fig. 4.11b shows $T_{\mathrm{IA}}$ when the UE uses hybrid beamforming to receive from half of the available directions at any given time (i.e., $L = 2$), while in Fig. 4.11c the UE receives from all available directions at any given time. This leads to an increased number of configurations which are able to complete a sweep in an SS block, even with a large number of antennas at the gNB and the UE.

Finally, Fig. 4.11d shows the performance of an uplink-based scheme, in which the SRSs are sent in the same time and frequency resource in which the SS blocks would be sent, and the gNB uses digital beamforming. It can be seen that there is a gain in performance for most of the configurations, because the gNB has to sweep more directions than the UE (since it uses narrower beams), thus using digital beamforming at the gNB-side makes it possible to reduce $T_{\mathrm{IA}}$ even more than when it is used at the UE-side.

**SS burst periodicity** – For the setup with hybrid beamforming at the UE, that generally requires more than one SS burst periodicity, we show in Fig. 4.12 the dependency of $T_{\mathrm{IA}}$ and $T_{\mathrm{SS}}$. It can be seen that the highest periodicities are not suited for a mmWave deployment, and that in general it is better to increase the number of SS blocks per burst in order to try to complete the sweep in a single burst.

**Subcarrier spacing** – Another parameter that has an impact on $T_{\mathrm{IA}}$ is the subcarrier spacing $\Delta_f$. As shown in Fig. 4.13, when the larger spacing is used the OFDM symbols have a shorter duration and the transmission of the SS blocks in the directions of interest can be completed earlier.

**Table 4.8:** Reactiveness performance for beam reporting operations considering an SA or an NSA architecture. Analog or digital beamforming is implemented at the gNB side, while the UE configures its optimal beamformed direction. $T_{\mathrm{SS}} = 20$ ms, $\Delta_f = 120$ KHz.

| | $T_{\mathrm{BR,SA}}$ [ms] | | | |
| | $N_{\mathrm{SS}} = 8$ | | $N_{\mathrm{SS}} = 64$ | |
| $M_{gNB}$ | Analog | Digital | Analog | Digital |
|---|---|---|---|---|
| 4 | 0.0625 | 0.0625 | 0.0625 | 0.0625 |
| 16 | 0.5 | 0.0625 | 0.5 | 0.0625 |
| 64 | 40.56 | 0.0625 | 1.562 | 0.0625 |
| $T_{\mathrm{BR,NSA}} \in \{10, 4, 0.8\}$ ms, according to [21]. | | | | |

**Impact of Beam Reporting** − For initial access, in addition to the time required for directional sweeping, there is also a delay related to the allocation of the resources in which it is possible to perform initial access, which differs according to the architecture being used. As introduced in Sec. 4.4, 3GPP advocates the implicit reporting of the chosen direction, e.g., the strongest SS block index, through contention-based random access messages, agreeing that the network should allocate multiple RACH transmissions and preambles to the UE for conveying the optimal SS block index to the gNB [7, 234]. When considering an SA configuration, beam reporting might require an additional sweep at the gNB side while, if an NSA architecture is preferred, the beam decision is forwarded through the LTE interface (and requires just a single RACH opportunity) which makes the beam reporting reactiveness equal to the latency of a legacy LTE connection. Assuming a 0% BLER data channel, the uplink latency in legacy LTE, including scheduling delay, ranges from 10.5 ms to 0.8 ms, according to the latency reduction techniques being implemented [21].

In Table 4.8, we analyze the impact of the number of SS blocks (and, consequently, of RACH opportunities) in a burst, with a fixed burst periodicity $T_{\mathrm{SS}} = 20$ ms and for a subcarrier spacing of $\Delta_f = 120$ KHz. The results are independent of the antenna configuration at the UE side, since the mobile terminal steers its beam through the previously determined optimal direction and does not require a beam sweeping operation to be performed. It appears clear that the SA scheme presents very good reactiveness for most of the investigated configurations and, most importantly, outperforms the NSA solution even when the LTE latency is reduced to 0.8 ms. The reason is that, if the network is able to allocate the needed RACH resources within a single SS burst, then it is possible to limit the impact of beam reporting operations on the overall initial access reactiveness, which is instead dominated by the beam sweeping phase. In particular, when considering small antenna factors and when digital beamforming is employed, beam reporting can be successfully completed through a single RACH allocation, thus guaranteeing very small delays.

### 4.6.3 Reactiveness Results for Beam Tracking

**Analysis** − For tracking, we define the reactiveness as the average time needed to receive the first CSI-RS after the end of each SS burst.

We assume that the $N_{\mathrm{user}}$ UEs are uniformly distributed in the space covered by the $k = N_{\theta,\mathrm{gNB}} N_{\phi,\mathrm{gNB}}$ beams available at the gNB. Moreover, each UE has to monitor $N_{\mathrm{CSI,RX}}$ directions. Given that a UE may or may not be in LOS, it is not obvious that these directions will be associated to the closest beams with respect to the one selected during the initial access. Therefore, we also assume that this scenario is equivalent to a scenario with $n = N_{\mathrm{user}} N_{\mathrm{CSI,RX}}$ uniformly distributed UEs, each of them monitoring a single direction. We will refer to $n$ as the

number of measures.

Consequently, on average there are $n/k$ measurements for the area belonging to each beam, if the beams divide the space into equally sized regions. Therefore, if $n \geq k$, a CSI-RS is needed in each beam, otherwise it is sufficient to send at least $n$ CSI-RSs, and thus the total number of CSI-RS that need to be transmitted is on average $Z_{\mathrm{CSI}} = \min\{n, k\}$.

Depending on the combination of $T_{\mathrm{SS}}$, $T_{\mathrm{CSI}} = T_{\mathrm{CSI,slot}} T_{\mathrm{slot}}$ and $Z_{\mathrm{CSI}}$, it may not be possible to allocate all the CSI-RS transmissions between two consecutive SS bursts. Notice that after the end of an SS burst, there are $T_{tot,\mathrm{CSI}} = T_{\mathrm{SS}} - D_{\mathrm{max,SS}}$ ms available for the CSI-RS transmission. Then, the number $N_{\mathrm{CSI}}$ of CSI-RS that can be allocated between two SS bursts may depend on which of the options shown in Fig. 4.2 is chosen.

*Option 1:* the first CSI-RS is transmitted $T_{\mathrm{CSI}}$ ms after the transmission of the SS burst. In this case, $N_{\mathrm{CSI}} = \lfloor T_{tot,\mathrm{CSI}}/T_{\mathrm{CSI}} \rfloor$, and single periodicity is not enough if $Z_{\mathrm{CSI}} > N_{\mathrm{CSI}}$. For option 1, the metric $T_{\mathrm{tr,opt1}}$ is given by

$$T_{\mathrm{tr,opt1}} = \frac{\sum_{p=0}^{\left\lfloor \frac{Z_{\mathrm{CSI}}}{N_{\mathrm{CSI}}} \right\rfloor - 1} \left( \sum_{i=1}^{N_{\mathrm{CSI}}} (pT_{\mathrm{SS}} + iT_{\mathrm{CSI}}) \right) + \sum_{i=1}^{Z_{\mathrm{CSI}} \bmod N_{\mathrm{CSI}}} \left( \left\lfloor \frac{Z_{\mathrm{CSI}}}{N_{\mathrm{CSI}}} \right\rfloor T_{\mathrm{SS}} + iT_{\mathrm{CSI}} \right)}{Z_{\mathrm{CSI}}}. \quad (4.7)$$

The last sum accounts for the case $Z_{\mathrm{CSI}} < N_{\mathrm{CSI}}$ and for the CSI-RS in the last SS burst periodicity when $Z_{\mathrm{CSI}} > N_{\mathrm{CSI}}$. The sum over $p$, instead, accounts for $Z_{\mathrm{CSI}} \geq N_{\mathrm{CSI}}$.

*Option 2:* thanks to the additional parameter $O_{\mathrm{CSI}}$ it is possible to transmit

$$N_{\mathrm{CSI}} = \lceil T_{tot,\mathrm{CSI}}/T_{\mathrm{CSI}} \rceil, \quad (4.8)$$

as shown in Fig. 4.2b. The offset is computed as

$$O_{\mathrm{CSI}} = \frac{T_{tot,\mathrm{CSI}} - (N_{\mathrm{CSI}} - 1)T_{\mathrm{CSI}}}{2}. \quad (4.9)$$

The metric $T_{\mathrm{tr,opt2}}$ is computed as for option 1, but taking into account also $O_{\mathrm{CSI}}$:

$$T_{\mathrm{tr,opt2}} = \frac{\begin{array}{c} \sum_{p=0}^{\left\lfloor \frac{Z_{\mathrm{CSI}}}{N_{\mathrm{CSI}}} \right\rfloor - 1} \left( \sum_{i=0}^{N_{\mathrm{CSI}}-1} (pT_{\mathrm{SS}} + iT_{\mathrm{CSI}} + O_{\mathrm{CSI}}) \right) + \\ \sum_{i=0}^{Z_{\mathrm{CSI}} \bmod N_{\mathrm{CSI}} - 1} \left( \left\lfloor \frac{Z_{\mathrm{CSI}}}{N_{\mathrm{CSI}}} \right\rfloor T_{\mathrm{SS}} + iT_{\mathrm{CSI}} + O_{\mathrm{CSI}} \right) \end{array}}{Z_{\mathrm{CSI}}}. \quad (4.10)$$

Notice that if $Z_{\mathrm{CSI}} > N_{\mathrm{CSI}}$, a signal in a certain direction could be either received as SS block in the next burst, or as CSI-RS, depending on how the transmission of SS blocks and CSI-RSs is scheduled.

**Scheduling options, number of users and CSI-RS periodicity** – Fig. 4.14a shows the value of $T_{tr}$ for different parameters, such as the different scheduling option 1 or 2, the number of users per gNB $N_{\mathrm{user}}$ and of directions of interest $N_{\mathrm{CSI,RX}}$, for SS burst periodicity $T_{\mathrm{SS}} = 20$ ms and 64 antennas at the gNB. The fundamental parameter is the periodicity of the CSI-RS transmission: only a small CSI-RS periodicity makes it possible to sweep all the directions to be covered during a relatively short interval, and to avoid the dependency on $T_{\mathrm{SS}}$. Moreover, if the periodicity is small (i.e., $T_{\mathrm{CSI}} = 0.625$ ms, or 5 slots with $\Delta_f = 120$ kHz), then there is no difference between the two scheduling options, while this becomes notable for $T_{\mathrm{CSI}} = 10$ ms, as expected.

**SS burst periodicity** – Fig. 4.14b compares two different $T_{\mathrm{SS}}$ periodicities, i.e., 10 and 40 ms, using the smallest $T_{\mathrm{CSI,slot}}$ available (i.e., 5 slots, or 0.625 ms at $\Delta_f = 120$ kHz). It can be

seen that using a higher $T_{\text{SS}}$ would allow a decreased $T_{tr}$, since more CSI-RSs can be scheduled between two SS bursts and consequently a larger number of directions can be swept. For the sake of completeness, Fig. 4.15 shows the number of CSI-RSs that can be scheduled in between two SS bursts as a function of $T_{\text{SS}}$ and of the different scheduling options and periodicities. Since in a mmWave scenario there may be a need to scan a large number of CSI-RSs, it is advisable to either use an adaptive scheme for the scheduling of CSI-RSs, which adapts the periodicity according to the number of users in the different directions, or adopt a conservative approach and use a short $T_{\text{CSI}}$ interval.

**Limits on the CSI-RS periodicity** – Since the CSI-RSs that a user receives from multiple base stations should not overlap in time and frequency (otherwise the RSRP value would be overestimated), there is a maximum number of neighboring cells that a gNB can support. According to [18], there are 4 symbols per slot in which a CSI-RS can be sent (additional symbols are under discussion), and a CSI-RS can last 1, 2 or 4 symbols, each with bandwidth $\rho B$. Assuming a common configuration for the gNBs deployed in a certain area, the total number of orthogonal CSI-RS transmission opportunities is

$$N_{\text{CSI},\perp} = \frac{T_{\text{SS}} - D\text{max}, \text{SS}}{T_{\text{slot}}} \frac{4}{N_{\text{symb,CSI}}} \left\lfloor \frac{1}{\rho} \right\rfloor, \tag{4.11}$$

where the first ratio is the number of slots in the time interval in which CSI-RSs can be scheduled, and the second and third express the number of CSI-RSs per slot (there are at most 4 OFDM symbols per slot for CSI-RSs). Then, the maximum number of neighbors that a gNB can support is

$$N_{\text{max,neigh}} = \left\lfloor \frac{N_{\text{CSI},\perp}}{N_{\text{CSI}}} \right\rfloor - 1, \tag{4.12}$$

with $N_{\text{CSI}}$ computed as in the previous paragraphs.

Fig. 4.16 reports the value of $N_{\text{max,neigh}}$ for a different number of OFDM symbols for the CSI-RSs and bandwidth scaling factor $\rho$, which ranges from 0.1 to 1, and represents also the bandwidth values corresponding to 240 subcarriers with $\Delta_f \in \{120, 240\}$ kHz, i.e., the bandwidth occupied by an SS burst. Notice that for the frequencies in the mmWave spectrum it is advisable not to use the entire bandwidth for CSI-RSs [212], and the number of neighbors of a mmWave



**(a)** $M_{\text{gNB}} = 64$, analog beamforming, $T_{\text{SS}} = 20$ ms

**(b)** $M_{\text{gNB}} = 64$, analog beamforming, $T_{\text{CSI}} = 0.625$ ms

**Figure 4.14:** Performance of tracking using CSI-RSs for Option 1 and Option 2, as described in Fig. 4.2, as a function of different parameters (e.g., $T_{\text{CSI}}$, $T_{\text{SS}}$), for $\Delta_f = 120$ kHz.

**Figure 4.15:** $N_{\mathrm{CSI}}$ as a function of the $T_{\mathrm{SS}}$ and $T_{\mathrm{CSI}}$ periodicities.



**Figure 4.16:** $N_{\mathrm{max,neigh}}$ as a function of $N_{\mathrm{symb,CSI}}$ and $\rho$ for different $T_{\mathrm{CSI}}$ periodicities, with $T_{\mathrm{SS}} = 20$ ms and $\Delta_f = 120$ kHz.

gNB will be limited, given the short propagation distance typical of these frequencies. If $T_{\mathrm{CSI}} = 10$ ms, then even when using 4 OFDM symbols and the whole bandwidth it is possible to support only 14 neighbors. Instead, when $T_{\mathrm{CSI}} = 0.625$ ms it is not feasible to use the whole bandwidth and 4 symbols, but more conservative configurations should be adopted. For example, with $\rho = 0.072$ (i.e., 240 subcarriers with $\Delta_f = 120$ kHz) it is possible to support 15 or 31 neighbors, respectively with 2 or 1 OFDM symbols.

**Standalone vs non-standalone** – Notice that when the standalone scheme is used and the UE experiences a link failure on all the $N_{\mathrm{CSI,RX}}$ directions it is monitoring, then the UE has no choice but using the SS blocks in the SS burst to perform either a link recovery or a new initial access, and meanwhile it is not able to transmit or receive data or control information [388]. When a non-standalone architecture is used, instead, the UE could signal this event to the RAN on the lower-frequency control link, and the data plane can be switched to the sub-6-GHz RAT, and faster recovery options could be designed, for example, by instructing the UE to monitor additional CSI-RSs.

**Downlink vs uplink and beamforming architecture** – Finally, we observe that, when a digital architecture is chosen, there exist some specific configurations in which a UL-based

**Table 4.9:** RLF recovery delay considering the SA or the NSA measurement frameworks, for different values of $N_{\mathrm{SS}}$, $T_{\mathrm{SS}}$ and for different beamforming configurations. $\Delta_f = 120$ kHz. ABF stands for Analog Beamforming, and DBF for Digital.

| Antenna | | $T_{\mathrm{RLF,SA}}$ [ms] | | |
|---|---|---|---|---|
| $M_{gNB}$ | $M_{UE}$ | $N_{\mathrm{SS}} = 8$, $T_{\mathrm{SS}} = 20$ $gNB$ ABF, UE ABF | $N_{\mathrm{SS}} = 64$, $T_{\mathrm{SS}} = 40$ $gNB$ DBF, UE ABF | $N_{\mathrm{SS}} = 64$, $T_{\mathrm{SS}} = 80$ $gNB$ DBF, UE ABF |
| 4 | 4 | 30.2322 | 20.3572 | 40.3572 |
| 64 | 1 | 130.1072 | 20.0535 | 40.0535 |
| 64 | 16 | 5250 | 22.6072 | 42.6072 |
| $T_{\mathrm{RLF,NSA}} \in \{10, 4, 0.8\}$ ms, according to the considerations in [21]. | | | | |

measurement framework can ensure more efficient tracking operations than its DL counterpart. In fact, due to the gNB's less demanding space constraints with respect to a mobile terminal, a larger number of antenna elements can usually be packed at the base station side, resulting in a larger number of directions that can potentially be scanned simultaneously through a digital beamforming scheme. Moreover, hybrid or fully digital receivers are more costly in terms of power consumption, and hence are more likely to be implemented in a gNB rather than in a UE.

**RLF recovery** – Another important factor that affects the reactiveness of beam management schemes is the time it takes to recover from an RLF. As assumed by 3GPP [204, 235], RLF occurs when the quality of an associated control channel falls below a certain threshold. As soon as the failure is detected, mechanisms to recover acceptable communication capabilities (e.g., by determining an alternative suitable direction of transmission or possibly handing over to a stronger and more robust gNB) need to be quickly triggered upon notifying the network. Natural candidates for monitoring the link quality and detect the link failure are the SS blocks in a burst [204, 236]. Assume that an object blocks the propagation path of the transceiver at time $T \sim \mathcal{U}[t, t + T_{\mathrm{SS}}]$, i.e., on average at time $\bar{T} = T_{\mathrm{SS}}/2$ within two consecutive SS bursts.

- When implementing an SA architecture, as soon as an impairment is detected, the UE may no longer be able to communicate with its serving gNB since the optimal directional path connecting the endpoints is affected by the failure. The recovery phase is most likely triggered at the beginning of the subsequent SS burst (i.e., on average after $T_{SS} - \bar{T} = T_{\mathrm{SS}}/2$ seconds) and at least after the completion of an IA operation of duration $T_{\mathrm{IA}}$ seconds.[9] Table 4.9 reports the RLF recovery delay $T_{\mathrm{RLF,SA}}$ for some network configurations when an SA architecture is implemented. We observe that the latency is quite high for all the investigated settings and is dominated by the IA delay, as illustrated in Fig. 4.11. Moreover, in some circumstances (e.g., $N_{\mathrm{SS}} = 8$, $T_{\mathrm{SS}} = 20$ ms, $M_{gNB} = 64$, $N_{gNB} = 16$ and when analog beamforming is implemented), the RLF recovery delay assumes unacceptably high values.

- Much more responsive RLF recovery operations may be prompted if the failure notification is forwarded through the LTE overlay (i.e., by implementing an NSA-based measurement framework), which may also serve the UE's traffic requests until the mmWave directional communication is successfully restored. If an NSA-DL framework is designed, the RLF recovery delay $T_{\mathrm{RLF,NSA}}$ is equal to the latency of a traditional LTE connection (which depends on the implemented latency reduction technique, as assessed in [21]). Alternatively, the gNB can autonomously declare an RLF event (without the user's notification) and

---

[9]In some circumstances, the UE can autonomously react to an RLF event by selecting an alternative direction of communication, as a sort of backup solution before the transceiver fully recovers the optimal beam configuration [163]. Although having a second available link, when the primary path is obstructed, adds diversity and robustness to the communication, it may not always guarantee sufficiently good communication performance.

react accordingly by monitoring the SRS messages. Without loss of generality, assuming that SRSs are uniformly allocated within two SS bursts with periodicity $T_{\text{SRS}}$, an RLF is detected as soon as the gNB is not able to correctly receive $N_{\text{SRS}}$ consecutive SRSs from its reference user. In this case, the reactiveness of the RLF recovery operation depends on the periodicity of the sounding signals and is equal to

$$T_{\text{RLF,NSA}} = \frac{T_{\text{SRS}}}{2} + (N_{\text{SRS}} - 1) T_{\text{SRS}}. \tag{4.13}$$

Analogously, if an NSA-UL framework is designed, the recovery may be immediately triggered by the gNB by switching the traffic to the LTE eNB in $T_{\text{RLF,NSA}}$ seconds, as given by Eq. (4.13). From the results in Table 4.9, it appears that fast and efficient RLF recovery operations can be guaranteed if an NSA solution is preferred over an SA one for all the investigated network configurations.

### 4.6.4   Overhead Results

In this section, we characterize the overhead for IA and tracking in terms of the ratio between the time and frequency resources that are allocated to SS bursts and the maximum duration of the SS burst (i.e., 5 ms), or the entire $T_{\text{SS}}$ interval.

**Analysis** − The total number of time and frequency resources $R_{\text{SS}}$ scheduled for the transmission of $N_{\text{SS}}$ SS blocks, each spanning 4 OFDM symbols and 240 (or multiple of 240) subcarriers, is given by

$$R_{\text{SS}} = N_{\text{SS}} \, 4T_{\text{symb}} \, 240 N_{rep} \Delta_f, \tag{4.14}$$

where $T_{\text{symb}}$ is expressed in ms and $\Delta_f$ in kHz. The overhead for the 5 ms time interval with the SS burst transmission and total bandwidth $B$ (in Hz) is then given by

$$\Omega_{\text{5ms}} = \frac{N_{\text{SS}} \, 4T_{\text{symb}} \, 240 N_{rep} \Delta_f}{5B}, \tag{4.15}$$

and the overhead considering the total burst periodicity $T_{\text{SS}}$ is

$$\Omega_{\text{T}_{\text{SS}}} = \frac{N_{\text{SS}} \, 4T_{\text{symb}} \, 240 N_{rep} \Delta_f}{T_{\text{SS}} B}. \tag{4.16}$$

Moreover, additional overhead is introduced by the transmission of CSI-RSs after the SS burst. The value of the overhead $\Omega_{\text{CSI}}$ depends on the number of symbols $N_{\text{symb,CSI}}$ and the bandwidth $\rho B$ for each CSI-RS, as well as on the number of CSI-RSs $N_{\text{CSI}}$ computed as in Sec. 4.6.3 for the two CSI-RS scheduling options:

$$\Omega_{\text{CSI}} = \frac{N_{\text{CSI}} N_{\text{symb,CSI}} T_{\text{symb}} \rho B}{(T_{\text{SS}} - D_{\text{max,SS}}) B} = \frac{N_{\text{CSI}} N_{\text{symb,CSI}} T_{\text{symb}} \rho}{(T_{\text{SS}} - D_{\text{max,SS}})}. \tag{4.17}$$

Finally, the total overhead $\Omega$ takes into account both the SS bursts and the CSI-RSs in $T_{\text{SS}}$:

$$\Omega_{tot} = \frac{N_{\text{CSI}} N_{\text{symb,CSI}} T_{\text{symb}} \rho B + R_{\text{SS}}}{T_{\text{SS}} B}. \tag{4.18}$$

**Subcarrier spacing and frequency diversity** − Fig. 4.17 reports the overhead related to the maximum duration of the SS burst (i.e., 5 ms) for different subcarrier spacings and repetition strategies. It can be seen that if no repetitions are used (i.e., $D = 0$) then the overheads for the configurations with $\Delta_f = 120$ kHz and $\Delta_f = 240$ kHz are equivalent. In fact, when configuring

**(a)** $\Omega_{5\mathrm{ms}}$ as a function of $N_{\mathrm{SS}}$, for different subcarrier spacings $\Delta_f$ and repetition strategies.

**(b)** $\Omega_{T_{\mathrm{SS}}}$ as a function of $T_{\mathrm{SS}}$, for different subcarrier spacings $\Delta_f$ and repetition strategies. $N_{\mathrm{SS}}$ is set to the maximum value, i.e., 64.

**Figure 4.17:** Overhead for initial access, introduced by the transmission of the SS blocks. Notice that the number of repetitions for the different subcarrier spacings $\Delta_f$ is chosen to send as many repetitions of the SS blocks as possible.



**(a)** Overhead $\Omega_{\mathrm{CSI}}$ as a function of $N_{\mathrm{symb,CSI}}$ and $\rho$, for different $T_{\mathrm{CSI}}$ periodicities, with $T_{\mathrm{SS}} = 20$ ms.

**(b)** Overhead $\Omega_{tot}$ as a function of $N_{\mathrm{symb,CSI}}$ and $\rho$, for different subcarrier spacings $\Delta_f$ and repetition strategies. $N_{\mathrm{SS}}$ is set to the maximum value, i.e., 64, and $T_{\mathrm{CSI,slot}} = 5$ slot.

**Figure 4.18:** Overhead for the CSI-RS transmission and total overhead, with $T_{\mathrm{SS}} = 20$ ms. Notice that the number of repetitions for the different subcarrier spacings $\Delta_f$ is chosen to send as many repetitions of the SS blocks as possible.

large subcarrier spacings (i.e., $\Delta_f = 240$ kHz), the OFDM symbols used for the SS blocks have half the duration, but they occupy twice the bandwidth of the systems with narrower subcarrier spacings (i.e., $\Delta_f = 120$ kHz), given that the same number of subcarriers are used. Instead, when a repetition strategy is used (i.e., $D = 1$), the overhead is different. As mentioned in Sec. 4.5.2, we consider 5 repetitions for $\Delta_f = 240$ kHz and 11 for $\Delta_f = 120$ kHz. Therefore, the actual amount of bandwidth that is used is comparable, but since the OFDM symbols with $\Delta_f = 120$ kHz last twice as long as those with the larger subcarrier spacing, the overhead in terms of resources used for the SS burst is higher with $\Delta_f = 120$ kHz.

**SS burst periodicity** − Fig. 4.17b shows the dependency of the overhead for initial access on $T_{\mathrm{SS}}$, which follows an inverse proportionality law. In particular, for very small $T_{\mathrm{SS}}$ (i.e., 5 ms) the impact of the SS bursts with repetitions in frequency is massive, with up to 43% of the resources allocated to the SS blocks. For $T_{\mathrm{SS}} = 20$ ms or higher, instead, the overhead is always

**Table 4.10:** Overhead for beam reporting operations considering an SA architecture. Analog or digital beamforming is implemented at the gNB side, for different antenna array structures.

| | $\Omega_{\text{BR,SA}} \cdot 10^{-3}$ | | | |
| | $\Delta_{f,\text{RACH}} = 60$ kHz | | $\Delta_{f,\text{RACH}} = 120$ kHz | |
| $M_{gNB}$ | Analog | Digital | Analog | Digital |
|---|---|---|---|---|
| 4 | 0.0894 | 0.0894 | 0.0894 | 0.0894 |
| 16 | 0.7149 | 0.0894 | 0.7149 | 0.0894 |
| 64 | 2.2341 | 0.0894 | 2.2341 | 0.0894 |

below 10%.

**CSI-RS periodicity** – The overhead due to the transmission of CSI-RSs is shown in Fig. 4.18a for different $T_{\text{CSI}}$ periodicities and time and frequency resource allocation to the CSI-RSs. It is always below 0.008 with $T_{\text{CSI}} = 5$ ms, and below 0.06 for $T_{\text{CSI}} = 0.625$ ms. However, for practical values of the configuration of the CSI-RSs, in which the bandwidth for the reference signal is smaller than half of the entire bandwidth, then also for $T_{\text{CSI}} = 0.625$ ms the overhead reaches very small values, i.e., below 0.028.

**Impact of IA and tracking** – The trend of $\Omega_{tot}$ is shown in Fig. 4.18b, where it can be immediately seen that the largest impact is given by the term $R_{\text{SS}}$ at the numerator and not by the CSI-RS-related overhead. The parameters on the $x$ and $y$ axes have indeed a limited effect on the gradient of the surfaces, which are almost horizontal. The main difference is introduced by the different subcarrier spacings and repetition strategies. Notice that, contrary to what is shown in Fig. 4.17a, there is a difference between the two different subcarrier spacings for the total overhead $\Omega_{tot}$ and for the CSI-RS-related overhead $\Omega_{\text{CSI}}$, because we consider a different $T_{symb}$ in Eq. (4.17), but the same $\rho$ factor, thus a different number of subcarriers for the different values of $\Delta_f$.

**Impact of beam reporting** – For the SA case, as reported in Table 4.10, the completion of the beam reporting procedure for initial access may require an additional overhead, due to the need for the system to allocate possibly multiple RACH resources[10] for the reporting operations. Conversely, for the NSA case, the beam decision is forwarded through the LTE overlay and requires a single RACH opportunity, with a total overhead of $0.0894 \cdot 10^{-3}$. Nevertheless, from Table 4.10, we notice that the SA additional reporting overhead is quite limited due to the relatively small number of directions that need to be investigated at this stage, especially when designing digital beamforming solutions.

## 4.7 Comprehensive Considerations on 3GPP NR Beam Management

This section provides some comprehensive considerations on the metrics which we separately studied in Sec. 4.6. The goal is to highlight which are the main trade-offs between accuracy, responsiveness and overhead and the design parameters for beam management in NR. Moreover, we provide some suggestions and guidelines to optimally dimension a measurement framework for NR at mmWave frequencies.

**Subcarrier spacing** $\Delta_f$    When using a smaller subcarrier spacing (i.e., $\Delta_f = 120$ kHz) it is possible to achieve a higher accuracy (i.e., smaller misdetection probability), either because the impact

---

[10]According to the 3GPP agreements [237], a bandwidth of 10 MHz (for $\Delta_{f,\text{RACH}} = 60$ kHz) or a bandwidth of 20 MHz (for $\Delta_{f,\text{RACH}} = 120$ kHz) is reserved for the RACH resources.

of the noise is less relevant, when frequency diversity is not used, or because it is possible to allocate a larger number of repetitions, when frequency diversity is used. This last option comes however at the price of an increase in the overhead in the order of 2 times, while the accuracy gain for the configuration with $\lambda = 30$ gNB/km$^2$ and the $4 \times 4$ antenna arrays is in the order of 23%, according to Fig. 4.10. A smaller subcarrier spacing has also a negative effect on the reactiveness, as shown in Fig. 4.13, since the OFDM symbols last longer and the SS blocks sweep takes more time.

**Frequency diversity**   The repetition in frequency of multiple SS signals for the same OFDM symbol results in an increased accuracy (e.g., up to 45%, when $\lambda = 60$ gNB/km$^2$ and considering the $4 \times 4$ array configuration). The overhead is, however, from 5 to 11 times higher in our setup (according to the $\Delta_f$ used), thus there is a trade-off between the amount of resources to allocate to the users that are already connected (which is higher if frequency diversity is not used) and the opportunity to discover new users (which increases with frequency diversity for the SS blocks). However, notice that the accuracy gain reduces when increasing the array dimension (e.g., when $\lambda = 60$ gNB/km$^2$ and considering the $64 \times 4$ array configuration, a gain of just 15% is achieved, as seen from Fig. 4.10). In those circumstances, it may not be desirable to adopt a frequency diversity scheme which would inevitably increase the overhead while only providing marginal accuracy gain.

**Number of SS blocks in a burst** $N_{\mathrm{SS}}$   This parameter has a fundamental impact on the reactiveness, since a higher number of SS blocks per burst increases the probability of completing the sweep in a single burst and thus prevents $T_{\mathrm{IA}}$ from being dependent on $T_{\mathrm{SS}}$. The number of SS blocks per burst, however, increases also the overhead linearly. $N_{\mathrm{SS}}$ has a strict relationship with the number of directions to be swept, i.e., with both the beamforming architecture and the number of antennas: if, for example, hybrid or digital beamforming is used at the receiver, a larger number of antennas (i.e., narrower beams) can be supported even with a smaller $N_{\mathrm{SS}}$, as shown in Fig. 4.11

**SS burst periodicity** $T_{\mathrm{SS}}$   The periodicity of a burst has an impact on the reactiveness for initial access, since a smaller $T_{\mathrm{SS}}$ enables a larger number of opportunities in which a UE can receive synchronization signals. However, if the beam sweeping procedure is completed in a single burst, $T_{\mathrm{SS}}$ does not impact $T_{\mathrm{IA}}$ as previously defined. The overhead is inversely proportional to $T_{\mathrm{SS}}$, which has a major impact also on the reactiveness related to the tracking and the transmission of CSI-RSs, as shown in Fig. 4.14. Overall, if the sweep can be completed in a single burst, a higher $T_{\mathrm{SS}}$ would decrease the overhead and increase the reactiveness for the CSI-RSs.

**CSI-RS periodicity** $T_{\mathrm{CSI}}$   A short $T_{\mathrm{CSI}}$ allows an improved reactiveness for the beam tracking of connected users. In particular, when the number of users per gNB is high then a short CSI-RSs periodicity enables a much shorter $T_{tr}$. On the other hand, the overhead related to the CSI-RSs is small if compared with that of the SS bursts.

**Number of CSI-RSs to be monitored at the UE side** $N_{\mathrm{CSI,RX}}$   The impact of this parameter on the reactiveness is related to both the number of users per gNB and the total number of directions to be swept with the reference signals. If there is a limited number of directions and a large number of users, uniformly distributed in the available directions, then the monitoring of additional CSI-RSs does not impact $T_{tr}$ or the overhead at the network side. The UE may, however, be

impacted by the energy consumption related to the monitoring of too many directions, i.e., by a needlessly high $N_{\mathrm{CSI,RX}}$.

**gNB density** $\lambda_b$    As the network density increases, the accuracy and the average received power increase, and this allows a larger number of users to be served by a mmWave network. Besides the cost in terms of equipment and energy, a higher density has also a negative effect on the interference [47]. Moreover, the number of neighbors of each single gNB increases, and this limits the available configurations for the CSI-RSs.

**Beamforming architecture** $K_{\mathrm{BF}}$    A digital beamforming architecture at the receiver side would improve the reactiveness of the measurement scheme and decrease the overhead, without penalizing the accuracy. The same improvement in terms of reactiveness and overhead can be achieved with an omnidirectional receiver, but the accuracy would decrease with a loss of around 30% (when $\lambda = 30$ gNB/km$^2$) with respect to the $M_{gNB} = 64$ configuration, as displayed in Fig. 4.9. The complexity of the transceiver implementation and the energy consumption [238] are, however, two important parameters that must be taken into account. A hybrid configuration could represent a trade-off between an improved reactiveness and a simpler and less consuming transceiver design. Finally, notice that a digital architecture allows a higher gain with respect to the reactiveness if used at the gNB in an uplink-based framework, since the directions to be swept at the gNB are usually more than at the UE.

**Antenna Arrays** $M_{gNB}, M_{UE}$    The antenna array is one of the parameters that has the largest impact on the accuracy. A larger number of antennas enable narrower beams and higher accuracy, since the received power at the UE (in downlink) or at the gNB (in uplink) increases. The width of the beam has, however, an inverse relationship with the number of directions to scan, thus configurations that provide a higher accuracy perform worse in terms of reactiveness and overhead. Notice that the choice of the antenna array and of the beam design is strictly tied to the beamforming architecture (if digital or hybrid beamforming is used then narrower beams can be supported without penalizing reactiveness and overhead) and the configuration of the SS bursts (a large number of directions to be swept with a limited number of SS blocks per bursts has a negative impact on the reactiveness). In Fig. 4.19 a direct comparison among three different schemes is shown. It can be seen that the $M_{gNB} = 4 \times M_{UE} = 4$ configuration presents a smaller overhead and an improved reactiveness with respect to the $M_{gNB} = 64 \times M_{UE} = 16$ configuration. Moreover, both configurations with 64 antennas at the gNB have the same overhead, but there is a trade-off between the reactiveness (the configuration with the omnidirectional UE has the best reactiveness) and accuracy (using 16 antennas at the UE provides the best accuracy, at the cost of a higher energy consumption because of digital beamforming).

**Measurement Framework**    As far as initial access is concerned, the implementation of a standalone scheme generally guarantees more reactive access capabilities. The reason is that faster beam reporting operations are ensured if multiple SS blocks and RACH opportunities can be allocated within a single SS burst. On the other hand, a non-standalone framework may be preferable to: (i) reduce the impact of the overhead in the beam reporting phase; (ii) in connected mode, implement efficient and reactive recovery operations as soon as a radio link failure event is detected; (iii) guarantee a more robust control signaling exchange (e.g., when forwarding the beam reporting messages). Moreover, a non-standalone architecture is also better than an SA one when it is not possible to allocate in the same SS burst the SS blocks for the first sweep and the subsequent RACH opportunities, because for example there are too many directions to monitor

**Figure 4.19:** Comparison of three different configurations for accuracy, reactiveness and overhead. The common parameters are $\Delta_f = 120$ kHz, $N_{user} = 10$, $N_{symb,CSI} = 2$, $\rho = 0.072$, $\lambda_b = 30$ gNB/km$^2$, $N_{CSI,RX} = 3$, $T_{SS} = 20$ ms, $T_{CSI} = 0.625$ ms.

at the gNB. Finally, NSA enables a centralized beam decision: unlike in traditional attachment policies based on pathloss measurements, by leveraging on the presence of an eNB operating at sub-6 GHz frequencies, an NSA-based beam association can be performed by taking into account the instantaneous load conditions of the surrounding cells, thereby promoting fairness in the whole cellular network [163].

Overall, it is possible to identify some guidelines for the configuration of the measurement framework and the deployment of a NR network at mmWave frequencies. First, a setup of $N_{SS}$, the RACH resources, the beamforming and the antenna array architectures that allows the completion of the beam sweeping and reporting procedures in a single burst is preferable, so that it is possible to increase $T_{SS}$ (e.g., to 20 or 40 ms), and consequently allocate a larger number of CSI-RSs for the connected users (to guarantee more reactive tracking operations) and reduce the overhead of the SS blocks.

Second, the adoption of a frequency diversity scheme for the SS blocks depends on the load of the gNBs: if many users are connected to a certain gNB, this could disable the frequency diversity to both reduce the overhead and avoid discovering new users. Third, with low network density, larger antenna arrays make it possible to detect farther users, and provide a wider coverage but, as $\lambda_b$ increases, it is possible to use a configuration with wide beams for SS bursts (so that it is more likely to complete a sweep in a single burst) and narrow ones for CSI-RS, to refine the pointing directions for the data transmission and achieve higher gains.

Finally, when considering stable and dense scenarios which are marginally affected by the variability of the mmWave channel, a standalone architecture is preferable for the design of fast

initial access procedures, since it enables rapid beam reporting operations. Conversely, an NSA configuration should be employed by users in connected mode to guarantee higher resilience and an improved reactiveness in case a radio link failure occurs. A downlink configuration is in line with the 3GPP design for NR and reduces the energy consumption at the UE side (since it has just to receive the synchronization or reference signals), but is less reactive because the gNBs have a larger number of directions to sweep with downlink SS blocks or CSI-RSs.

## 4.8 Experimental Evaluation of Position-Aided Beam Management Schemes for UAVs

In the context of 5G networks, mobile base stations are considered as a way to address subscribers coverage and increasing capacity demands [239, 240]. These nomadic cells can be mounted, for instance, on UAVs, and deployed in areas where the cellular network service is unavailable (e.g., because of natural calamities) [241], to offload ground infrastructure during peak demand hours, or to scale up the network capacity during temporary events (e.g., rallies and concerts), while helping mobile operators save on the deployment costs of additional ground infrastructure. Ultimately, nomadic-cell-based solutions will shift the cellular networking paradigm toward proactive, cost-effective, and elastic resource deployment strategies.

The deployment of nomadic cells, however, demands a high-bandwidth and low-latency wireless backhaul infrastructure, as the fiber optic cables typically employed at ground-based cells are not a feasible solution for UAV-based aerial cells. On the other hand, while mmWaves point-to-point or point-to-multipoint links have been proposed as a solution for high-capacity wireless backhaul in cellular networks [80, 242], they face more challenges when applied to nomadic cells [243–245] than in the context of cellular networks discussed in the first part of this chapter, given that the high level of mobility of the endpoints of the connection requires precise and fast tracking of the best beam pair to be used for the communication [246, 247]. Consequently, the key to mobile mmWave backhaul connectivity of UAV-based aerial cells is an efficient, low-overhead beam management. As discussed in the first part of this chapter, however, most of the beam management solutions for mmWave networks lead to significant initial access delays, which increases with the number of available beams and, consequently, directions to scan. These solutions are unfit for highly mobile cells, where UAV-based aerial cells frequently relocate to satisfy the users' service demands. In these scenarios, in fact, achieving a low-latency backhaul link establishment becomes paramount to avoid the overhead that could worsen the user experience and prevent timely communications. In conclusion, there is a need for a fast and efficient beam identification procedure toward rapid backhaul link establishment and management for UAV-based aerial cells.

In this section, we propose mmBAC, a low-overhead location-aided mmWave backhaul beam management scheme for UAV-based cellular networks. We envision a beam management strategy in which the relative positioning between the UAV and the ground mmWave backhaul endpoints is leveraged to perform an efficient beam scan aimed at identifying the strongest beam path. Specifically, by using a sub-6 GHz link dedicated to the control plane, the aerial cell coordinates



**Figure 4.20:** UAV-based aerial cell with mmWave backhaul.

with the ground backhaul endpoint, removing the need for a blind scan that would inherently increase the link establishment latency. We prototyped mmBAC mounting a 60 GHz mmWave radio on a UAV, which continuously relocates to emulate different user traffic demands. Our results show that the proposed beam management scheme reduces the average latency of the mmWave link establishment by 66% with respect to a state-of-the-art iterative scan, leading to up to a 10× spectral efficiency gain in highly mobile scenarios.

The remainder of this section is organized as follows. In Sec. 4.8.1, we propose our location-aided beam management scheme. Then, we describe the mmBAC prototype in Sec. 4.8.2 and report the experimental results in Sec. 4.8.3.

### 4.8.1 Location-aided beam management

As discussed in Sec. 4.4, beam management solutions proposed for wireless cellular networks such as 3GPP NR [7] usually envision an exhaustive sequential scan of the available beam directions, searching for the optimal transmitter-receiver beam pair.

Other possible and more advanced solutions are based on a faster multi-tier scan. In this case, a small subset $S$ of all the available beam pairs $D_{tx} \times D_{rx}$ (with $D_i$ the directions available at endpoint $i$) is tested during a first initial phase, where each of these beam pairs $\mathbf{b} = [b_{tx}, b_{rx}]$ is ranked by its measured SNR $\Gamma_{\mathbf{b}}$. The best ranked beam pair $\tilde{\mathbf{b}} = \arg\max_{\mathbf{b} \in S}(\Gamma_{\mathbf{b}})$ is then used as a starting point for a more refined test involving $(2N+1)^2 - 1$ beam pairs around $\tilde{\mathbf{b}}$, with $N$ a tunable parameter.

These schemes, however, lead to very high initial access delays and do not scale well to UAV-based aerial cell deployments. In such scenarios, in fact, beam management is challenged by faster and more erratic mobility than in traditional cellular networks, and the initial link establishment latency may prevent users from accessing the network for long periods of time. Other more refined techniques have been proposed for cellular networks, but they often require significant changes in the signal processing chain, and/or in the radio hardware implementation [191]. On the other hand, works such as [192–194] have shown the benefits of leveraging side information (e.g., the relative location of the two endpoints of the link) for beam tracking, toward a faster, more reliable beam management for mobile cellular networks.

In this section we (i) present a location-aided beam tracking algorithm by exploiting the GPS coordinates of a mobile cell, and (ii) experimentally evaluate the benefits of employing side information for beam management of UAV-based aerial cells with mmWave wireless backhaul.

**Aerial cellular networks.** The mmBAC link establishment procedure starts when the UAV reaches a target location to supply service to mobile users. In order to do so, it needs to quickly establish a high-data-rate backhaul link using a 60 GHz mmWave connection with the ground radio. Given the relative locations of the ground radio $(x_g, y_g)$ and of the UAV $(x_u, y_u)$, mmBAC computes the angle between the radio on the UAV and the one on the ground, relative to a North-South axis (Fig. 4.21a) as follows :

$$\Theta_{g,u} = \tan^{-1}\left(\frac{x_g - x_u}{y_g - y_u}\right). \tag{4.19}$$

We consider the antenna arrays of the UAV facing North, and those of the ground station facing South. This can be generalized by adding an offset to $\Theta_{g,u}$ and extended to a scenario with multiple ground radios covering different sectors. mmBAC maps $\Theta_{g,u}$ to a beam index $b_{g,u}$, and defines a set of directions to be scanned at the two endpoints of the mmWave link as $D_{g,u} = \{b_{g,u} - N, \dots, b_{g,u}, \dots, b_{g,u} + N\}$. Considering that the computed angle $\Theta_{g,u}$ might not offer enough precision to identify the strongest beam path, the scan width $N$ can be tuned based

**(a)** Geometry.

**(b)** Example of beam scan patterns with baseline and mmBAC schemes.

**Figure 4.21:** mmBAC and baseline scan examples.

on the GPS localization inaccuracy and the flying stability of the UAV.

This beam tracking logic is implemented at the UAV, which coordinates with the ground station about the set of beams to scan by leveraging multi-connectivity (i.e., through a sub-6 GHz control link). In this way, the ground station transmits over the beams in $D_{g,u}$ known pilot symbols which are, then, used at the UAV for channel quality evaluation as in 3GPP NR. The radio at the UAV side measures the channel quality (e.g., SNR) for every beam combination and selects the best beam pair to use. Specifically, given the set of beam pairs $D_{g,u} \times D_{g,u}$ and the monitored SNR $\Gamma_{\mathbf{b}}$, $\mathbf{b} \in D_{g,u}$, mmBAC selects as the optimal beam pair

$$\hat{\mathbf{b}} = \underset{\mathbf{b} \in D_{g,u} \times D_{g,u}}{\arg \max} \Gamma_{\mathbf{b}}. \tag{4.20}$$

Data transmission on the backhaul link follows the link establishment phase. The experimental evaluation of this paper focuses on the latter.

By employing relative positioning information to find an initial master beam pair, the proposed algorithm significantly reduces the initial scanning overhead towards a more efficient and lightweight beam management. A qualitative example of the difference between a state-of-the-art multi-tier blind scanning procedure (which we consider as baseline) and mmBAC is shown in Fig. 4.21b. The multi-tier blind scheme starts with an initial fast scan, measuring the SNR of one out of every 25 beam pairs at regular angular intervals and identifies a master beam pair, which is used in a second phase as the starting point for a refined search. The left part of Fig. 4.21b reports the results in terms of SNR and beam pairs for the initial scan, which identifies $\tilde{\mathbf{b}} = [0,0]$ as the master beam pair. Then, as shown in the right part of Fig. 4.21b, a second, more refined scan is performed around $\tilde{\mathbf{b}} = [0,0]$ ($N = 5$), which terminates with the final selection of the best beam pair being $\hat{\mathbf{b}} = [2,0]$. On the other hand, mmBAC avoids the initial time-demanding search by identifying at a glance the master beam pair, thanks to the knowledge of the approximate relative locations of the two mmWave link endpoints. This translates into a faster best beam pair search which is highlighted at the bottom of Fig. 4.21b.

### 4.8.2 mmBAC Prototype

We prototyped mmBAC on a DJI Matrice 600 (M600) Pro UAV [248] mounting a 60 GHz mmWave Facebook Terragraph radio [249] and an Intel NUC 7i7DNKE Mini PC. The mmBAC

(a) Prototype components.

(b) Prototype diagram.

**Figure 4.22:** mmBAC prototype: hardware and components diagram.

aerial mmWave endlink prototype is shown in Fig. 4.22.

The DJI M600 Pro is a professional 6-motor UAV built for industrial applications. Its unfolded dimensions are $1.67 \times 1.52 \times 0.73$ m including propellers, arms and landing gear, it weighs 9.5 kg (including batteries) and can reach a maximum speed of 18 m/s. This UAV model supports six 4500mAh TB47S batteries that guarantee 16 minutes of hovering time at its full payload capacity of 6 kg. The DJI M600 Pro houses three Inertial Measurement Units (IMUs) sensors, employing a combination of accelerometer, gyroscope, and magnetometer; and three Global Navigation Satellite Systems (GNSS) units allowing centimeter scale localization precision. Both sensors and motors interface with the A3 Pro DJI flight controller that guarantees a stable and precise UAV navigation. The A3 unit exposes an API to an Onboard SDK that enables drone control by uploading flying missions with specific GPS waypoints, as well as the monitoring of telemetry readings such as the drone altitude and its GPS location.

The mmBAC code runs on the Intel NUC Mini PC mounted on the M600 Pro. The Intel NUC is a commercial Mini PC, whose limited dimensions ($101.60 \times 101.60 \times 25.69$ mm), light weight (0.61 kg), and good computational capabilities (Intel Core i7 processor with 32 GB RAM) make it particularly suitable to be carried on board. The Intel NUC is powered by the UAV batteries through a DC-DC step-up power supply module and interfaces the DJI A3 flight controller through a JTAG-USB cable and the 60 GHz radio through a 1 Gbps Ethernet cable. The compute board runs Ubuntu 16.04 LTS and executes the Python 3.7 implementation of the mmBAC algorithm which (i) interacts with the DJI APIs through the Onboard SDK to read the UAV location, and (ii) performs the location-aided scanning procedure described in Sec. 4.8.1.

The mmWave transceiver carried by the UAV is a Facebook Terragraph mmWave radio. This radio is optimized for working in the 60 GHz frequency band and is equipped with TX and RX arrays of $36 \times 8$ antenna elements each. Each array covers an angular space of 90° with a total of 64 beam directions. Specifically, the beams from 0 to 31 and from 32 to 63 cover the angular spaces $(0°, +45°)$ and $(0°, -45°)$, respectively. Each beam is as fine as 2.8° and the radio has an EIRP of 45 dBm. A second Terragraph mmWave radio is located on an adjustable tripod placed on the ground. The UAV-mounted and the ground radios communicate through a mmWave aerial link. We employed a Dell Latitude 3550 laptop with Ubuntu 16.04 LTS to drive the ground radio through a 1 Gbps Ethernet interface. The two controllers exchange target beam information over an out-of-band Wi-Fi channel. The mmWave radios were provided with a limited closed-source firmware (i.e., they can only be used for channel sounding) and could not be used for data transmission, which we consider outside the scope of this work. Moreover, at the current state of the firmware, the time scale at which the Facebook Terragraph radios perform channel sounding (approximately 70 ms per beam pair) is not comparable to that of commercial solutions [7]. Thus, the duration of the beam search procedures reported in Sec. 4.8.3 is longer than that of a commercial deployment.

109

**Figure 4.23:** Line experiment. The red boxes represent the overhead to find the best beam pair for different schemes.

**Figure 4.24:** Trapezoid experiment. The red boxes represent the overhead to find the best beam pair for different schemes.

### 4.8.3 Experimental Evaluation

**Experimental Setup.** We deployed one mmWave radio as fixed ground transmitter and let the UAV-mounted mmWave receiver operate in full line-of-sight conditions typical of aerial links. We operated the UAV-mounted mmWave receiver following a predefined flying mission with intermediate stopping points to mimic an aerial cell following time-varying user traffic demands. We performed two sets of experiments envisioning two different aerial trajectories, both starting from an initial waypoint located 5 m away from the ground station, mounted on an adjustable tripod. We set the tripod and the UAV hovering height to 3 m, the flying speed to 2 m/s, the hovering duration to 30 s over each waypoint, and operated in line-of-sight conditions. For each mission, we evaluated our *location-aided beam tracking algorithm* against the *fast multi-tier beam scan algorithm* described in Sec. 4.8.1 and accounted for three different metrics: (i) the total beam searching time $\Delta_s$ to find the best beam pair; (ii) the SNR $\Gamma_{\hat{b}}$ corresponding to the selected beam pair; and (iii) the spectral efficiency $S$ of the mmWave backhaul link. The latter is computed starting from the Shannon capacity equation and accounting for the overhead of the link establishment. Given the relocation interval $\Delta_h$, the spectral efficiency is

$$S = \frac{\Delta_h - \Delta_s}{\Delta_h} \log_2 \left(1 + \Gamma_{\hat{b}}\right). \tag{4.21}$$

**Line mission.** The first set of experiments concerns a UAV-based aerial cell relocating across 5 waypoints forming a line trajectory and stopping for 30 s at every waypoint throughout the mission. The nomadic cell trajectory is shown in the top part of Fig. 4.23. Upon each cell relocation, the mobile cell performs a beam tracking procedure following the two techniques described above. The average performance in terms of SNR and beam tracking overhead for the

110

**Figure 4.25:** Backhaul link spectral efficiency $S$ as a function of the relocation interval.

two algorithms are shown in Fig. 4.23. While the multi-tier iterative beam scanning algorithm has an average beam search overhead of 64.1%, the proposed location-aided beam tracking algorithm operates with an average 23.6% overhead, achieving a 72.6% higher link efficiency with a minimal SNR drop.

**Trapezoid mission.** In the second set of experiments, the UAV-based aerial cell relocates across 6 waypoints drawing a trapezoid trajectory (top part of Fig. 4.24) covering an approximate area of $58\,\mathrm{m}^2$. The cell stops for $30\,\mathrm{s}$ at each waypoint, as for the line experiment. The performance of the baseline and the location-aided beam tracking algorithm is shown in Fig. 4.24. As for the previous case, the proposed algorithm ensures higher link spectral efficiency by reducing the beam search overhead by 2.34 times ($8.1\,\mathrm{s}$ vs. $19.01\,\mathrm{s}$) while maintaining approximately the same signal quality.

**Takeaways.** The beam management experiments conducted on mmBAC highlight the importance of a fast, lightweight beam tracking solution. The proposed GPS-aided beam tracking algorithm ensures a 66% overhead reduction in link establishment compared to a state-of-the-art blind beam management scheme, while guaranteeing minimal link quality loss, which results from the strong secondary paths at the minimal height of 3 m. Reducing mmWave communication link establishment overhead leads to a higher spectral efficiency which is particularly important in mobile mmWave networks. Fig. 4.25 shows the spectral efficiency $S$ as a function of the mobile cell relocation interval $\Delta_h$. Smaller values of $\Delta_h$ represent a more dynamic network environment, where service demands rapidly change over time and call for frequent cell relocations, while large values of $\Delta_h$ represent a slowly varying service demand scenario characterized by longer relocation intervals. Fig. 4.25 highlights how location-aided beam management schemes lead to significant performance gains for highly dynamic network scenarios. For instance, the proposed algorithm outperforms the baseline by 10.7 times for $\Delta_h = 20\,\mathrm{s}$ and by 1.65 times for $\Delta_h = 30\,\mathrm{s}$. This article presents some real-world experimental results of dynamic beam-tracking algorithms, and envisions possible gains coming from side-information-aided beam management for future continuous-motion mmWave aerial links. A video demonstration of mmBAC can be found at [250].

## 4.9 Conclusions

In this chapter, we have presented a tutorial on beam management frameworks for mmWave communications in 3GPP NR and an experimental performance evaluation of a beam maganement scheme for UAVs. The harsh propagation at mmWave frequencies requires the implementation of directional transmissions supported by beamforming techniques to increase the link budget. Therefore, control procedures such as initial access must be updated to account for the lack of an omnidirectional broadcast channel, and the optimal beam pair with which a base station and

a UE communicate should be tracked when needed. Consequently, the design and configuration of efficient IA and tracking procedures is of extreme importance in cellular and UAV networks operating at mmWaves.

After a brief overview of the literature on beam management at mmWave frequencies, we described the frame structure and reference signals in 3GPP NR, focusing on the settings for communication at frequencies above 6 GHz. Then, we described several beam management procedures according to different network architectures (standalone and non-standalone) and signal transmission directions (downlink or uplink). We also evaluated the impact of several parameters (specified by 3GPP for NR) on their performance. We showed that there exist trade-offs among better detection accuracy, improved reactiveness and reduced overhead. We also provide insights and guidelines for determining the optimal initial access and tracking strategies in different network deployments, according to the need of the network operator and the specific environment in which the nodes are deployed.

Finally, we proposed mmBAC, a GPS-aided beam management algorithm for mmWave aerial links. We first prototyped a mmWave aerial link, employing a DJI M600 Pro UAV and Facebook Terragraph mmWave radios, and then evaluated mmBAC against a state-of-the-art iterative beam scanning algorithm over two sets of flying experiments, achieving an average 66% link establishment overhead reduction and up to $10\times$ higher spectral efficiency for highly mobile scenarios. Future work will focus on side-information-aided beam management schemes for mmWave aerial links in non-line-of-sight conditions employing 3D-maps of the surroundings and aerial cell relocation strategies.

# 5

# Integrated Access and Backhaul at mmWave Frequencies

## 5.1 Introduction

As highlighted in Chapter 1, operating in the mmWave spectrum comes with its own set of challenges, severe path and penetration losses being one of them [9, 42]. One promising approach to overcome such limitations is using high gain antennas to help close the link, thus introducing directionality in the communication, with electronic beamforming to support mobile users. Network densification is also used to improve the performance by reducing inter-site distance to establish stronger access channels. An ultra-dense deployment, however, involves high capital and operational expenditures (capex and opex) for network operators [251], because high capacity backhaul connections have to be provided to a larger number of cellular base stations than in networks operating at lower frequencies.

Network disaggregation (i.e., the separation of the layers of the protocol stack into different physical equipments) [27] and virtualization (i.e., the usage of software- and not hardware-based protocol stack implementations) [252] can lower capex and opex by reducing the complexity of individual base stations. Some researchers have also started investigating the feasibility of Integrated Access and Backhaul (IAB), in which only a fraction of gNB connect to traditional fiber-like infrastructures, while the others wirelessly relay the backhaul traffic, possibly through multiple hops and at mmWave frequencies [253]. The importance of the IAB framework as a cost-effective alternative to the wired backhaul has been recognized by the 3GPP. Indeed, it has recently completed a Study Item for 3GPP NR Release 16 [254], which investigates architectures, radio protocols, and physical layer aspects for sharing radio resources between access and backhaul links. Although the 3GPP LTE and LTE-Advanced standards already provide specifications for base stations with wireless backhauling capabilities, the Study Item on IAB foresees a more advanced and flexible solution, which includes the support of multi-hop communications, dynamic multiplexing of the resources, and a plug-and-play design to reduce the deployment complexity. However, despite the consensus about IAB's ability to reduce costs, designing an efficient and high-performance IAB network is still an open research challenge.

In this chapter, we present a selection of results related to the integration of IAB in 3GPP NR at mmWave frequencies. In particular, the contributions, also presented in [399, 417, 418][1]

---

[1]Part of this chapter is based on joint work with Marco Giordani.

are four-fold:

- we review the 3GPP standardization activities on IAB, with details on the Study Item (SI) [254]

- we present novel results related to the choice of the backhaul path in an IAB setup, using a mmWave channel model based on real measurements, with realistic beamforming and a sectorized deployment. We compare how different greedy policies perform with respect to the number of hops and the bottleneck SNR, i.e., the SNR of the weakest wireless backhaul link, relying only on local information, without the need for a centralized coordinator. Moreover, we discuss the use of a function that biases the link selection towards base stations with wired backhaul to the core network, and show that, for a certain set of parameters for this bias, it is possible to decrease the number of hops without affecting the average bottleneck SNR. This study can be used as a guideline for the choice and the design of backhaul path selection policies in IAB mmWave networks.

- we introduce the extension with the IAB features of the mmWave module for ns-3, which already models the mmWave channel and the PHY and MAC layers of the mmWave protocol stack, as discussed in Chapter 2. This extension can support both single- and multi-hop deployments and autonomous network configuration, and features a detailed 3GPP-like protocol stack implementation. Moreover, new scheduling mechanisms have been developed in order to support the sharing of access and backhaul resources.

- we perform the first end-to-end evaluation of IAB networks at mmWave frequencies. In particular, we compare network scenarios in which a percentage of gNBs (i.e., the IAB-nodes) use wireless backhaul connections to a few gNBs (i.e., the IAB-donors) with a wired connection to the core network against two baseline solutions, i.e., a network with only the IAB-donors, and one in which all the gNBs have a wired connection. We also investigate how to efficiently forward the backhaul traffic from the wireless IAB-nodes to the core network and demonstrate the impact of topology setup strategies on the overall throughput and latency performance. Unlike traditional performance analyses, e.g., [255, 256, 418], which are focused on PHY or MAC layer layer protocols, we also investigate the impact of upper layers, thereby providing a more comprehensive network-level analysis. Moreover, we consider both traditional User Datagram Protocol (UDP) services and more realistic applications including web browsing and Dynamic Adaptive Streaming over HTTP (DASH) for high-quality video streaming.

Our results demonstrate that, while wired backhaul implementations deliver improved overall throughput in conditions of highly saturated traffic, the IAB configuration promotes fairness for the worst users by associating to relay nodes (IAB-nodes) the UEs which otherwise would have a poor connection to the wired donor. Moreover, the performance of an IAB network depends on several factors, including the cross-layer interactions, the traffic patterns, and the attachment policies. Despite these encouraging features, real benefits of the IAB architecture and questions on network behavior under different traffic conditions have been largely left unanswered so far. Accordingly, in this chapter we study the performance of a typical network under different traffic considerations and provide insights on the observed performance gains and shortcomings under different scenarios.

The remainder of the chapter is organized as follows. In Sec. 5.2 we describe the 3GPP activities related to IAB. Then, in Sec. 5.3, we present the distributed path selection policies and the related performance evaluation. Sec. 5.4 introduces the ns-3 mmWave extension for IAB networks, with results reviewed in Sec. 5.5. Sec. 5.6 identifies the potentials and challenges of IAB networks, and Sec. 5.7 concludes the chapter.

## 5.2 Integrated Access and Backhaul in 3GPP NR

Research on wireless backhaul solutions has spanned the last two decades, with the goal of replacing costly fixed links with more flexible wireless connections. For example, mesh and multihop wireless backhaul architectures have been extensively studied for IEEE 802.11 networks [257, 258]. However, in the cellular domain, integrated solutions that provide both access and backhaul functionalities have not been widely adopted yet. There exists a relay functionality integrated in the LTE specifications, which however has not been extensively deployed due to its limited flexibility [259]: the resource configuration is fixed, it supports only single-hop relaying, and there is a fixed association between the relay and the parent base station that connects it to the wired core network. On the other hand, the wireless backhaul links that are actually used to complement fiber optic cables for backhauling traffic in sub-6 GHz cellular networks are usually custom point-to-point solutions, not integrated with the RAN.

Nonetheless, the integration of the wireless backhaul with the radio access is being considered as a promising solution for 5G cellular networks. Papers [79, 260] provided preliminary results on wireless backhaul for 5G, using also mmWave links, and showed that such solutions can meet the expected increase in mobile traffic demands. However, they did not consider a tight integration between the access and the backhaul, which is instead the focus of the more recent 3GPP SI on IAB for NR [261], recently finalized by the 3GPP. In this case, the main objective was to assess the feasibility of integrated access and wireless backhaul over NR (i.e., the 5G radio interface), and to propose potential solutions to ensure efficient backhauling operations. This Study Item led to a Work Item, and is expected to be integrated in future releases of the 3GPP specifications.

The Study Item considered fixed wireless relays with both in-band (i.e., the access and the backhaul traffic are multiplexed over the same frequency band) and out-of-band backhauling capabilities (i.e., the access and the backhaul traffic use separate frequency bands), with a focus on the former, which is more challenging in terms of network design and management but maximizes the spectrum utilization. For the in-band scenario, half-duplex relaying will be supported, although the 3GPP SI does not exclude support also for full-duplex [261]. According to [254], IAB operations are spectrum agnostic, thus the relays can be deployed in a *plung-and-play* manner either in the above-6 GHz or sub-6 GHz spectrum, and can operate both in SA (connected to the 5G core network) or NSA modes (connected to the 4G EPC). The possible topologies for an IAB network are (i) a Spanning Tree (ST), in which each IAB-node is connected to a single parent, or (ii) a Directed Acyclic Graph (DAG), in which each IAB-node may be connected to multiple upstream nodes. Moreover, IAB relays will present a higher flexibility in terms of network deployment and configuration with respect to LTE. Finally, as stated in [261, 262], 5G IAB relays will be used in both outdoor and indoor scenarios, also with multiple wireless hops, in order to extend the coverage, and should be able to reconfigure the topology autonomously in order to avoid service unavailability. Moreover, a flexible split between the access and the backhaul resources is envisioned, in order to increase the efficiency of the resource allocation.

In the following sections, we will review the main innovations introduced in [254] for the network architecture, the procedures for network management, and the resource multiplexing through scheduling, and then discuss the relevant state of the art with respect to IAB at mmWave frequencies.

### 5.2.1 Architecture

As shown in Fig. 5.1, the logical architecture of an IAB network is composed of multiple IAB-nodes, which have wireless backhauling capabilities and can serve UEs as well as other IAB-nodes,

**Figure 5.1:** Protocol stack and basic architecture of an IAB network. The Uu interface represents the interface between the UE and the DU in the IAB-node, while the F1* interface is used between the IAB DU and the upstream CU.

and IAB-donors, which have fiber connectivity towards the core network and can serve UEs and IAB-nodes.

The Study Item initially proposed five different configuration options for the architecture, with different levels of decentralization of the network functionalities and different solutions to enable backhauling. The final version, selected for future standardization, was preferred because it had limited impact on the core network specifications, had lower relay complexity and processing requirements, and had more limited signaling overhead.

According to the chosen architecture, each IAB-node hosts two NR functions: (i) a Mobile Termination (MT), used to maintain the wireless backhaul connection towards an upstream IAB-node or IAB-donor, and (ii) a DU, to provide access connection to the UEs or the downstream MTs of other IAB-nodes. The DU connects to a CU hosted by the IAB-donor by means of the NR F1* interface running over the wireless backhaul link. Therefore, in the access of IAB-nodes and donors there is a coexistence of two interfaces, i.e., the Uu interface (between the UEs and the DU of the gNBs) and the aforementioned F1* interface.

With this choice it is possible to exploit the functional split of the radio protocol stack: the CU at the IAB-donor holds all the control and upper layer functionalities, while the lower layer operations are delegated to the DUs located at the IAB-nodes. The split happens at the RLC layer, therefore RRC, Service Data Adaptation Protocol (SDAP) and PDCP layers reside in the CU, while RLC, MAC and PHY are hosted by the DUs. An additional adaptation layer is added on top of RLC, which routes the data across the IAB network topology, hence enabling the end-to-end connection between DUs and the CU.

### 5.2.2 Network Procedures and Topology Management

An important element to be considered in an IAB deployment is the establishment and management of the network topology. This is because the end-to-end performance of the overall

network strongly depends on the number of hops between the donor and the end relay, on how many relays the donor needs to support, and strategies adopted for procedures such as network formation, route selection and resource allocation. To ensure efficient IAB operations, it is necessary to optimize the performance of various network procedures involving topology and resource management.

The topology establishment is performed during the IAB-node setup, and is a critical step. When an IAB-node becomes active, it first selects the upstream node to attach to. To accomplish this, the MT performs the same initial access procedure as a UE, i.e., it makes use of the synchronization signals transmitted by the available cells (formally called synchronization signal block (SSB) in NR) to estimate the channel and select the parent. Moreover, although not currently supported by the specifications, we argue that it would be beneficial if the MT could retrieve additional information (e.g., the number of hops to reach the donor, the cell load, etc.), and then select the cell to attach to, based on more advanced path selection metrics [395] than just the Received Signal Strength (RSS), as will be discussed in Sec. 5.3. Then, the IAB-node configures its DU, establishes the F1* connection towards the CU in the remote IAB-donor, and is ready to provide services to UEs and other IAB-nodes. During this initial phase, the IAB-node may transmit information to the IAB-donor about its topological location within the IAB network.

The topology management function then dynamically adapts the IAB topology in order to maintain service continuity (e.g., when a backhaul link is degraded or lost), or for load balancing purposes (e.g., to avoid congestion). In addition to the information provided during the initial setup procedure, the IAB-nodes may also transmit periodic information about traffic load and backhaul link quality. This allows the CU to be aware of the overall IAB topology, find the optimal configuration, and adapt it by changing network connectivity (i.e., the associations between the IAB-nodes) accordingly.

In case the IAB-nodes support a DAG topology with multi-connectivity towards multiple upstream nodes, it is also possible to provide greater redundancy and load balancing. In this case, the addition/removal of redundant routes is managed by the CU based on the propagation conditions and traffic load of each wireless backhaul link.

### 5.2.3 Scheduling and Resource Multiplexing

For in-band IAB operations, the need to multiplex both the access and the backhaul traffic within the same frequency band forces half-duplex operations. This constraint has been acknowledged in the 3GPP Study Item report [254], although full-duplex solutions are not excluded. Therefore, the radio resources must be orthogonally partitioned between the access and the backhaul, either in time (Time Division Multiplexing (TDM), which is the preferred solution in [254]), frequency (Frequency Division Multiplexing (FDM)), or space (Space Division Multiplexing (SDM)), using a centralized or decentralized scheduling coordination mechanism across the IAB-nodes and the IAB-donor.

Despite the limitations imposed by the half-duplex constraint, the IAB network is required to address the access traffic requirements of all the users. For this reason, the available resources should be allocated fairly, taking into account channel measurements and topology-related information possibly exchanged between the IAB-nodes. Furthermore, both hop-by-hop and end-to-end flow control mechanisms should be provided to mitigate the risk of congestion on intermediate hops, which might arise in case of poor propagation conditions.

### 5.2.4   IAB at mmWave Frequencies

The usage of mmWave frequencies for IAB nodes introduces new opportunities and challenges. In particular, the directionality of mmWave links implies a higher spatial reuse, possibly enabling a spatial division multiple access scheme and a higher throughput, as discussed in [51]. On the other hand, the harsh propagation environment in the mmWave band requires a prompt adaptation of the topology and a fast link selection in case of outage, together with a dynamic scheduling process that adjusts the resource partition between access and backhaul according to the respective load. Therefore, mmWave IAB nodes can fully benefit from the flexibility and the self-organizing properties envisioned in the 3GPP SI for IAB.

In this regard, some papers recently analyzed the performance of IAB deployments at mmWaves, focusing however primarily on scheduling. In [263], the authors consider a centralized scheduling and routing problem, and show its performance in terms of throughput and complexity required to find the optimal solution. Similarly, paper [264] considers a joint optimization of scheduling and power, with the energy efficiency of the system as a target. In [265], the authors focus on the resource split between access and backhaul, without considering link selection for IAB nodes. None of these works, however, considers a channel characterized by the full channel matrix, with large and small scale fading phenomena, nor realistic beamforming, in the performance evaluation. In [242], the authors demonstrated that the noise-limited nature of large-bandwidth mmWave networks offer interference isolation, thereby providing an opportunity to incorporate self-backhauling in a mesh small-cell deployment without significant throughput degradation. Paper [260] showed that wireless backhaul over mmWave links can meet the expected increase in mobile traffic demands, while paper [266] evaluates the energy efficiency of mmWave backhaul at different frequencies.

Despite its clear strengths, the design of IAB solutions in mmWave systems is a research challenge that is still largely unexplored. Most of the existing literature does not consider a channel characterized by the full channel matrix, nor realistic beamforming patterns. Moreover, the prior art lacks considerations on the end-to-end performance of the self-backhauling architectures, which are in turn part of our original contributions.

## 5.3   Path Selection Policies for IAB at mmWaves

In this section, we use the NYU channel model for mmWave frequencies described in [37] to analyze the performance of different path selection policies for the backhaul. In the following paragraphs, we will use the term (i) *wired* gNB or *donor* to identify gNBs which are connected to the core network with a wired backhaul; (ii) IAB node or *relay* to label gNBs which do not have a wired backhaul link; and (iii) *parent* gNB to name a gNB which provides a wireless backhaul link to an IAB node. The parent can be itself a wireless IAB node, or a wired gNB.

For all of the policies, the IAB node that has to find the path towards the core network initiates the procedure by applying the selection policy on the first hop, and then the procedure continues iteratively at each hop until a suitable wired gNB is reached. Therefore, the strategies we evaluate are greedy, i.e., consider local information[2] to perform the hop-by-hop link selection decisions, and do not need a centralized controller. These policies can be used to re-route backhaul traffic on the fly, in case of a link failure, and to connect (possibly via multiple hops) an IAB node which is joining the network for the first time to a suitable wired gNB in an autonomous and non-coordinated fashion.

---

[2]With the exception of information related to the position and the backhaul technology, which can however be shared in advance.

**Table 5.1:** Comparison between the different link selection policies studied in this paper.

| Policy | Metric | Selection rule | Pros | Cons |
|--------|--------|----------------|------|------|
| **HQF** | SNR | Select the link with the highest SNR | High bottleneck SNR | High probability of not reaching a wired gNB |
| **WF** | SNR | Select the wired gNB, if available, otherwise apply HQF | Low number of hops | Low bottleneck SNR |
| **PA** | SNR | Select the link with the highest SNR among those with parents which are closer to a wired gNB | Low number of hops | Possible ping-pong effects |
| **MLR** | Load and Shannon rate | Select the link with the highest achievable rate | High bottleneck rate, traffic balancing | High probability of not reaching a wired gNB |

In Sec. 5.3.1, we will describe each of these strategies, while in Sec. 5.3.2 we will introduce the bias functions that we designed in order to improve the forwarding performance in terms of number of hops. Then, in Sec. 5.3.3 we will present the system model and simulation assumptions, and then review the results of the performance evaluation. Some takeaways will be given in Sec. 5.3.4.

### 5.3.1 Distributed Path Selection Policies

The considered policies differ from one another because of the metric used to measure the link quality (SNR or rate), and because of the ranking criterion of the different available links at each hop. For every policy, and at each hop, we consider an SNR threshold $\Gamma_{\text{th}}$, i.e., for the link selection, we compare only backhaul connections with an SNR $\Gamma$ higher than or equal to this threshold. If $\Gamma_{\text{th}}$ is small, then it is possible to select and compare a larger number of base stations as parent candidates, and possibly increase the probability of successfully reaching a wired gNB, at the price of a lower data rate on the bottleneck link. For the access network, $\Gamma_{\text{th}}$ is usually set to $-5$ dB [416], i.e., access links with an SNR smaller than $-5$ dB are usually considered in outage. However, this choice is not valid in a backhaul context, where the link is required to reliably forward high-data-rate traffic from the relay to its parent gNB. Therefore, we select a higher value for $\Gamma_{\text{th}}$, i.e., 5 dB, which corresponds to a theoretically achievable Shannon rate of 830 Mbps, on a single carrier with a bandwidth $B = 400$ MHz [18]. Moreover, we avoid loops, i.e., if an IAB node was used as a relay in a previous hop, it cannot be selected again.

Table 5.1 sums up the main properties of each policy, which are described in detail in the following paragraphs.

Highest-quality-first (HQF) policy    At each hop, the HQF strategy compares the SNR $\Gamma$ of the available links towards each possible parent gNBs (either wired or wireless), and selects that with the highest SNR, without considering any additional information. It is a very simple selection rule, which can be implemented only by measuring the link quality using synchronization signals. Moreover, by always selecting the best SNR, the bottleneck link, i.e., the link with the lowest SNR among the hops towards the wired gNB, will have a high SNR when compared to other policies. On the other hand, given that this policy follows a greedy approach, it may happen that

the parent gNB with the best SNR leads further away from a wired gNB, thus increasing the number of hops. Moreover, in some cases, the highest SNR leads to the choice of another relay gNB which however is not within reach of any other possible wireless parent or wired donor, thus failing to connect to a wired gNB.

**Wired-first (WF) policy** The WF policy is designed to reduce as much as possible the number of hops needed to reach a wired gNB. Indeed, if at a given hop one of the available backhaul links is toward a wired gNB, i.e., if a wired gNB is reachable from the current IAB node with an SNR higher than the threshold $\Gamma_{\text{th}}$, then the wired gNB is selected even if it is not associated to the connection with the highest SNR. If instead no wired gNB is available, then the HQF policy is applied. The IAB node would need to know which candidate parents are wired or wireless, and this can be done by extending the information directionally broadcast (using SS blocks [204]) by each gNB in the Master Information Block (MIB) or Secondary Information Block (SIB). While this policy increases the probability of reaching a wired gNB, even with a greedy approach, it may cause a degradation in the quality of the bottleneck link.

**Position-aware (PA) policy** This strategy uses additional context information related to the position of the IAB node that has to perform the link selection and the wired gNB in the scenario. This information can be available in advance and pre-configured in the relays (especially if non-mobile relays are considered [261]), or shared on directional broadcast messages. The goal is to avoid selecting a parent gNB that is more distant from the closest wired gNB than the current IAB node. Therefore, the IAB node divides the neighboring region into two half-planes, identified by a line which (i) passes through the position of the IAB node and (ii) is perpendicular to the line that passes through the positions of the IAB node and the closest wired gNB. Then, it considers for its selection only the candidate parents which are in the half-plane containing the wired gNB, and selects that with the highest SNR. This policy should strike a balance between HQF and WF.

**Maximum-local-rate (MLR) policy** The MLR policy does not consider the SNR as a metric, but at each hop selects the candidate parent with the highest achievable Shannon rate. Consider IAB node $i$, and the candidate parent $j$, and let $N_j$ be the number of users and IAB nodes currently attached to $j$. Then, given a bandwidth $B$ and the SNR $\Gamma_{i,j}$ between the IAB node and the candidate parent, the Shannon rate is computed as $R_j = B/N_j \log_2(1 + \Gamma_{i,j})$. Finally, the IAB node selects the parent with the highest achievable rate $R$. Once again, we assume that the information on the load (in terms of number of users $N_j$) of candidate parent $j$ is known to the IAB node, for example through extension of the MIB or SIB, or with a passive estimation of the power ratio between the resources allocated to synchronization signals and to data transmissions. This strategy is designed to take into account the load information in the decision, but has the same drawbacks of the HQF policy, i.e., it may yield a high number of hops and/or connection failures.

## 5.3.2 Wired Bias Function

For multi-hop scenarios, one of the Key Performance Indicators considered in the 3GPP SI for IAB is the number of hops from a certain wireless IAB node to the first wired gNB it can reach. However, as discussed in the previous section, some of the proposed policies may need a high number of hops, or even never reach the target wired gNB. In order to solve this issue, it is possible to apply a Wired Bias Function (WBF) to the SNR of the wired gNBs during the

evaluation of the metric for the link selection. Consequently, a wired gNB may be chosen as parent even though it is not the candidate with the highest considered metric.

The bias is not fixed, but is a function $W(N)$ of the number of hops $N$ traveled from the IAB node that is trying to connect to a wired gNB. The idea is that as $N$ increases, it becomes more and more convenient to select as a parent a wired gNB with respect to another wireless IAB node (that would otherwise add up to the number of hops) even though the wired gNB is not the best according to the metric considered. The WF policy is a particular case of a decision with bias, with $W(N)$ large enough so that the wired gNB is always selected if above the $\Gamma_{\text{th}}$ threshold.

We compare two different WBFs, which are respectively polynomial and exponential in the number of hops $N$. The first is defined as follows:

$$W_p(N) = \left(\frac{N}{N_{h,t}}\right)^k \Gamma_{gap} + \Gamma_H, \tag{5.1}$$

where $k$ is the degree of the polynomial, $N_{h,t}$ is a threshold on the number of hops, $\Gamma_{gap}$ a tolerable SNR gap, and $\Gamma_H$ an SNR hysteresis. The idea is that, if $N$ is smaller than $N_{h,t}$, then the SNR gap parameter $\Gamma_{gap}$ is multiplied by a number smaller than 1, and the WBF $W(N)$ does not impact too much the link choice. When the number of hops $N$ reaches the threshold $N_{h,t}$, then $W(N)$ takes values which are greater than or equal to $\Gamma_{gap}$, increasing the weight of the bias in the link selection. The SNR hysteresis $\Gamma_H$ is set to 2 dB, and slightly offsets the choice towards a wired gNB in case the best wireless relay candidate and the wired gNB have a very similar SNR. Very conservative WBF would use a large $N_{h,t}$, and small $k$ and $\Gamma_{gap}$, and vice versa for an aggressive parameter tuning.

Similarly, the exponential WBF is defined as

$$W_e(N) = \gamma^{\left(\frac{N}{N_{h,t}}\right)} \Gamma_{gap} + \Gamma_H. \tag{5.2}$$

Notice that $\gamma$ must be greater than or equal to 1, otherwise $\gamma^{\left(\frac{N}{N_{h,t}}\right)}$ would decrease with the number of hops. Moreover, for any $\gamma$, the exponential WBF $W_e(N)$ is larger than the polynomial $W_p(N)$, for the same choice of the other parameters. For example, if $N_{h,t} = 6$ and $N = 1$, with $\gamma = 1.5$ we have $\gamma^{\left(\frac{N}{N_{h,t}}\right)} = 1.07$, while with $k = 1$ we have $\left(\frac{N}{N_{h,t}}\right)^k = 0.17$.

### 5.3.3 Performance Evaluation

In this section, we first provide some details on the system model used for the performance evaluation and then discuss the simulation results and compare the different policies described in Sec. 5.3.

#### System Model

The performance evaluation for this paper is done via Monte Carlo simulations with 20000 independent repetitions for each configuration. The main parameters for the simulations are reported in Table 5.2.

The gNBs (both wired and wireless) are deployed according to a Poisson Point Process (PPP) with density $\lambda_g \in \{30, 60\}$ gNB/km², and a fraction $p_w \in \{0.1, 0.3\}$ is configured with a wired backhaul link to the core network. Therefore, the density of the wired gNBs is $\lambda_{w,g} = p_w \lambda_g$ gNB/km², while the IAB nodes have a density $\lambda_{i,g} = (1 - p_w)\lambda_g$ gNB/km². For the evaluation

**Table 5.2:** Simulation parameters.

| Parameter | Value | Description |
|---|---|---|
| $B$ | 400 MHz | Bandwidth of mmWave gNBs |
| $f_c$ | 28 GHz | mmWave carrier frequency |
| $P_{TX}$ | 30 dBm | mmWave transmission power |
| NF | 5 dB | Noise figure |
| $M$ | $\{8 \times 8, 16 \times 16\}$ | gNB UPA MIMO array size |
| $S$ | 3 | Number of sectors for each gNB |
| $\lambda_g$ | $\{30, 60\}$ gNB/km$^2$ | gNB density |
| $p_w$ | $\{0.1, 0.3\}$ | Fraction of wired gNB |

**Table 5.3:** WBF parameters.

| Configuration | Parameters |
|---|---|
| Aggressive $W_e(N)$ | $N_{h,t} = 1$, $\gamma = 3$, $\Gamma_{gap} = 15$ dB, $\Gamma_H = 2$ dB |
| Conservative $W_e(N)$ | $N_{h,t} = 6$, $\gamma = 1.5$, $\Gamma_{gap} = 5$ dB, $\Gamma_H = 2$ dB |
| Aggressive $W_p(N)$ | $N_{h,t} = 1$, $k = 3$, $\Gamma_{gap} = 15$ dB, $\Gamma_H = 2$ dB |
| Conservative $W_p(N)$ | $N_{h,t} = 6$, $k = 1$, $\Gamma_{gap} = 5$ dB, $\Gamma_H = 2$ dB |

of the MLR policy, we also deploy UEs according to a PPP with density of $\lambda_{UE}$ UE/km$^2$, and associate them to the gNB with the smallest pathloss, in line with previous studies [265].

We assume that the IAB nodes are equipped with $S$ uniform planar antenna arrays, with the same number $M \in \{64, 256\}$ of isotropic antenna elements at both endpoints of the connection. Each antenna array covers a sector of $2\pi/S$ degrees. Moreover, node $i$ can monitor the link quality of the neighboring gNB $j \in \mathcal{N}_i$, where $\mathcal{N}_i$ is the set of wired or wireless gNBs whose reference signals can be received by node $i$. The IAB node can then select the best beam to communicate with $j$ using the standard beam management procedures of 3GPP NR[3].

Table 5.3 summarizes the main parameters used for the WBF. In particular, we identify a conservative policy, with $N_{h,t} = 6$, $\Gamma_{gap} = 5$ dB and $k = 1$ or $\gamma = 1.5$ for the polynomial and the exponential policies, respectively, and an aggressive one, with $N_{h,t} = 1$, $\Gamma_{gap} = 15$ dB and $k = 3$ or $\gamma = 3$.

### Results and Discussion

The performance of the IAB path selection schemes will be evaluated by comparing the CDFs of (i) the number of hops required to forward the backhaul traffic from a wireless to a wired gNB, and (ii) the bottleneck SNR, i.e., the SNR of the weakest link.[4]

**Antenna and deployment configurations** − In Fig. 5.2 we investigate how the relaying performance evolves as a function of different setup configurations, i.e., the number of antenna elements $M$ each gNB is equipped with and the gNB density $\lambda_g$. The WF strategy is considered. As expected, increasing the MIMO array size has beneficial effects on both the number of hops and the bottleneck SNR. In the first case, the narrower beams that can be steered and the resulting higher gains that are produced by beamforming enlarge the discoverable area of each gNB, thereby increasing the probability of detecting a wired gNB with sufficiently good signal

---

[3]One of the goals of the IAB SI, indeed, is to reuse the NR specifications for the access links also for the backhaul. In any case, enhancements related to the backhaul functionality can be introduced, thanks to the more advanced capabilities of an IAB node with respect to a mobile UE [261].

[4]When considering the bottleneck SNR for the policies with WBF, we report the actual SNR, i.e., without bias.

**Figure 5.2:** Performance of the WF policy with different values of the number of antennas $M$ and gNB density $\lambda_g$.



**Figure 5.3:** Comparison of WF, HQF and PA policies, without WBF, for $M = 64$ antennas at the gNBs, $\lambda_g = 30$ gNB/km$^2$ and $p_w = 0.3$.

quality and through a limited number of hops. In the second case, sharper beams guarantee better signal quality and, consequently, stronger received power.

Similarly, enhanced backhauling performance is achieved by densifying the network since the gNBs are gradually closer and thus establish more precise alignment and, in general, connections with a higher link budget. Of course, increasing $\lambda_g$ beyond a point has a negative impact on the performance due to higher interference from the surrounding base stations.

Finally, notice that the $M = 64$, $\lambda_g = 60$ gNB/km$^2$ and the $M = 256$, $\lambda_g = 30$ gNB/km$^2$ configurations show, on average, comparable performance in terms of bottleneck SNR. However, for low SNR regimes, i.e., when considering farther nodes and more demanding signal propagation characteristics, densification is more effective than directionality.

**Path selection policies** – Fig. 5.3 compares the performance of the different path selection algorithms presented in Sec. 5.3 for different values of $p_w$, without WBF. In general, increasing $p_w$ makes it possible to minimize the number of hops required to forward the backhaul traffic from a wireless node to the core network and, at the same time, guarantees more efficient relaying operations. However, the trade-off oscillates between more robust backhauling and more expensive network deployment and management. Moreover, although the HQF policy delivers the best bottleneck SNR performance, it exhibits the worst behavior in terms of number

**(a)** Comparison of WF, and HQF policies with and without WBF.



**(b)** Comparison of WF, and PA policies with and without WBF.



**(c)** Comparison HQF policies with polynomial and exponential WBF.

**Figure 5.4:** Impact of WBF (aggressive or conservative, polynomial or exponential) on the performance.

of hops, as it greedily selects the strongest available gNB as a relay regardless of the nature (i.e., wired or wireless) of the destination node. On the other hand, both WF and PA mechanisms have the potential to reduce the number of hops since the selection is biased by the availability of the wired gNB (independent of the quality of other surrounding cells) and by context information related to the position of the wired nodes, respectively. Conversely, both approaches degrade the quality of the bottleneck link as they may end up selecting a suboptimal node among all the candidate relays within reach.

Interestingly, we observe that, when the number of available wired gNBs is very low (i.e., $p_w = 0.1$ and for low SNR regimes), the PA policy performs better than WF in terms of both number of hops and bottleneck SNR. As can be seen in Fig. 5.3, indeed, the PA policy needs a smaller number of hops than WF (and also HQF) for the paths with 4 or more hops. In low SNR and $\lambda_{w,g}$ regimes, the WF scheme asymptotically operates as HQF and, therefore, the best choice is to select the parent which is geographically closest to a wired gNB with the PA strategy.[5]

**WBF configurations** − In Fig. 5.4a we compare the behavior of the HQF and the WF policies when considering different WBF configurations to bias the path selection results. First, we see that, since the WF approach is designed to minimize the number of hops to reach a wired gNB, it generally outperforms any other architecture for the hop-count metric. However, the quality of the bottleneck link inevitably decreases (on average by more than 4 dB compared to its

---

[5]For $p_w = 0.3$ this phenomenon is obviously less pronounced but still the PA and WF paradigms reveal comparable performance in low SNR regimes.

HQF counterpart), thereby increasing the risk of communication outage between the endpoints. Moreover, for bad SNR regimes (i.e., as the probability of detecting valid wired nodes reduces) the HQF scheme implementing aggressive WBF achieves the best performance in terms of both number of hops and bottleneck SNR.

Second, we observe that, although a conservative WBF applied to an HQF scheme does not provide any significant performance improvements with respect to a pure HQF approach, a more aggressive design of the bias function has the ability to remarkably reduce the number of hops required to forward the backhaul traffic to a wired gNB, without any visible degradation in terms of SNR. We deduce that it is highly convenient to configure very aggressive[6] WBF functions since, for a multi-hop scenario, they deliver more efficient relaying operations without affecting the communication quality.

The same conclusions can be drawn by comparing the performance of the PA and the WF policies as a function of the different WBF configurations. In this regard, Fig. 5.4b illustrates how the biased PA approach guarantees very fast and high-quality backhauling thanks to the low number of hops that need to be made before successfully forwarding the traffic to the core network, and the relatively large bottleneck SNR that is experienced. In particular, the reduction in the number of hops is even beyond the capabilities of the biased HQF counterpart, at the cost of a slight reduction in the bottleneck SNR (in the order of 2 dB on the 50% percentile). Moreover, as already mentioned before, both biased and unbiased PA architectures outperform the WF scheme in the case of low SNR regimes.

Finally, in Fig. 5.4c we compare the behavior of the HQF policy with polynomial and exponential WBFs. Based on the design choices presented in Table 5.3 and according to Eqs. (5.1) and (5.2), the exponential bias function is more aggressive than the polynomial one for all values of $N$, i.e., the current number of hops. However, the exponentially-biased HQF approach, because of its inherently aggressive nature, is affected by SNR deterioration, though moderate (i.e., smaller than 1 dB on average), with respect to its polynomially-biased counterpart.

**MLR performance** − While the IAB results presented in the previous paragraphs were based on SNR considerations, i.e., the candidate parent is chosen according to the instantaneous quality of the received signal, the CDF curves displayed in Fig. 5.5 analyze the performance of the MLR backhauling approach which relies on the instantaneous cell load and the Shannon rate as a metric for the path selection operations. We observe that Fig. 5.5 leads to the same conclusions previously set out, i.e., the design of aggressive polynomial bias functions has the potential to significantly reduce the number of hops without affecting the quality of the communication (in terms of bottleneck SNR). Aggressive exponential WBFs are able to further reduce the number of relaying events, though this may slightly undermine the quality of the weakest link.

### 5.3.4 Final Considerations

Based on the above discussion, in the following we provide some guidelines on how to optimally configure the path selection policies presented in the previous sections to maximize the performance of the IAB traffic relaying operations.

We state that a WF approach, although minimizing the number of hops required to connect to a wired gNB, is affected by performance degradation in terms of bottleneck SNR. Moreover, this scheme has proven particularly inefficient when reducing the number of wired nodes (i.e., for low values of $p_w$) and for low SNR regimes (i.e., when configuring very wide beams and considering sparsely deployed networks).

---

[6]Of course, if the WBF parameters are too aggressively configured, the HQF approach will more likely operate as a WF policy, with all that this implies (including, but not limited to, a detrimental degradation of the bottleneck SNR).

**Figure 5.5:** Comparison of MLR policy with and without WBF.

In this context, a PA strategy may deliver improved performance leveraging on context information (e.g., the position of the surrounding wired gNBs) that is periodically distributed throughout the network. Furthermore, it is possible to design aggressive polynomial and exponential WBFs to bias the relay selection procedures and further reduce the overall number of hops without significant performance degradation in the quality of the weakest link.

## 5.4 IAB in ns-3 mmWave

The ns-3 mmWave module, described in Chapter 2, enables the simulation of end-to-end cellular networks at mmWave frequencies. It features a complete stack for UEs and gNBs, with a custom PHY layer, described in [110], the 3GPP mmWave channel model and, thanks to the integration with ns-3, a complete implementation of the TCP/IP protocol stack.

As mentioned in the previous sections, IAB will be important for NR ultra-dense mmWave deployments[7]. Therefore, in order to increase the realism and the modeling capabilities of the ns-3 mmWave module, we implemented an IAB framework[8] that will be described in the following sections. It features a new ns-3 `NetDevice`, the `MmWaveIabNetDevice` with a dual stack for access and backhaul, an extension of the ns-3 mmWave module schedulers, and network procedures to support IAB nodes in a simulation scenario. Moreover, we simulate the wireless relaying of both data and control plane messages, in order to accurately model the IAB operations.

An example of IAB network that can be now supported by ns-3 is shown in Fig. 5.6. In particular, we consider a tree architecture, with the root being a donor gNB, i.e., a base station with a wired connection to the core network. Therefore, this is not a traditional mesh architecture, which is used, for example, for random-access-based backhaul technologies such as IEEE 802.11 [258], in which there is no strict parent/child relationship between network nodes. In a cellular context, it is necessary to define a tree structure because every communication is scheduled [7], i.e., the base station assigns specific time and frequency resources for downlink or uplink communication with any connected UE. Therefore, given that the access and the backhaul share

---

[7]The 3GPP work item on IAB, which follows the SI, is still ongoing and is scheduled for completion as part of Release 16. We therefore do not preclude in the future to further extend the features of the ns-3 IAB module to make it fully compliant with the latests 3GPP specifications on this topic.

[8]The code can be found at `https://github.com/signetlabdei/ns3-mmwave-iab`.

**Figure 5.6:** Example of IAB architecture, with a single donor and multiple downstream IAB nodes.

the same resources, then also communication between the gNB and any IAB node must be scheduled. Notice that the connection between a parent and a child node can change with handover procedures[9], for example if the link quality between them degrades because of blockage.

In the following paragraphs we will describe the protocol stack that is deployed in each IAB node, the scheduling mechanism, and how to set up a simulation with IAB features.

### 5.4.1 IAB node

As mentioned in [261], the IAB nodes should re-use the specifications for the access stack of NR as much as possible. At the moment, there are a few protocol stacks being discussed in the 3GPP [267]. All of them, however, include PHY, MAC and RLC layers, and differ because of the support of layer-2 (i.e., RLC or PDCP) or layer-3 relaying. Given the need for a flexible solution, able to adapt to the direction that the 3GPP will take, we have implemented a light layer-3 relaying solution, i.e., each backhaul radio bearer is set up locally, and an adaptation

---

[9]We will introduce support for this functionality in the next iteration of the module.



**Figure 5.7:** Protocol stack and organization of the ns-3 classes for an IAB node.

layer above the PDCP handles the forwarding of the packets from the access to the backhaul PDCPs. Fig. 5.7 shows the protocol stack for an IAB node and the classes that model it.

The main novelties are the `MmWaveIabNetDevice` and the `EpcIabApplication` classes. The first is an extension of the ns-3 `NetDevice` class, and, similarly to the `NetDevice` implementations of the UE and gNB, holds pointers to all the objects that model the other layers of the protocol stack. Moreover, it is internally used in the ns-3 model to forward packets between an instance of the `EpcUeNas` class in the backhaul stack and the `EpcIabApplication` in the access stack.

The `EpcIabApplication`, instead, implements the main logic related to the control and data plane management in the IAB node. In particular, for the data plane, the `EpcIabApplication` class is in charge of applying the forwarding rules for local UEs, i.e., those directly connected to the IAB node this class belongs to, and for remote UEs, i.e., those connected to downstream IAB nodes. In this case, the traffic will be forwarded to the local bearer mapped to the downstream IAB device. More details on how the routing is performed will be given in Sec. 5.4.2. This class is also responsible for the processing and forwarding of control packets for the interfaces toward the core network and the other neighboring gNBs. When a control message is received on either the access or the backhaul interface, the `EpcIabApplication` checks if it is a local message, i.e., if the destination is the RRC layer of the current IAB node, and, if this is the case, forwards the packet to the RRC. Otherwise, as done in the data plane, the packets are relayed via one of the downstream IAB nodes.

The other classes are the same as those used in the UE protocol stack (for the backhaul) and gNB protocol stack (for the access). The consequence is that, in the access, the UEs in the scenario consider the IAB node as a normal gNB, and, similarly, in the backhaul, the parent gNBs and/or IAB nodes consider the IAB child as a UE. Therefore, there is no need to adapt the UE and gNB ns-3 implementations to support the IAB feature. The only change is the extension of the gNB schedulers, to support the multiplexing of access and backhaul in the same resources, and the introduction of a new interface between the access and backhaul MAC layers. These extensions will be described in Sec. 5.4.3. Nonetheless, additional enhancements can be introduced in future releases, to improve the overall performance of the IAB protocol stack and track the 3GPP SI and specifications on IAB.

### 5.4.2 Single- and multi-hop control procedures

Given that the 3GPP is still considering IAB as an SI, there are no standard specifications yet on control procedures to support IAB networks. Nonetheless, the SI [261] specifies that both single- and multi-hop topologies should be considered, and that the IAB node should be able to autonomously connect to the network, adapt the access and backhaul resource partitioning and, eventually, independently update the parent node in case of blockage. All these features require specific control procedures, and, given the high level of detail of the ns-3 model, we implemented a number of realistic control procedures, which involve an exchange of messages on the wireless backhaul links to set up and automatically configure the network. These can be easily updated to implement different procedures that the 3GPP may specify in the future.

In particular, we assume that the parent IAB node for a backhaul link terminates the NG control interface to the core network (i.e., the NR equivalent of the LTE S1 interface) [25], and that it takes care of forwarding the control messages towards the network servers that host the Access and Mobility Management Function (AMF). Moreover, the IAB node has a similar role with respect to the UEs connected to it, as would happen with a traditional wired gNB. Thanks to this design, the differences with respect to the 3GPP specifications for the access stack are minimized. This configuration makes it possible to seamlessly support both single- and multi-hop deployments, given that the architecture of the upstream portion of the network is transparent

to each IAB node, which will simply relay all of its packets to the parent. Furthermore, for the purpose of packet transport in the backhaul network, we exploit GPRS Tunneling Protocol (GTP) tunnels from each IAB node to the relevant element in the core network (i.e., the server with control functions or the packet gateway). Each data bearer of all the UEs (and IAB nodes, for the backhaul part) is associated with a unique tunneling ID, and all the packets sent on backhaul links will be associated with a GTP header carrying that ID.

We also implemented realistic autonomous access and configuration procedures for the IAB nodes. When the IAB selects its parent node during the IA procedure, the parent sends an initial message to the AMF, which will reply with the configuration for the backhaul bearer between the IAB node and its parent. These messages will be relayed by all the IAB nodes in the path between the parent and the donor gNB, and each of them will register the presence of an additional downstream IAB device. Notice that there may be multiple IAB children for each parent, therefore the parent has to match the new downstream node to the correct child to correctly route the other control and data packets.

For the UEs, there is no difference between a wireless relay and a gNB with a wired connection to the core network. Therefore, the UE's IA procedure does not change, and the IAB node will take care of forwarding the relevant control messages to the AMF and the other network functions involved in the IA. Moreover, the upstream relays and the donor gNB will exploit the control messages for the UE's IA to associate to each IAB bearer the total number of downstream UEs. For example, by considering the deployment in Fig. 5.6, if UEs 1 and 2 connect to IAB node 3, and UE 3 connects to IAB node 4, then IAB node 1 will know that the backhaul bearer towards IAB node 2 will carry the traffic for 2 UEs, and that towards node 4 will account for a single UE. This information could be exploited by advanced IAB MAC schedulers. Finally, during the UE IA procedure, each gNB associates the GTP tunneling ID of the bearers of downstream UEs to a local IAB child, so that, when a backhaul packet is received, the gNB uses the information in the GTP header to correctly route the packet.

### 5.4.3 Backhaul-aware dynamic scheduler

The MAC and the associated scheduler are a key component in the design of scheduled wireless relay architectures in which the resources between the access and the backhaul are shared. In order to avoid self-interference between access and backhaul, indeed, there is a need to multiplex the two interfaces. In our implementation, we consider TDMA, but we plan to extend the support to spatial division multiplexing in future releases, to harness the directionality of mmWave communications. Moreover, the scheduler is usually not part of the 3GPP specifications, and, therefore, equipment vendors have the possibility of designing custom solutions in this domain.

We opted for a distributed scheduling solution, in order to minimize the difference in the scheduling mechanism with respect to a traditional access-only scenario, and to limit the amount of control overhead that a centralized solution would require. Therefore, in the ns-3 mmWave IAB module, each gNB (either wired or wireless) schedules the resources for its access interface (i.e., for both UEs and IAB children) independently of the other gNBs, as would happen in a traditional network without IAB. In a TDMA setup, however, the IAB node cannot schedule resources in the time and frequency slots already allocated to the backhaul by their parent. Therefore, if at time $t$ the relay has to perform a scheduling decision for subframe $t + \eta$, then it has to be already aware of the scheduling decision of its parent for $t + \eta$. Given a delay $\epsilon$ for the communication of scheduling information between the parent and the current relay, then the parent should perform its scheduling decisions for $t + \eta$ at time $t - \epsilon$.

In order to efficiently address this issue, we implemented a *look-ahead backhaul-aware scheduling* mechanism. The backhaul-aware component is given by a new interface between the access

**Figure 5.8:** Example of resource allocation for time $T$ with a look-ahead scheduler at the donor gNB and IAB nodes 1, 2 and 3 in the deployment of Fig. 5.6.

and the backhaul MAC layers. The backhaul MAC layer is seen as a UE by the parent node, and thus will receive DCIs with the scheduling and modulation and coding scheme information for $\eta$ subframes in advance. Then, the backhaul MAC shares DCI with the scheduler of the IAB node (in the access stack), which registers the resources occupied by backhaul transmissions for the relevant subframe.

The look-ahead mechanism, additionally, makes it possible to adjust the value of $\eta$ according to the maximum number of downstream relaying hops $N$ from the current gNB to the farthest IAB node: the gNB schedules ahead by $\eta = N+1$ subframes[10], and propagates this information with a DCI to the UEs and IAB nodes connected to it. In turn, these IAB nodes will schedule ahead by at most $\eta = N$ subframes. Each of them will consider the time and frequency resources allocated for their downlink or uplink backhaul transmission as busy, and will schedule access resources for their UEs and, eventually, for IAB nodes in unallocated resources. For example, by considering Fig. 5.6, the maximum number of hops from the donor gNB is 3. Therefore, the donor will schedule ahead by 4 subframes. On the other hand, IAB node 2 has a single hop to the farthest relay, thus it will schedule ahead by 2 subframes. Notice that, thanks to the procedures introduced in Sec. 5.4.2, there is no need to manually tune the $\eta$ parameter, which is automatically configured according to the structure of the IAB tree, and can be updated in case of variations in the architecture of the network.

We added the look-ahead and backhaul-aware capabilities to two of the ns-3 mmWave module schedulers, i.e., the `MmWaveFlexTtiMacScheduler` class, which models a RR scheduler, and the `MmWaveFlexTtiPfMacScheduler` class, which implements a PF scheduling algorithm. Moreover, in a TDMA setup, with shared resources between the access and the backhaul, it is important to make sure that the parent gNB does not schedule all of the available resources to a single IAB node (e.g., if it is the only active terminal connected to the parent). Otherwise, the child IAB node would not be able to allocate any resource to the access. Therefore, we limit the maximum number of time and frequency resources that can be allocated to an IAB device to half of the total available resources.

An example of resource allocation is shown in Fig. 5.8, where a total number of 10 time and frequency resources are dynamically allocated to access and backhaul links. In particular, we refer to the deployment in Fig. 5.6, and present a possible resource partitioning for the donor

---

[10]The additional subframe with respect to $N$ is needed because the farthest IAB node (without IAB children) has to schedule its resources at least one subframe in advance, in order to transmit the DCI beforehand to its UEs

gNB, IAB nodes 1, 2 and 3 and the UEs connected to these gNBs. As can be seen, each IAB node does not allocate access transmission in the resources reserved for its backhaul, but can exploit all of the other resources for communication with other relays and the UEs, including those allocated by one of the upstream nodes to other backhaul links. While in general this may increase the interference, it must be noticed that, at mmWave frequencies, the large antenna arrays that can be built and the resulting directional transmissions that can be established have the potential to provide increased spatial reuse and isolation, thereby guaranteeing reduced interference [226]. Moreover, interference-aware schedulers can be designed and tested with the simulator. Finally, it is possible to update the allocation on-the-fly, to dynamically adapt to changed channel conditions and traffic requirements from the different connected terminals.

Fig. 5.8, however, also highlights one of the main bottlenecks of an IAB architecture, i.e., the fact that the donor gNB needs to serve not only its own users, but also all the downstream relays, carrying traffic from many other UEs. On one hand, the amount of data that can be exchanged on a backhaul link in each time and frequency resource is generally higher than the equivalent for a gNB-UE link, thus the backhaul will probably require fewer resources. Indeed, the backhaul link between two gNBs has usually a better quality than that between a gNB and a UE, given that a backhaul link is expected to be in LOS, and that a larger number of antennas can be deployed in a relay than in a UE. On the other hand, the scalability of an IAB deployment has an intrinsic limitation due to the resource sharing between the access and the backhaul link. Therefore, efficient scheduling algorithms will be key for high-performance IAB networks. This makes the ns-3 mmWave module with the IAB integration a valuable platform for researchers interested in IAB networks, given that it offers a lean interface to the scheduler implementations, which can be easily extended to test new IAB scheduling paradigms in realistic end-to-end scenarios.

### 5.4.4  Simulation setup

The setup of a simulation with the IAB feature resembles that of a simulation with traditional wired-only backhaul. An extensive description of how to configure an ns-3 simulation script for the mmWave module is provided in the tutorial in [390]. We added two auxiliary methods in the `MmWaveHelper` class, which hides from the ns-3 user much of the complexity related to the configuration of the mmWave RAN and core network. Similarly to the methods used to set up UEs and gNBs, the `InstallIabDevice` method returns a `NetDevice` properly configured, with the stack described in Fig. 5.7.

The initial attachment of each IAB node to its parent gNB is performed by the methods `AttachIabToClosestWiredEnb` or `AttachIabToBestNodeHQF`. The latter scans the signal quality of the available IAB nodes or wired donors, and selects that with the highest SNR. Moreover, it avoids the creation of loops in the network tree. These helper methods, moreover, automatically register the new IAB nodes to the control entities in the core network, and define the default radio bearer that will be used for the backhaul link. Finally, by default, the UEs in the ns-3 mmWave module perform the initial attachment as soon as the simulation starts, i.e., at simulation time $t_s = 0$. Therefore, we added the `AttachToClosestEnbWithDelay` method that delays by $D$ seconds the initial attachment of UEs to the chosen gNBs, either wired or wireless. This method can be used to let the UEs perform IA only after the IAB nodes have completed their IA and backhaul bearer setup.

| Parameter | Value |
|---|---|
| mmWave carrier frequency | 28 GHz |
| mmWave bandwidth | 1 GHz |
| 3GPP Channel Scenario | Urban Micro |
| mmWave max PHY rate | 3.2 Gbps |
| MAC scheduler | Round Robin |
| Subframe duration | 1 ms |
| Donor gNB to remote server latency | 11 ms |
| RLC buffer size $B_{RLC}$ for UEs | 10 MB |
| RLC buffer size $B_{RLC}$ for IAB nodes | 40 MB |
| RLC AM reordering timer | 2 ms |
| UDP rate $R$ | $\{28, 224\}$ Mbps |
| UDP packet size | 1400 byte |
| Number of independent simulation runs | 50 |

**Table 5.4:** Simulation parameters

## 5.5 End-to-end Evaluation of IAB

### 5.5.1 Single-Hop Scenario

In this section, we validate the implementation of the IAB features for the ns-3 mmWave module through simulations in a single-hop scenario. We illustrate some preliminary results related to the coverage performance of an IAB deployment in a Manhattan grid, with blocks of 50 m for each side, and with 10 m between each block, for a total area of 0.053 km². A gNB with a wired connection to the core network is placed at the center of the scenario, while the number of IAB nodes (i.e., gNBs with wireless backhaul functionalities) varies from 0 to 4. The IAB nodes are one block in each direction away from the donor (i.e., at a distance of 85 m), and they are in LOS (e.g., placed on the building rooftops). Each relay directly connects to the wired donor wirelessly, thus this scenario only considers single-hop transmissions. 40 users are randomly placed outdoors using the new ns-3 `OutdoorPositionAllocator` method, and connect to the closest gNB, either wired or wireless. Each UE downloads content from a remote server at a constant rate $R = \{28, 224\}$ Mbps using UDP as the transport protocol. These two different source rates are used to test the network in different congestion conditions. Finally, the MAC layer performs HARQ retransmissions, and the RLC layer uses the AM to provide additional reliability. The scheduler is Round Robin, with the look-ahead backhaul-aware mechanisms described in Sec. 5.4.3. The other simulation parameters are in Table 5.4.

We consider two different end-to-end metrics, i.e., the experienced throughput and the application-layer latency averaged over multiple independent runs. Fig. 5.9 investigates three different throughput values for different source rates $R$ and varying the number of IAB relays. We observe that, for the low source rate scenario (i.e., $R = 28$ Mbps), the total throughput remains almost constant, while, in the congested scenario (i.e., $R = 224$ Mbps) the rate progressively increases with the number of relays. This shows that, in the considered Manhattan scenario, the relays extend the area in which the mobile terminals can benefit from the coverage of their serving infrastructures and, in particular, have the potential to improve the quality of the access link between the cell-edge users and the donor gNB, thereby guaranteeing higher capacity.

The average latency is shown in Fig. 5.10. We see that, in a Manhattan grid scenario, the average latency of the UEs directly connected to the wired gNB decreases as a result of increasing the number of wireless relays. Indeed, if the relays are used, the wired gNB will serve fewer users, i.e., those with the best channel quality, and will avoid allocating resources to cell-edge users which, generally, require a high number of HARQ and RLC retransmissions. Although these

**Figure 5.9:** Sum end-to-end throughput for different source rate $R$ and number of relays. The total throughput is the sum of the throughput of all the users, while the wired (or IAB nodes) sum throughput refers to the aggregate throughput of UEs connected to the donor (or the relays, respectively).



**Figure 5.10:** Average end-to-end latency for different source rate $R$ and number of relays. We report the average latency considering all the UEs, or only those connected to the wired gNB or wireless relays. The dotted black line represents the average latency of the configuration with 0 relays.

benefits are particularly evident in the $R = 224$ Mbps case, a latency improvement is also observed for the non-congested scenario (i.e., $R = 28$ Mbps) when four relays are deployed.

On the other hand, from Fig. 5.10 we notice that the average latency of the users attached to IAB nodes increases with respect to the configuration without relays, especially when just one or two wireless relays are deployed. This is mainly due to the buffering that occurs in the backhaul. In an IAB context, indeed, the backhaul and access resources are shared, thus the IAB nodes and the UEs attached to the donor contend for the same resources. With an RR scheduler, a similar number of transmission opportunities is allocated to the IAB nodes and to the UEs, but the relays generally have more data to transmit than each single UEs. Consequently, the buffering latency at the RLC layer of the relays increases. Nonetheless, for the congested scenario (i.e., $R = 224$ Mbps), the overall average latency when more than three relays are deployed (i.e., 287 and 250 ms for three and four relays, respectively) is equivalent or lower than that in the configuration with the donor gNB only (i.e., 292 ms), as shown in Fig. 5.10.

The above discussion exemplifies how an IAB architecture introduces both opportunities and challenges. From one side, the deployment of wireless relays is a viable approach to increase

**(a)** Throughput.

**(b)** Latency.

**Figure 5.11:** Throughput and latency comparison of IAB path selection policies varying the percentage of IAB-donors $p$ for a density of 45 gNB/km$^2$ and a constant bitrate traffic.

the coverage of cell-edge users, i.e., the most resource-constrained network entities, thereby promoting fairness in the whole network. Moreover, the presence of the wireless backhaul nodes has the potential to reduce the communication latency in case of congested networks. From the other side, the IAB nodes may degrade the throughput and latency performance of some UEs, i.e., those with the best channel quality, whose traffic would have been successfully handled even in traditional wired backhaul scenarios. It becomes therefore fundamental to determine the optimal number of wireless backhaul nodes to be deployed and to design efficient scheduling policies, according to the context and considering the constraints imposed by the available network and economic resources. This research challenge will be part of our future work.

### 5.5.2 Multi-Hop Scenario

We have also evaluated the end-to-end performance of an IAB mmWave network in a generic multi-hop, end-to-end scenarios with the TCP/IP stack and realistic applications, such as the 3GPP HyperText Transfer Protocol (HTTP) model. In the scenario we investigated, the base stations are deployed following a PPP with density $\lambda$ BS/km$^2$, and a fraction $0 \leq p \leq 1$ of the $N$ base stations have wired backhaul connections (i.e., the IAB-donors), while the others (i.e., the IAB-nodes) are wirelessly connected to the IAB-donors, perhaps over multiple hops. The network implements in-band backhaul, at 28 GHz, with TDM of the radio resources among the access and the backhaul links. We consider uniform rectangular antenna arrays in the base stations and UEs, with 64 and 16 elements, respectively, and the beamforming model described in [390]. The base stations use the backhaul-aware round robin scheduler presented in [417]. The UEs are also deployed with a PPP with density $\lambda_u = 10\lambda$ UE/km$^2$, although we only evaluate the performance of the subset of users connected to a target base station, which is either the first gNB deployed in a baseline scenario in which all nodes have a wired connection to the core network, or the first IAB-node that performs the initial access in an IAB scenario.

**Backhaul path selection policies.** The first set of results, reported in Fig. 5.11, sheds light on the impact of different backhaul path selection policies in an IAB setup. As introduced in Sec. 5.2.2, path selection refers to the procedure by which IAB-nodes find the path towards an IAB-donor, possibly through multiple hops. In our previous work [418], we investigated two different policies to forward the backhaul traffic: (i) a *HQF* approach which selects, as a parent, the gNB with the highest quality, i.e., the highest SNR, and (ii) a *WF* approach which selects a direct link to the IAB-donor with the best signal, even if an IAB-node with better

134

channel quality is available, provided that some minimum channel quality criterion is satisfied. The first approach facilitates a best-quality wireless backhaul connection in the first hop but, in turn, may increase the number of hops required to forward the traffic to an IAB-donor. The second approach, while minimizing the number of end-to-end hops, may choose backhaul links with poorer channel quality. The HQF policy may also leverage a function that biases the link selection towards gNBs with wired backhaul to decrease the number of hops to the core network. The bias computed by the function is not fixed, but depends on the number of hops from the IAB-node to the candidate parent it is trying to connect to [418]. Moreover, both conservative and aggressive bias functions can be designed (aggressive HQF policies will progressively operate like WF policies). Fig. 5.11 demonstrates that the WF approach should be preferred since it offers lower end-to-end latency and higher total throughput compared to the other investigated policies. The results show that minimizing the number of hops required to connect to an IAB-donor improves throughput and reduces latency by reducing relaying overhead and congestion at intermediate IAB-nodes.

**IAB deployment scenarios.** We also tested three different deployment scenarios. The best case is when all the $N$ base stations in the network are equipped with a wired connection to the core network (i.e., the *all wired* scenario). This represents the most expensive solution, in terms of density of fiber drops, but permits the whole bandwidth to be used for access traffic. With the *IAB-nodes* option, $pN$ base stations are IAB-donors, i.e., have a wired connection and $(1-p)N$ have wireless backhaul. Finally, the baseline is the one that 3GPP considers for comparisons with IAB solutions, described in [254], i.e., a deployment with only $pN$ wired base stations and no IAB-nodes (the *only donors* configuration).

*UDP user traffic.* In Fig. 5.12, we consider an IAB network where each user downloads content from a remote server with a constant bitrate of 220 Mbps, using UDP as the transport protocol, thus introducing a full buffer source traffic model. The flow of each end-to-end connection does not self-regulate to the actual network conditions, thus congestion arises. This experiment aims to test the performance of an IAB setup in a saturation regime, where the access and backhaul links are constantly used. As expected, the best performance is provided by the all wired configuration, given that it provides the same access point density as the IAB setup, but avoids the multiplexing of resources between access and backhaul. On the other hand, it is possible to identify two advantages and one drawback of the IAB configuration with respect to



(a) Fifth percentile throughput.

(b) Throughput CDF for $p = 0.3$.

**Figure 5.12:** Fifth percentile and CDF of the throughput for the users of a target IAB-node, varying the percentage of IAB-donors $p$ and the deployment strategies, for a density of 45 gNB/km$^2$.

**(a)** Average rebuffering for DASH clients.



**(b)** Average delay to retrieve an HTTP web page.

**Figure 5.13:** Performance for users in a target IAB-node, with different applications, for a density of 30 gNB/km$^2$.

the only donors one. A higher throughput for the worst users is achieved when using IAB-nodes, as shown by the fifth percentile throughput plot in Fig. 5.12a. In particular, for $p = 0.5$ (i.e., when the number of relays is equal to the number of IAB-donors), IAB has only 13% less fifth percentile throughput than the all wired configuration. Moreover, the usage of IAB-nodes likely offloads the worst users from the IAB-donors, and this frees up resources for users with the best IAB-donor channel quality, thereby enabling a higher throughput, as illustrated in Fig. 5.12b. The IAB solution, however, requires multiplexing of the wireless resources between access and backhaul. In a scenario where the links are always saturated, this results in a worse performance for the average users connected to the relays, which are throttled on the backhaul links by the round robin scheduler at the donors and have a smaller throughput than with the only donors setup.

*DASH, HTTP user traffic.* The next set of results considers a more common use case, in which the users either stream video using DASH [54] or access web pages using HTTP from a remote server. This kind of source traffic is asynchronous and bursty, and, in the DASH case, the flow adapts itself to the varying capacity offered by the network, after some delays due to the signaling and convergence of the algorithm. Therefore, the network is not as stressed as in the previous experiment, and in this case the advantage of IAB is more visible. Indeed, thanks to the better channel seen on average by the user due to more numerous nodes compared to the only donor case, and thanks to the asynchronous and independent nature of the traffic at each user, which provides greater multiplexing gains, the performance of the IAB network is not far from that of the network with all wired access points. In particular, Fig. 5.13a reports the average duration of a rebuffering event for a DASH stream, for all the users in a target base station. The rebuffering happens when the DASH framework does not adapt fast enough to the network conditions, or if the network capacity is not sufficient to sustain even the minimum video quality available in the DASH remote server. As can be seen, the only donors setup has the worst performance, with a 5 and 2 times higher rebuffering than the all wired configuration, for $p = 0.3$ and 0.5, respectively. The IAB deployment, instead, degrades the performance of the all wired only by 1.4 and 1.3 times, for $p = 0.3$ and 0.5, respectively. Likewise, Fig. 5.13b shows the average time it takes to completely download a web page, from the first client HTTP request to the reception of the last object, and, as can be seen, the trend is similar to that of the DASH rebuffering. Finally, for this kind of traffic, the improvement introduced by the densification of IAB-donors (i.e., by increasing $p$ from 0.3 to 0.5) is less marked than with the constant bitrate

traffic shown in Fig. 5.12.

## 5.6  Potentials and Challenges of IAB

As highlighted by the results presented in Sec. 5.5, IAB networks present both benefits and limitations with respect to deployments where the radio resources are not multiplexed between the access and the backhaul. First, the IAB solution may present lower deployment costs and complexity with respect to the all wired setup, but, at the same time, splitting the available resources between access and backhaul traffic makes the overall network performance worse than in the all wired case under heavily loaded network scenarios. However, for bursty traffic the performance of the IAB solution approaches that of the all wired case. This shows that when evaluating the performance of IAB networks it is important to consider the specific use case and end-to-end applications that run on top of the network. Moreover, the results suggest that the main advantages of an IAB deployment, when compared to the only donors setup, come from an improvement in channel quality for cell edge users, on average, which consequently improves the area spectral efficiency.

On the other hand, the deployment of an IAB network presents challenges related to the design and interactions at different layers of the protocol stack. An important issue is related to the enforcement of QoS guarantees in single and multi hop scenarios, so that mixed IAB traffic flows for end-to-end applications can safely coexist. Additionally, the resources in the IAB network are limited and shared between the access and the backhaul. Therefore, the admission and bearer configuration should take this into consideration, in order to avoid overbooking the available resources and introducing congestion in the network. As shown in Fig. 5.12b, this may indeed worsen the experience of the average users. Similarly, during the setup phase, in which the IAB-nodes join the network by performing initial access to their IAB parents, it is important to consider the attachment strategies to avoid overloading some IAB-donors, or excessively increasing the number of hops. Even though we demonstrated in Fig. 5.11 that reducing the number of relay operations is always beneficial in terms of both end-to-end latency and throughput, how to design path selection strategies which are robust to network topology changes and end terminals' mobility is still an open research challenge which deserves further investigation.

Most of these system-level challenges are strictly related to the design of ad hoc scheduling procedures at the MAC layer, able to efficiently split the resources between the access and the backhaul and provide interference management. Another important challenge is related to cross-layer effects emerging from retransmissions at multiple layers, and the configuration of RLC and transport layer timers may need to account for the additional delays related to the retransmissions over multiple hops and the reordering of packets at the receiver. At the physical layer, it will be interesting to evaluate the gain of the spatial multiplexing of the access and the backhaul, by using digital or hybrid beamforming, which could avoid the time or frequency multiplexing that are needed when using single-beam analog beamforming.

Overall, these challenges represent promising research directions to enable self-configuring, easy-to-deploy and high-performing IAB networks, which could represent a cost-effective solution for an initial ultra-dense NR deployment at mmWave frequencies.

## 5.7  Conclusions and Future Work

High-density deployments of 5G cells operating at mmWaves call for innovative solutions to reduce capital and operating costs without degrading the end-to-end network performance. In

this context, IAB has been investigated as an approach to relay access traffic to the core network wirelessly, thereby removing the need for all base stations to be equipped with fiber backhaul.

In this chapter we have reviewed the characteristics of IAB capabilities that are currently being standardized in 3GPP NR Release 16, and evaluated the performance of different distributed path selection strategies to efficiently forward the backhaul traffic (possibly through multiple hops) from a wireless gNB to a wired gNB connected to the core network. The investigated policies may or may not leverage a function that biases the link selection towards base stations with wired backhaul capabilities, to minimize the latency of the relaying operations. We have shown through simulations that it is always possible to decrease the number of hops required to connect to a wired gNB by designing aggressive bias functions without affecting the average bottleneck SNR (i.e., the quality of the weakest link).

Additionally, we have presented the first implementation of IAB for the ns-3 mmWave module. The simulator, which features the 3GPP mmWave channel model and a complete characterization of the TCP/IP protocol stack, now also implements the wireless relaying operations on both the data and the control planes, thereby accurately modeling the operations of an IAB network. We believe that this tool can be used by researchers to understand the main limitations and the performance gains that IAB networks can provide, and to evaluate new integrated scheduling algorithms and multi-hop routing strategies with a realistic, end-to-end protocol stack.

By using this tool, we evaluated IAB networks for different applications and traffic types such as Internet browsing (i.e., HTTP) and video streaming (i.e., DASH). We showed that IAB represents a viable solution to efficiently relay cell-edge traffic, although the benefits decrease for more congested networks. We have also highlighted the limitations of the IAB paradigm and provided guidelines on how to overcome them.

IAB standardization is, however, still an on-going process. As part of our future work, we will validate wireless backhaul solutions considering recently proposed 3GPP scenarios and investigate the impact of mobility and network reconfiguration on the network performance. Moreover, we will further extend the ns-3 mmWave module with additional IAB features, in order to address mobility scenarios, and keep track of the 3GPP specifications on this topic.

# Part III

# The Protocols: End-to-End and Cross-Layer Analysis of 5G mmWave Networks

# 6

# TCP performance in mmWave Networks

## 6.1 Introduction

End-to-end connectivity over the internet largely relies on transport protocols that operate above the network layer and are in charge of delivering packets between remote nodes. The most widely used transport protocol is the Transmission Control Protocol (TCP), designed in the 1980s [268] to offer reliable packet delivery and sending rate control to prevent congestion in the network. Reliability is accomplished with receiver's acknowledgments (ACKs) fed back to the sender, which retransmits packets if needed, while rate control is achieved by dynamically adjusting the congestion window, i.e., the maximum amount of data that the sender can transmit without receiving ACKs. Several Congestion Control (CC) algorithms have been proposed in order to improve the goodput (defined as the application layer throughput) and latency of TCP over different types of networks [269].

In wireless networks, however, the loss of a packet is not necessarily caused by congestion, but may instead be due to a sudden (and possibly only temporary) drop in signal quality. In [138,270], the authors study the behavior of TCP in relation to a complex mobile network such as LTE, showing (i) that as the distance between the User Equipment (UE) and the evolved Node Base (eNB) increases the TCP throughput degrades, and (ii) how TCP is affected by network events such as handovers.

However, the next generation of cellular networks will present new challenges for TCP and, in general, for transport protocols that implement congestion control mechanisms based on an abstract view of the end-to-end network (e.g., Quick UDP Internet Connections (QUIC) [271], Stream Control Transmission Protocol (SCTP) [272]). In particular, these issues are specifically related to the presence of mmWave links in the radio access network, which, as discussed in the previous chapters, exhibit an erratic propagation behavior. This technology is seen as a promising enabler for the 5G targets of multi-gigabit/s data rates and ultra-low latency [273], but the end-to-end performance perceived by the user will eventually depend on the interaction with transport protocols such as TCP.

An example of the difference in propagation conditions between an LTE (at 2.1 GHz) and a mmWave (at 28 GHz) system is shown in Fig. 6.1, for a UE that moves at 2 m/s at an average distance of 75 meters from the eNB, and switches from a LOS to a NLOS condition. The main differences between the mmWave and the LTE channels are that:

- the LOS to NLOS pathloss transitions are deeper for mmWave. At sub 6 GHz frequencies,

**Figure 6.1:** ns–3 simulated SINR for a mmWave and LTE link, with the UE moving at 2 m/s from (45, 0) to (55,0), with the eNB at coordinates (75, 50). The traces are generated using the channel model described in [42, 388] for the mmWave channel and the ns–3 LTE channel model.

these are of the order of 10-15 dB, while for the mmWave channel the fluctuation can exceed 30 dB. Therefore the available capacity changes dramatically. Moreover, mmWave networks will be small cell networks, and mmWave links are sensitive to blockage from foliage, the human body, moving obstacles and so on, thus the transitions from LOS to NLOS for mmWave connections will be much more frequent than it is in LTE [42]. Furthermore, given the shorter range of mmWave communication, it is more probable for mmWave links than for LTE ones to experience an outage (i.e., no signal is received) because of shadowing;

- the mmWave channel has a shorter coherence time, which results in variations of the channel of the order of hundreds of microseconds [42], that are faster than those in current mobile networks. As shown in Fig. 6.1, the transitions of the LTE channel due to fading are much smoother than those of the mmWave link.

In this chapter, which is based on [389, 393, 404], we will highlight that the extreme variability of the signal quality over mmWave links yields either a degraded TCP goodput and a very low utilization of the resources at mmWave frequencies, or, in the presence of link-layer retransmissions, high goodput at the price of high latency. This is detrimental to the experience of the end user in mmWave networks, which will not benefit from the resources available at the physical layer. In the next chapters, instead, we will discuss possible strategies and mechanisms to improve the end-to-end performance when TCP is used as transport protocol, or, alternatively, with UDP.

This chapter aims at describing how TCP will perform in mmWave 5G cellular networks. In particular, we compare the performance of different TCP congestion control algorithms over simulated 5G end-to-end mmWave networks considering a number of realistic deployment scenarios, as further described in Sec. 6.3. Our detailed simulation study demonstrates that the performance of TCP over mmWave depends critically on several aspects of the network:

1. **Edge vs. Remote Server:** By comparing the end-to-end performance at varying server's location, we show that for a shorter control loop, i.e., when the server is placed at the cellular network edge, TCP can react faster to link impairments.

2. **Handover and Multi-Connectivity:** Due to unreliability of individual mmWave links, dense deployments of small cells with fast handover protocols are critical in maintaining stable connections and avoiding TCP timeouts.

3. **CC Algorithms:** With remote servers, we observe higher performance variations across different congestion control algorithms, while the difference is almost negligible with edge servers. Overall, BBR outperforms loss-based TCP in terms of both rate and latency.

4. **TCP Packet Size:** We quantitatively compare the benefits of transmitting larger TCP packets in LTE versus mmWave networks, and show that, given the fluctuating Gbps data rates offered at mmWave frequencies, a larger packet size provides a faster growth of the congestion window and higher achievable rate.

5. **Link-level Solutions:** The retransmission at the link-level (e.g., RLC AM and HARQ in a 3GPP stack) play a fundamental role in the tradeoff between throughput and latency.

6. **RLC Buffer Size:** Due to the erratic nature of mmWave channels, buffers need to be carefully sized. We analyze TCP performance over small and large buffers. While the TCP goodput degradation caused by buffer overflow in undersized buffers is difficult to mitigate, the problem of bufferbloating, i.e., large buffer occupancy leading to delays, can be approached by appropriately designing cross-layer algorithms [64].

The rest of the chapter is organized as follows. We first discuss the main features of TCP that are relevant to the interplay with the lower layers of the protocol stack. Then, we describe the scenarios of interest in Sec. 6.3, and report the main results and observations in Sec. 6.4. Sec. 6.5 reports our conclusions and suggestions for research directions related to the improvement of the end-to-end performance in mmWave cellular networks.

## 6.2 TCP Congestion Control Protocols

In this section, we will describe the congestion control protocols and the TCP performance enhancement techniques considered in this performance evaluation.

### 6.2.1 TCP Congestion Control Algorithms

TCP was designed in the 1980s as a connection-oriented and reliable protocol that provides end-to-end connectivity over multiple hops and congestion control (CC). Reliability is enabled by a retransmission mechanism, based on the acknowledgments received by the TCP transmitter. Congestion control is implemented by different algorithms that increase and/or decrease the maximum amount of unacknowledged data that the sender is allowed to transmit (*congestion window*), reacting to network events such as packet losses. There have been several evolutions of the original congestion control algorithms: the latest Request For Comments (RFC) describing them is [53], and the survey in [274] lists 13 TCP variants implemented in the Linux kernel. In this chapter, we consider three commonly used CC algorithms (TCP NewReno, HighSpeed TCP, and TCP CUBIC), and a recent addition, i.e., TCP BBR.[1]

**TCP NewReno** [275] has been the default algorithm for the majority of communication systems. In the congestion avoidance phase, the congestion window `cwnd` is updated after the reception of every ACK. The update is based on the Additive Increase Multiplicative Decrease (AIMD) design: `cwnd` is increased by summing a term $\alpha/$`cwnd` for each received ACK, and divided by a factor $\beta$ for each packet loss. For NewReno these parameters are fixed to $\alpha = 1$ and $\beta = 2$.

**HighSpeed TCP** [276] is designed for high Bandwidth-Delay Product (BDP) networks, in which NewReno may exhibit a very slow growth of the congestion window. HighSpeed behaves the same as NewReno when the congestion window is small, but when it exceeds a predefined threshold the parameters $\alpha$, $\beta$ become functions of the congestion window, in order to maintain

---

[1]We refer the interested reader to our survey [396] for an overview of the most recent novelties related to transport protocols.

a large `cwnd`. Moreover, the window growth of NewReno and HighSpeed depends on the ACK reception rate, thus a shorter Round Trip Time (RTT) increases the ACK frequency and further speeds up the window growth.

**TCP CUBIC** [277], instead, increases the congestion window over time, without considering the ACK reception rate but rather capturing the absolute time since the last packet loss and using a cubical increase function for `cwnd`. It has been designed to increase the ramp-up speed of each connection while maintaining fairness with other users.

**TCP BBR**, recently presented by Google [278], measures bottleneck bandwidth and round-trip propagation time, or BBR, to perform congestion control. It strives to match the sending rate to the estimated bottleneck bandwidth by pacing packets and setting the congestion window to `cwnd` gain $\times$ BDP, where the `cwnd` gain is a factor ($\leqslant 2$) that is used to balance the effects of delayed, stretched and aggregated ACKs on bandwidth estimation.

### 6.2.2 TCP Performance Enhancement Techniques

The performance of TCP has been the object of many studies over the last decades, and, besides new CC algorithms, many other techniques have been proposed and deployed either at the endpoints of the connection (TCP sender and receiver) or inside the network.

In case of multiple packet losses, the TCP Selective Acknowledgment (SACK) option [279] allows the receiver to inform the sender which packets were received successfully, so that the sender can retransmit only those which were actually lost. This dramatically improves the efficiency of the TCP retransmission process.

AQM schemes [280], instead, are deployed in network devices (e.g., routers, gateways, gNBs), to control the behavior of their queues and buffers. The size of these buffers plays an important role in the end-to-end performance. These buffers enqueue packets before they are transmitted on the next hop of the end-to-end connection, and, if this link acts as a bottleneck, or if its capacity suddenly drops (as may happen with a mmWave link), the buffer occupancy grows. If the buffer is too small, many packets may be dropped when the buffer is full, according to the drop tail policy. Conversely, if the buffer is too large, then the bufferbloat phenomenon occurs [280]. , i.e., too many packets are buffered, causing an increase of the end-to-end latency, unfairness across short and long flows, and a late reaction from the TCP sender, which is unaware of the congestion for a long interval and may end up triggering a retransmission timeout. AQM techniques can be deployed at the buffers to drop packets before the queue is full, so that the TCP sender can proactively react to the congestion that could arise in the near future. Optionally, packet losses can be avoided by using Explicit Congestion Notifications (ECNs), i.e., a congestion flag in the IP header of the packet is activated, and when a TCP endpoint receives a packet with such a flag it propagates this information to the other endpoint.

Finally, there are some techniques that are typically used in combination with wireless links. Current and future mobile networks deploy different retransmission mechanisms in order to prevent packet loss and increase the throughput at the mobile devices. When using mmWave links, these retransmission protocols become a key element in hiding the highly dynamic and consequently unstable behavior of the channel to the higher layer transport protocols. HARQ is used at the MAC layer, as discussed in Chapter 2. When the PHY layer at the receiver receives a packet, but detects the presence of some errors that prevent reliable decoding, it asks for a retransmission. The sender then transmits additional redundancy that helps retrieve the correct version of the packet [117]. Moreover, the RLC layer can perform additional retransmissions with the acknowledged mode. Since the number of retransmissions at the MAC layer is usually limited (typically only 3 attempts are performed), RLC AM offers another way of recovering lost packets. Thanks to periodic reports from the receiver, the RLC AM sender knows which packets

**Figure 6.2:** High speed and urban deployment scenarios

are missing and can retransmit them. The number of attempts that RLC AM can perform is also limited, and, if some packets are still missing, a Radio Link Failure is declared. RLC UM, instead, does not perform any retransmission in addition to those of the HARQ at the MAC layer. These retransmission mechanisms operate based on information related to the link and with a greater timeliness with respect to TCP, which instead uses packet losses to detect congestion and operates on the larger timescale of retransmission time-outs (RTOs), of the order of a second.

Another technique which is often used in wireless networks is proxying [269], i.e., the connection is split into two at some level in the mobile network (e.g., at the gateway with the internet, at the gNB, etc), and different CC techniques are deployed over the two parts of the connection. We will discuss proxy-based solutions for mmWaves in the next chapters.

## 6.3   5G Deployment Scenarios

In order to assess how TCP will perform in mmWave cellular networks, we consider two of the most challenging scenarios among those specified by the 3GPP in [15], i.e., a high speed train and a dense urban scenario, represented in Fig. 6.2. They were implemented in the ns–3-based mmWave end-to-end simulation framework described in Chapter 2.

**High speed scenario:** In this scenario, shown on the left side of Fig. 6.2, we test the performance of TCP over a channel that varies frequently in time and under realistic mobility conditions. Multiple gNBs provide coverage to the railway, which is mostly LOS: even if the current gNB is blocked by obstacles placed between gNBs 2 and 3, the UE can quickly perform a handover to another LOS gNB. The gNBs are at a height of 35 meters, with an intersite distance of 580 meters. The train moves at a speed of 108 km/h, and, as a result, the channel experienced by the UE varies very quickly because of severe fading and the Doppler effect, and, on a longer time scale, due to obstacles, as shown in the SINR plot of Fig. 6.2. The UE performs handovers across the different gNBs in order to get the highest SINR possible. For example, during the interval between $t = 23$ s and $t = 32$ s, the closest gNBs are blocked by some obstacles, thus the UE switches to the farther base stations, gNB 1 and gNB 4, which are still in LOS. We use the channel tracking and mobility scheme described in Chapter 3, which features fast and locally coordinated handovers for devices that are dual-connected to a mmWave gNB and a sub-6 GHz gNB (e.g., an LTE base station).

**Dense urban scenario:** In this deployment, shown on the right side of Fig. 6.2, we study

145

the fairness of TCP flows over multiple UEs with different channel conditions. A single mmWave gNB placed at a height of 25 meters serves a group of ten users moving at walking speed. They are located in different positions, in order to account for a mixture of channel conditions: four UEs are in LOS, thus perceiving a very high SINR, four are in NLOS and the last two are inside a building, so that the received power is additionally attenuated by the building penetration loss.

For both scenarios we consider two deployments of the TCP server which acts as the endpoint of the connection. The first is a traditional setup in which the server is hosted in a remote data center, with a minimum RTT in the order of 40 ms, accounting for the latencies of both the core network and the public internet. The second is a Mobile Edge Cloud (MEC) scenario [281], in which the server is located close to the gNBs with smaller latency (of the order of 4 ms).

## 6.4 TCP Performance on mmWave links

In the following paragraphs we will report the performance of the TCP congestion control algorithms presented in Sec. 6.2 over the 5G mmWave deployment scenarios described in Sec. 6.3, focusing on both goodput and latency. The results are averaged over multiple independent simulation runs, so that the confidence intervals are small (they are however not shown to make the figures easier to read). In all the simulations, we use full buffer traffic with the TCP SACK option and disable the TCP delayed ACK mechanism, thus each received packet will generate an ACK. The minimum retransmission timeout is set to 200 ms. Moreover, we always consider RLC AM with HARQ retransmission, except than in Sec. 6.4.3, where we compare the performance with different link-level retransmission schemes.

### 6.4.1 High Speed Deployment Scenario

In this scenario we compare different combinations of the RLC buffer size $B$ and the Maximum Segment Size (MSS) $P$ with a single TCP connection from the UE. For both the remote and the edge server deployments the RLC buffer is 10% or 100% of the BDP computed considering the maximum achievable data rate (3 Gbps) and the minimum latency, i.e., $B$ equals 1.5 or 15 MB for the remote server deployment, and 0.15 or 1.5 MB for the edge server. We also consider two different MSS, i.e., a standard MSS of 1400 bytes (1.4 KB) and a large MSS of 14000 bytes (14 KB). The goodput of saturated UDP traffic is also provided as a reference for the maximum achievable rate, as shown in Fig. 6.3.

Notice that, thanks to the mobility management scheme based on dual connectivity and fast secondary cell handover, and despite the high mobility of the scenario, we never observed a TCP connection reset due to an outage, i.e., even if the closest two base stations are blocked, the UE is still capable of receiving signals from other nearby gNBs. Therefore, even if blockage events are still possible, in a scenario with a dense deployment (according to 3GPP guidelines), it is possible to provide uninterrupted connectivity to the final user [409].

In the following paragraphs we will provide insights on the effects of the different parameters on TCP performance over mmWave at high speed.

**Impact of the server deployment**

Loss-based TCP benefits from the shorter control loop related to an edge server deployment, as shown by comparing Figs. 6.3a and 6.3b. With the latter, indeed, the differences between the maximum goodput of the loss-based TCP versions are less marked, since the faster reaction time makes up for the differences among them. Moreover, the RTT difference between the large

**(a)** Remote server.

**(b)** Edge server.

**Figure 6.3:** Goodput and RTT for the high speed train scenario, with the remote and the edge server for different combinations of the buffer size and the MSS.

and the small RLC buffer is lower in absolute terms (milliseconds with edge server versus tens of milliseconds with remote server), but the ratio is approximately the same. However, for CUBIC and HighSpeed with the smallest buffer configuration, the goodput is lower with the edge than with the remote server, i.e., there is a 30% loss with the smallest MSS, and of 50% with the largest one. In this case, indeed, the buffer size is very small (i.e., $B = 0.15$ MB), thus incurring buffer overflows[2], which reduce the sending rate.

**Impact of the congestion control algorithm**

The congestion control algorithm has a stronger impact in the remote server scenario. The best performance, in terms of goodput, is achieved by BBR with large buffer size, but it is still 400 Mbps lower than the maximum achievable rate. Moreover, as observed in [57, 389], high goodput values also correspond to higher end-to-end latency. However, with small buffers, BBR produces the highest goodput (especially in the edge server scenario), with a latency comparable to loss-based TCP. BBR, indeed, regulates its sending rate to the estimated bandwidth and is not affected by packet loss, i.e., the congestion window dynamics of BBR, presented in Fig. 6.4a, matches the SINR plot in Fig. 6.2.

However, the loss-based versions of TCP cannot adjust their congestion window fast enough to adapt to the channel variations and perform worse than BBR, especially with small buffer, as seen in Fig. 6.4a. Among them, TCP HighSpeed provides the highest goodput because of the aggressive window growth in the high BDP region. TCP CUBIC performs better than NewReno in the remote server case, but worse in the edge server case. This is because CUBIC's window growth is not affected by the ACK rate, and therefore is more reliable over long RTT links.

---

[2]With large MSS just 11 packets are enough to cause a buffer overflow.

**(a)** TCP congestion window with mmWave

**(b)** TCP congestion window with LTE

**Figure 6.4:** Congestion window evolution over time for different CC algorithms. The scenario is configured with remote servers and small RLC buffers

**Impact of the MSS**

The MSS does not affect the performance of BBR, which probes the bandwidth with a different mechanism, whereas, for loss-based TCP, the impact of the MSS on the goodput is remarkable.[3] The standard MSS of $P =$ 1.4 KB exhibits much worse performance compared to a larger MSS of $P =$ 14 KB. This happens because, in congestion avoidance, the congestion window increases by MSS bytes every RTT, if all the packets are received correctly and delayed acknowledgment is not used, so the smaller the MSS the slower the window growth. Hence, the MSS dictates the congestion window's growth, which is particularly critical in mmWave networks for two main reasons: *(i)* The mmWave peak capacity is at least one order of magnitude higher than in LTE, so that the congestion window will take a much longer time to reach the achievable link rate. In this case, we can gain in performance by simply using a larger MSS, as depicted in Fig. 6.3. *(ii)* In addition, the channel fluctuations in the mmWave band will result in frequent quality drops, thus often requiring the congestion window to quickly ramp up to the link capacity to avoid underutilizing the channel.

**Large MSS − mmWave vs. LTE**: Aimed at better illustrating why larger packets are particularly important in 5G mmWave networks, we also provide a performance comparison

---

[3]Typically, TCP segments are mapped to multiple MAC/PHY data units, which complicates the dependence between a larger value of the TCP MSS and the correspondingly higher packet error probability over the wireless link. This non-trivial relationship, which would deserve a study by itself, has been properly captured in our numerical results.

|                | Buffer | Remote Server | | Edge Server | |
|----------------|--------|-----|---------|-----|---------|
|                |        | BDP | 10% BDP | BDP | 10% BDP |
| TCP NewReno    | LTE    | 1.06 | 1.17 | 0.80 | 0.65 |
|                | mmWave | 1.81 | 3.96 | 1.27 | 1.15 |
| TCP CUBIC      | LTE    | 1.06 | 1.15 | 1.03 | 0.89 |
|                | mmWave | 2.2  | 1.83 | 1.89 | 1.44 |
| HighSpeed TCP  | LTE    | 1.08 | 0.9  | 0.94 | 0.95 |
|                | mmWave | 1.09 | 1.69 | 1.05 | 0.98 |
| TCP BBR        | LTE    | 1.00 | 0.96 | 1.02 | 0.82 |
|                | mmWave | 1.14 | 0.97 | 1.06 | 1.06 |

**Table 6.1:** Ratio between the goodput achieved with $P = 14$ KB and with $P = 1.4$ KB, for different configurations of the simulated scenario.

against LTE in the same scenario[4], and report in Table 6.1 and Fig. 6.4 detailed results focusing on the impact of the TCP MSS on the congestion window growth and, consequently, on the goodput of the system. Only a single user is placed in the high-speed train scenario, thus the drops in the congestion window are due to the worsening of the channel quality and not to contention with other flows. Fig. 6.4 shows that the loss-based TCP congestion window with a small MSS grows very slowly in congestion avoidance, and consequently loss-based TCP does not fully exploit the available bandwidth during the intervals in which the received signal has a very high SINR (i.e., at $t = 20$ s and $t = 40$ s, as shown in Fig. 6.2). The large MSS helps speed up the congestion window's growth, which translates into higher goodput. Conversely, the goodput degradation associated with small packets is less relevant in LTE networks, given that the goodput is limited by the available bandwidth and not by the congestion window increase rate. These trends are reflected in Table 6.1. Among all loss-based TCP versions, only HighSpeed increases its congestion window fast enough even when transferring small packets. As a consequence, the goodput gain obtained with large MSS values is much smaller.

Large packets introduce an additional benefit: due to (1) a reduced TCP/IP header overhead and (2) a reduced number of TCP ACKs, there will be more available downlink/uplink resources, resulting in higher goodput values.

This solution may not be practical in an end-to-end network in which the Maximum Transmission Unit (MTU) is not entirely in control of the mobile network provider and is typically dictated by the adoption of Ethernet links (i.e., an MTU of 1500 bytes). By contrast, in a MEC scenario, in which the whole network is deployed by a single operator, it is possible to support a large MSS thanks to Ethernet jumboframes [282, 283].

**Impact of the buffer size and AQM**

The buffer size is also critical for the performance of TCP. As shown in Fig. 6.3, large buffers generally yield higher goodput, because the probability of buffer overflow is smaller, and they offer a more effective protection against rapid and temporary variations of the mmWave channel quality. However, when a large buffer is coupled with loss-based TCP, the latency inevitably increases. Conversely, smaller buffers provide lower latency at the price of lower goodput.

For loss-based TCP, an intermediate solution is provided by applying AQM to the largest buffer, especially in the remote server scenario. Controlled Delay Management (CoDel) is used as the default AQM in our simulation because of its simple configuration. It controls the buffering

---

[4]For the LTE setup the small buffer represents 50% of the BDP (i.e., 0.08 and 0.2 MB for edge and remote server, respectively), because a 10% BDP buffer would be too small to protect from random fluctuations of the channel.

**Figure 6.5:** Goodput vs RTT for ten UEs in the Urban Scenario, for different choices of the CC algorithm.

latency by dropping packets when the minimum queuing delay within an interval, starting from 100 ms, is larger than 5 ms. CoDel is picked as an example to show the trade-off between latency and goodput by using AQM. Our goal is not to select the best AQM scheme or optimize AQM, which in itself is a very interesting topic, and could be considered for future research. As shown in Fig. 6.3a, the goodput with the AQM option is larger than that with the smallest buffer, and in some cases (i.e., for the smallest packet size) is comparable to that of the BDP buffer without AQM, which in general yields the highest goodput. However, the latency is equivalent to the one associated with the small buffer, which is the lowest. In the edge server scenario the TCP control loop is short (the RTT is 4 ms) and the reaction to congestion is quick. Hence, its performance is indeed equivalent to having BDP buffers without AQM.

BBR tries to solve this problem without modifying the buffers in the routers by maintaining a congestion window equal to twice the BDP regardless of packet loss, as shown by Fig. 6.4. As a consequence, latency is only doubled in large buffers, and the goodput is slightly reduced in small buffers. These behaviors are also observed in the oversized and 10% BDP buffer cases of Fig. 6.3.

### 6.4.2 Urban Deployment Scenario

In this scenario we consider ten UEs attached to a single mmWave gNB. In particular, we position four UEs in LOS conditions, four in NLOS and two inside a building. The average SINR for each channel condition is provided in Fig. 6.5. Notice that, with low blockage density and walking speed, the channel condition is relatively stable over time. For each UE pair one is connected to an edge server, and the other is connected to a remote server. In this way, it is possible to test the performance of TCP over a mixture of different conditions. The gNB uses a RR scheduler, so that the resource management at the base station does not have an impact on the fairness among different flows. All the UEs use the same TCP version. We consider a standard MSS of 1400 bytes and an RLC buffer size of 1.5 MB for each UE.

Fig. 6.5 shows the average cell goodput (labeled in parentheses) and the goodput-latency trade-off for each type of user, separately, and for each CC algorithm, in order to evaluate the fairness and the overall performance of different TCP versions with respect to different user channel conditions.

All CC algorithms achieve the same average cell goodput, and similar goodput per UE. However, the RTT varies a lot among the CC algorithms. The reason is that all UEs use the same buffer size regardless of their channel conditions and network latency. As a consequence,

the RLC buffer size may be large for some UEs, such as those at the edge. Therefore, the CC algorithms that adopt a more aggressive window growth policy, such as CUBIC and HighSpeed, yield much higher latency. For the loss-based TCP, NLOS and indoor UEs suffer from a higher latency: given the same buffer size and backhaul data rate, a reduced bottleneck bandwidth results into an increased queueing delay in the buffers, until TCP settles to a steady state phase. BBR, instead, limits the congestion window to twice the estimated BDP, and results in a maximum latency of $2 \times$ minimum RTT. We also draw a gray area in the plot representing a typical 5G application requirement, i.e., goodput greater than 100 Mbps and delay lower than 10 ms. In this scenario, among all CC algorithms, only BBR meets this requirement for the UEs connected to an edge server, and only under good channel conditions.

### 6.4.3 Impact of the link-level retransmissions

In order to test the effectiveness of coupling TCP with lower-layer retransmission mechanisms, we performed some simulations using the framework described in Sec. 2.11, where we considered an uplink connection from a UE placed at different distances from a gNB. We use IPERF on top of the Linux implementation of TCP CUBIC, with the statistical channel model [37], and perform Montecarlo simulations for each distance $d \in \{50, 75, 100, 150\}$ m.

RLC AM introduces additional redundancy in order to perform the retransmissions, but, when the distance between the eNB and the UE is equal to $d = 50$ m and the UE is in LOS with very high probability, these retransmissions are not actually needed, because of the low packet error rate of the channel. Therefore, as also shown in Fig. 6.6, the throughput is lower when RLC AM is used (though by only a minimal amount). As the distance increases, instead, the performance of TCP without HARQ and without RLC AM collapses, because the TCP congestion control algorithm sees a very lossy link and triggers congestion avoidance mechanisms or, worse, a RTO.

If instead we compare the performance of HARQ with RLC UM and that of HARQ with RLC AM, it can be seen from Fig. 6.6 that the additional retransmissions given by RLC AM increase the throughput by 100 Mbps at $d = 75$ m and 50 Mbps at $d = 100$ m. For $d = 150$ m, instead, RLC AM does not improve the performance of RLC UM, showing that at such distance even further transmission attempts fail to successfully deliver packets (for example, because of extended outage events).

RLC AM at large distances instead increases the latency of successfully received packets, as



**Figure 6.6:** Latency throughput tradeoff for TCP CUBIC, with and without the different retransmission mechanisms of the mmWave protocol stack.

**Figure 6.7:** Download time as a function of the file size and of the distance, for TCP CUBIC with and without lower-layer retransmissions.

shown in Fig. 6.6, because of retransmissions and additional segmentation that may introduce Head of Line (HoL) blocking delays. The smallest latency is achieved without HARQ and with RLC UM, because no retransmissions are performed, but this option is not able to deliver a high TCP throughput in general.

Fig. 6.7 shows the download time for a file of different sizes (from 1 MB to 10 MB) using wget (the file is hosted in the UE and retrieved by the remote server, in order to be consistent with the previous uplink simulations). The results show that lower-layer retransmission mechanisms help decrease the download time, and that the performance gain increases as the distance and the file size increase. Moreover, the difference between the download times with RLC AM and with RLC UM (no retransmissions) is more noticeable than that between the throughput values of Fig. 6.6, showing that for short-lived TCP sessions it is important to perform retransmissions as fast as possible, i.e., at a layer as close to the radio link as possible.

These results are well known when applied to traditional LTE networks [270], but these are the first simulations that show how much TCP depends on lower-layer retransmissions in mmWave networks, using the real Linux TCP/IP implementation. They show that, also in mmWave networks, the support of lower-layer retransmission mechanisms is fundamental for reaching a high TCP throughput even at large distances between transmitter and receiver, at the price of additional latency. In particular, in the simulated scenario the most effective retransmission scheme is HARQ at the MAC layer, since it provides the greatest throughput gain, but also the acknowledged mode of the RLC layer helps improve the performance of the mmWave link by reducing the download time for short-lived TCP sessions.

## 6.5 Final Considerations

The massive but intermittent capacity available at mmWave frequencies introduces new challenges for all layers of the protocol stack, including TCP, the most widely used transport protocol. The interplay between congestion control algorithms and mmWave channel quality fluctuations makes the topic particularly complex, and represents the key driver behind this work. We have carried out a thorough simulation campaign, based on ns-3, across 3GPP-inspired scenarios, whose results are summarized in Table 6.2. The main findings and some relevant research questions are listed as follows:

| | Loss-based | MSS impacts goodput | Summary | Considerations over 5G |
|---|---|---|---|---|
| TCP NewReno | yes | yes | remote server: lowest goodput | need to move servers to the edge |
| TCP CUBIC | yes | yes | edge server: lowest goodput | need to increase MSS |
| HighSpeed TCP | yes | only remote server | big buffer: high goodput and high latency | need to mitigate latency with AQM |
| TCP BBR | no | no | *big buffer:* high goodput and high latency<br>*small buffer*: small rate reduction and low latency | small buffer is preferred<br><br>performs well over mixed UE conditions |

**Table 6.2:** Results of the CC algorithms over 5G deployments

1. TCP benefits from a shorter control loop, where the server is placed at the cellular network edge and can react faster to link impairments. *Should we (re)consider splitting TCP at some point?*

2. Moreover, when the RTT is high, loss-based TCP underutilizes the mmWave capacity, while those based on congestion (e.g., BBR) show an improved performance by estimating the bandwidth of mmWave links. This means that new approaches based on *more refined abstractions of the end-to-end network* can be studied for highly-variable and high-data-rate mmWave links.

3. Multi-connectivity and smart handovers, supported by advanced beamtracking and beam-switching techniques, will result in robust TCP connections. *How densely should we deploy mmWave cells? Should we deploy both LTE and mmWave as a multi-tier overlay network? How to support backhaul for densely deployed mmWave cells?*

4. We show very clearly how loss-based TCP over mmWave bands can greatly benefit from using larger datagrams. *Has the time come to break the legacy MTU value,* by natively supporting larger packets in a wider set of networks?

5. Finally, it is well known that buffer size must scale proportionally to BDP to achieve maximum TCP goodput. However, it is very challenging to properly dimension the buffers for mmWave links, given the rapid bandwidth variations between LOS and NLOS conditions, and to protect from link losses without introducing bufferbloat. Given the low latency requirement and massive available bandwidth, *is it beneficial to trade bandwidth for lower latency,* for example by running BBR over small RLC buffer configurations?

We believe that these insights will stimulate further exploration of this important topic, which is essential to fully exploit the true potential of mmWave communications. Moreover, the observations provided by this initial simulation-based study have been used to guide the design of novel techniques to improve the end-to-end user experience in mmWave cellular networks, as we will discuss in the next chapters, and of experimental activities, which are necessary to further validate the challenges that mmWave links pose to TCP.

# 7

# TCP performance enhancing techniques in mmWave Networks

## 7.1 Introduction

As discussed in Chapter 6, the characteristics of the mmWave channel significantly impact the performance of TCP, with a degradation in the throughput and an increase in buffering latency after LOS to NLOS transitions. In order to address these issues, in this chapter, which is based on [407, 411], we propose the design of two cross-layer solutions that improve the coordination between the transport layer and the wireless protocols stack.

First, we will describe *milliProxy*, a novel TCP proxy for mmWave mobile networks aimed at fully reaping the benefits of mmWave links to achieve high throughput with low latency. It is transparent to the end hosts of the connection, and respects the end-to-end connection semantics. The main rationale is to split the TCP control loop in the mmWave RAN to optimize the flow control over the wireless link. It is based on a cross-layer, data driven approach and enables a number of optimizations for the operation of TCP in mmWave networks. Then, we will introduce X-TCP, a cross-layer congestion control for TCP uplink flows on mmWave links. X-TCP exploits the knowledge of resource allocation and channel quality at the UE side to tune the congestion window. The main goal of both approaches is to avoid that the TCP sender sends more packets than those that can actually be delivered on the mmWave wireless link, preventing the increase of the queues occupancy and of the end-to-end packet delay. We test the effectiveness of the proposed schemes using the mmWave module of ns–3 introduced in Chapter 2, and compare the performance of these approaches against TCP CUBIC and other congestion control protocols in terms of latency, throughput and fairness in randomly generated scenarios.

The rest of the chapter is organized as follows. In Sec. 7.2, we provide an overview of the literature related to TCP proxies for traditional wireless networks. The architecture of milliProxy is described in Sec. 7.3, and the results of a performance evaluation campaign are reported in Sec. 7.4. X-TCP, instead, is introduced in Sec. 7.5, with the analysis of its performance in Sec. 7.6. Finally, conclusions and future extensions of this work are provided in Sec. 7.7.

## 7.2 Related Work on TCP Performance Improvement with Proxies

The performance of TCP on wireless networks has been under the spotlight since the 1990s, when the first cellular networks capable of data transmission were commercially deployed. Even though TCP faces more challenging conditions when running on top of mmWave cellular networks, as discussed in [284], it is worth describing the main approaches that can be found in the literature related to the enhancement of TCP performance on wireless links.

A first comprehensive review on the topic can be found in a paper by Balakrishnan *et al.* [285]. The authors claim that the poor performance of TCP in mobile networks is due to packet losses over an unreliable channel. However, as shown in [404], the channel losses can be masked by retransmission mechanisms. Moreover, the considered links have very low data rate and small buffers are used in the network. The settings in a mmWave networks are very different, since large buffers and retransmissions are already implemented in the wireless link to make up for packet loss at the price of increased latency and exposing more the network to the bufferbloat phenomenon. However, the authors of [285] provide a comparison of different strategies that can possibly be adapted to mmWave networks, using TCP Reno as a baseline, and including also TCP split approaches.

In a more recent paper [269], Liu *et al.* introduce a TCP proxy middlebox for the optimization of TCP performance without the need for any modification to the protocol stack of servers, clients and base stations. They observe that the adoption of a new end-to-end TCP congestion control mechanism may be useless in the presence of HTTP proxies, which are frequently used in mobile networks. Moreover, they design their solution for modern LTE networks, characterized by large buffers (in the order of 5 MB) and bandwidth fluctuations (even if not as wide as those in mmWave networks [42, 389]), and a fixed network which does not act as a bottleneck. Their solution is a middlebox that can be placed anywhere in the mobile operator core network, and breaks the TCP connection in two segments, i.e., it does not respect end-to-end connection semantics[1]. This box performs some optimizations on the fly, such as (i) not using the information of the receiver congestion window, which may be too small with respect to the actual rate available on the link, given that experimental evaluations on the receiver buffer in real devices have highlighted that it is never filled; (ii) changing retransmission patterns by intercepting duplicate ACKs; (iii) tuning the congestion window with a rate estimation algorithm. In this design, the TCP connection from the sender to the receiver is terminated at the middlebox, which buffers the packets for the final receiver until it can forward them.

A third approach is described in [288], where Ren *et al.* introduce a TCP proxy in the mobile network base station. This study, however, is focused on the UMTS architecture. Their approach is based on a queue control mechanism: by using the sliding mode variable structure (SMVS) control theory the buffer queue length at the base station is kept at the same size. This proxy does not respect the end-to-end TCP semantics, because it terminates the connection at the proxy. The advertised window at the proxy is used to limit the sending rate of the server and to avoid buffering delays. At the proxy, a control mechanism is used to keep the queue length at a reference value, by inferring the bandwidth available at the base station.

Some other interesting approaches that respect the end-to-end TCP semantics are (i) Mobile TCP (M-TCP) [289], which freezes the TCP sender when it senses imminent congestion, in order to avoid packet loss and connection timeouts; and (ii) Snoop [290], also from Balakrishnan *et al.*, which performs local retransmissions when TCP packet losses are sensed, in order to improve the

---

[1]According to [286, 287], the end-to-end principle, which states that certain functions in the internet are designed to work at the end hosts, is a founding paradigm of the internet. A proxy that splits the TCP connection into two independent segments does not respect the end-to-end TCP semantics, meaning that ACKs may be sent to the sender before the packet is actually received by the other end host.

connection reactiveness. I-TCP [291], instead, is a TCP split approach, not compliant with the end-to-end TCP semantics, that uses a traditional TCP congestion control also on the wireless link and does not yield a great performance improvement.

In [147, 292], a performance enhancement proxy for mmWave cellular networks is proposed. It is installed in the base stations, and breaks the end-to-end TCP semantics by sending early ACKs to the server. Moreover it performs batch retransmissions, i.e., it retransmits the packets that were detected as lost as well as the segments with a sequence number which is close to that of the lost packets. The detection of the lost packets however assumes that on the link only HARQ retransmissions are performed, while in general with TCP the RLC AM is also used. In the performance evaluation, moreover, the authors of [147] limit the application data rate to 100 Mbps, which can be usually sustained also in NLOS. Therefore, the performance analysis does not account for the very high data rates that can be achieved with mmWave and for the wide rate variations of the LOS to NLOS (or vice versa) transitions. Finally, it focuses only on the throughput and delivery ratio for a single user, without considering the latency and thus the bufferbloat problem. Finally, the authors of [64] propose a cross-layer adaptation scheme for downlink flows based on the estimation of the capacity of the mmWave link at the UE side and the adaptation of the receive window.

## 7.3   End-to-end Proxy Architecture for mmWave

In this section we describe our TCP proxy architecture for mmWaves, called milliProxy, and highlight the main innovations with respect to the solutions reported in Sec. 7.2. Importantly, by being transparent to both the end points of the TCP flow, milliProxy respects the end-to-end semantics of the TCP connection, as opposed to most of the proposed approaches cited in Sec. 7.2. The key functionalities of milliProxy are (1) the ability to split the control loop of the connection with a different and tunable Flow Window (FW) policy at the source server and at the proxy, as well as (2) the capability of controlling the MSS of the connection in the portion between the proxy and the UE.

### 7.3.1   Proxy Architecture

MilliProxy is a TCP proxy which can be implemented and deployed as a network function, composed of several modules that can be updated or changed. It can be placed in the gNB, fully benefiting from the interaction with the mmWave protocol stack, or in a node in the core network, sharing out-of-band information with the gNB to which the TCP receiver is connected. According to the position of the proxy, there may be the need to design a mechanism to cope with the user mobility. For example, if the proxy is in the gNB, when the UE performs a handover the network has to transfer the milliProxy's state from the source to the target gNB. If instead the proxy is in an edge node of the core network, then it can manage multiple cells without the need to forward the state for each UE handover. Additional considerations on this issue are left for future work.

The basic structure of the proxy is shown in Fig. 7.1. An instance of the proxy is created for each TCP flow which goes through the node in which it is installed, so that different policies can be enabled for different users, or different flows of the same user. Each instance has its own customizable buffer (set by default to 10 MB), flow window module and ACK management unit. The buffer is used to store the payload of the TCP packets before they can be delivered to the TCP receiver, and the ACK management unit checks for incoming ACKs to clear the contents of the buffer. The flow window policy is the equivalent at the proxy of the congestion window

**Figure 7.1:** Architecture of milliProxy

mechanism at the TCP sender, i.e., it controls the amount of data that can be forwarded by the proxy. The policy is not hard-coded into the proxy, but is loaded as a module, according to the implementation of the TCP congestion control mechanisms in the Linux kernel [85]. The ACK management unit of the proxy modifies the *advertised window* in each ACK that is relayed to the server in order to enforce the proxy flow window value also at the TCP sender. According to [53], the TCP sender selects the minimum between its congestion window and the received value of the advertised window as the maximum number of bytes it can send. Similarly to [64], the advertised window in the modified ACKs is set to be equal to the flow window determined at the proxy. This makes it possible to capture both components of the network, and adapt accordingly: the wired part is regulated by the classical TCP congestion control selected, while the wireless channel is used in a cross-layer fashion by the proxy, which selects the proper value of the advertised window.

The presence of the buffer makes it possible to tune the MSS of the connection between the proxy and the UE differently from that of the other part of the connection, enabling further optimizations. If the MSS of the overall connection is limited by the MTU of some intermediate networks using ethernet as link layer technology (i.e., the MSS is at most 1460 bytes), then the proxy buffers the 1460-byte payloads, and can send a larger segment which aggregates multiple payloads of the end-to-end connection. For example, fourteen 1460-bytes payloads received back-to-back in a small time interval can be combined into a single 20440-bytes segment which is sent from the proxy to the UE. This increases the efficiency of the transmission process in the last mile of the connection, i.e., in the mmWave wireless link, because of the smaller overhead of the TCP/IP headers (in the previous example, just one TCP/IP header is used instead of fourteen), and because fewer uplink resources have to be scheduled for the transmission of ACKs from the UE [393]. Notice that aggregation is generally performed also at the RLC and MAC layers of very high-bandwidth connections in order to improve the transmission efficiency [293, 294], and the larger MSS helps also this process, since fewer concatenation and segmentation operations are required at the transmitter and the receiver.

Fig. 7.2 depicts how a packet is processed by milliProxy. By design, it is completely transparent to the UE, i.e., the TCP receiver. It intercepts all the packets belonging to the flows it is handling, and the payload of data packets is stored in the proxy buffer. Any options in the packet header are processed, for example to estimate the RTT, as will be described in the following sections, or to handle the advertised window scaling. The payload will then be sent as part of a larger segment as soon as the flow window allows it. When an ACK is received, the proxy checks its sequence number, and marks the corresponding bytes in the buffer as received,

**Figure 7.2:** Packets processing in a milliProxy instance.

which will then be discarded, allowing the flow window to advance. Consequently, a number of ACKs corresponding to the number of original packets received (approximatively equal to the ratio between the MSS of the proxy-UE connection and that of the server-proxy connection) is sent to the TCP sender. In each ACK the *advertised window* value is overwritten with the value of the flow window in the proxy.

## 7.3.2 RTT estimation

The estimation of the RTT can be performed using the TCP *timestamp* option [295]. This option is symmetric, i.e., it is added both to data packets at the TCP sender and to ACKs at the receiver. It has a total length of 10 bytes, and contains two timestamps. The first ($TS_{\mathrm{val}}$) is that of the clock of the end host that transmits the packet, the second ($TS_{\mathrm{echo}}$) is the $TS_{\mathrm{val}}$ of a recently received packet from the other end host. Its usage is advised in [295] in order to improve the TCP performance, in terms of both throughput and security.

If both the end hosts share the same clock, the estimation of the RTT is composed by two phases as follows. In the first one, which is shown in Fig. 7.3, the milliProxy instance estimates the latency on the path from the UE to the server. The timestamp $TS_{\mathrm{echo}}$ of the data packet sent from the server to the UE corresponds to the time $t_{-1}$ at which the UE sent an ACK. Similarly, the timestamp $TS_{\mathrm{val}}$ in the same packet corresponds to the time instant $t_0$ at which the server transmitted the data packet. Given the very high packet rate that is sustained in mmWave networks, it is unlikely to observe a significant time interval between receiving the ACK corresponding to $TS_{\mathrm{echo}}$ and sending the data packet corresponding to $TS_{\mathrm{val}}$. Therefore, the latency of the uplink path can be estimated as $T_{\mathrm{UE}\rightarrow\mathrm{server}} = t_0 - t_{-1}$. In a similar fashion, it is possible to use the timestamp values carried by ACK packets to estimate the latency on the downlink path $T_{\mathrm{server}\rightarrow\mathrm{UE}}$. Finally, the RTT is estimated as $RTT_e = T_{\mathrm{server}\rightarrow\mathrm{UE}} + T_{\mathrm{UE}\rightarrow\mathrm{server}}$.

If instead the two end hosts do not have the same clock, or if the TCP timestamp option is not supported, other methods can be used to estimate the RTT as reported in [296].



**Figure 7.3:** RTT computation at the proxy.

159

### 7.3.3 Integration with the 5G protocol stack

The proxy is configured to collect some statistics from the connected 5G gNB. According to the location of the proxy, this data collection can be performed with or without delay. If the proxy is installed in the gNB, the information can be retrieved instantaneously, whereas if it resides in a node in the core or edge network some signaling is necessary, which would introduce some incremental latency. Thanks to this information it is possible to enable a cross-layer approach, which is useful for the design of flow window management algorithms driven by the performance and the statistics of the mmWave link.

More information associated with each user can be retrieved from the protocol stack of the gNB. The first is the RLC buffer occupancy $B$, which can be seen as a signal of a congestion event and a consequent increase in latency. The second is an estimate of the PHY layer data rate between the UE and the gNB. In [269] this is done by measuring the number of bytes transmitted in the previous slots, dividing it by the duration of the slots. This approach, however, is sensitive to the actual rate that is injected in the network by the TCP source, and can lead to an underestimation of the available rate if the source rate does not saturate the connection. This limitation is particularly relevant in mmWave networks, where it takes a long time for the TCP source to reach a full utilization of the available resources. For milliProxy and X-TCP (that will be introduced in the following sections), instead, we rely on the information provided by the AMC module and the scheduler at the MAC layer. By knowing the channel quality of a UE it is possible to compute the modulation and coding scheme, predict how many bytes the scheduler could allocate to the user (with full buffer assumption) in the next time slot, and divide by its duration to obtain an achievable data rate $R_e$ that is not influenced by the source rate. Another useful metric that can be acquired in a cross-layer setup is the SINR of the UE, which could give an indication on the link status: for example, if it is below a certain threshold, then the proxy will know that the UE is in outage.

### 7.3.4 Window Management

The management of the flow window is an essential component of milliProxy. In this chapter we propose a scheme based on the computation of the BDP. The implementation and testing of alternative FW management policies is left for future works.

In the BDP-based scheme, the FW management module uses three different kinds of cross-layer data: the RLC buffer occupancy $B$, the estimated data rate $R_e$ and the estimated RTT $RTT_e$. The BDP is then given by $R_e RTT_e$.

However, notice that, when the queueing delay in the RLC buffers starts increasing, then the RTT also increases and the estimate is artificially inflated. This worsens the performance of the proxy, since, if the flow window blindly follows the BDP estimate, the increase of BDP due to the longer queueing delay would increment the sender rate, thus further exacerbating the congestion. Therefore, following the approach described in [64, 297], we filter $RTT_e$ and use the minimum value observed in a certain time interval, $RTT_{\min}$, as an estimate of the latency without buffering delays. The mmWave link latency (without retransmissions) has a very limited impact on $RTT_{\min}$ since it is smaller than 1 ms, so that the forwarding delays introduced by the core network and the public internet are dominant. Therefore, the mobility of the UE in a cell or across different cells has almost no effect on $RTT_{\min}$.

The flow window is then computed as $w = \lfloor RTT_{\min} R_e \rfloor$. When the RTT estimate is not yet available (i.e., for the first ACK after the reception of the SYN packet), the flow window is arbitrarily initialized to a high value of 400 MB. Moreover, it is possible to make the policy more conservative when the RLC buffer occupancy exceeds a predefined value (e.g., 2 MB). In this case, the flow window is set to $w = \max\{\lfloor RTT_{\min} R_e \rfloor - 2B, 0\}$.

| Parameter | Value |
|---|---|
| mmWave carrier frequency | 28 GHz |
| mmWave bandwidth | 1 GHz |
| 3GPP Channel Scenario | Urban Micro |
| Max PHY layer rate | 3.2 Gbps |
| S1 link latency $D_{S1}$ | 1 ms |
| Latency from PGW to server $D_{RS}$ | $[1, 5, 10, 20]$ ms |
| RLC AM buffer size $B_{\mathrm{RLC}}$ | $[10, 20]$ MB |
| RLC AM Reordering Timer | 1 ms |
| RLC AM Report Buffer Status timer | 2 ms |
| UE speed $v$ | 5 m/s |
| TCP $\mathrm{MSS}_1$ (server - proxy) | 1400 byte |
| TCP $\mathrm{MSS}_2$ (proxy - UE) | 20000 byte |

**(a)** Simulation parameters



**(b)** Randomly generated simulation scenario. The three grey rectangles represent obstacles such as buildings, cars, trees.

**Figure 7.4:** Simulation parameters and scenario for the milliProxy evaluation.

## 7.4 milliProxy Performance Evaluation

### 7.4.1 Scenario and parameters

We implemented milliProxy in the ns-3 mmWave module described in Chapter 2. The main simulation parameters are reported in Table 7.4a. In this section we focus on testing the performance of milliProxy in a single user scenario, in order to evaluate the responsiveness of the proxy architecture to channel variations, from LOS to NLOS and viceversa. In order to model them, some obstacles are randomly deployed in the simulation scenario between the gNB (which is at coordinates $(25, 100)$ m) and the UE (moving from $(0, 0)$ m to $(50, 0)$ at speed $v$). As the user moves, it will experience multiple transitions, with a random duration of each LOS or NLOS phase in each different run of the simulation. An example of scenario is shown in Fig. 7.4b. All the results are averaged over 50 independent simulation runs.

### 7.4.2 Results

Fig. 7.5 shows a comparison of both goodput and RAN latency when milliProxy is deployed in the gNB or not, for different RLC buffer sizes $B_{\mathrm{RLC}}$ and fixed-network latencies. It can be seen that milliProxy performs better in terms of both goodput and latency, with a goodput gain of up to 2.24 times (combined with a latency reduction of 1.98 times) with the highest $D_{RS}$, or a latency reduction of 43 times with a similar goodput in the edge server scenario (i.e., $D_{RS} = 1$ ms). MilliProxy is therefore effective at reducing the impact of the bufferbloat issue: when the channel switches from a LOS to a NLOS state, milliProxy can reduce the TCP sending rate faster, and thereby avoid extra queuing latency. On the other hand, when the channel quality improves, milliProxy is able to (i) track the available data rate at the physical layer and (ii) promptly inform the TCP sender of the increased resource availability, which indeed results in higher goodput. The performance of milliProxy is independent on the buffer size, since it manages to keep the buffer occupancy and consequently the RLC queuing delay to a minimum. As shown in Fig. 7.5 and extensively discussed in Chapter 6, traditional approaches without proxy result in higher goodput at the price of increased RAN latency when using larger buffers.

A comparison between different configuration options for milliProxy is given in Fig. 7.6. In

**(a)** TCP goodput

**(b)** Latency from the PDCP at the eNB to that at the UE

**Figure 7.5:** Comparison of goodput and RAN latency with and without milliProxy, for different buffer sizes $B$.

| $D_{S1} + D_{RS}$ [ms] | 2 | 6 | 11 | 21 | $D_{S1} + D_{RS}$ [ms] | 2 | 6 | 11 | 21 |
|---|---|---|---|---|---|---|---|---|---|
| $B_{\mathrm{RLC}} = 10$ MB | 1.1941 | 1.6875 | 1.7202 | 2.2430 | $B_{\mathrm{RLC}} = 10$ MB | 11.8008 | 4.7547 | 2.5574 | 1.9888 |
| $B_{\mathrm{RLC}} = 20$ MB | 1.0135 | 1.1448 | 1.0765 | 1.9901 | $B_{\mathrm{RLC}} = 20$ MB | 43.3299 | 11.5578 | 5.8104 | 3.6988 |

**(a)** TCP goodput gain when using milliProxy, i.e., ratio between the goodput with milliProxy and with TCP NewReno.

**(b)** RAN latency reduction when using milliProxy, i.e., ratio between the latency with TCP NewReno and that with milliProxy.

**Table 7.1:** Goodput and latency performance gains with milliProxy.

particular, we are interested in studying the sensitivity of goodput and latency with respect to the delay $D_{\mathrm{info}}$ in the acquisition of the cross-layer information from the gNB: it is equal to 0 when milliProxy is deployed in the gNB, and greater than 0 when installed in a node in the core or edge network. We consider $D_{\mathrm{info}} = 3$ ms, i.e., we assume that the latency between the proxy deployed in the core/edge network and the gNB will be smaller than 3 ms. As shown in Fig. 7.6, the two tested configurations have a similar behavior in terms of both goodput and latency, showing that milliProxy is robust with respect to different possible deployments in the edge network or in the gNBs.



**(a)** TCP goodput

**(b)** Latency from the PDCP at the eNB to that at the UE

**Figure 7.6:** Comparison of goodput and RAN latency with different milliProxy configurations. $D_{\mathrm{info}}$ represents the latency needed to forward the cross-layer information from the gNB to milliProxy, $T_{\mathrm{info}}$ is the periodicity at which this information is collected.

## 7.5  X-TCP - Uplink Cross-Layer Congestion Control

Motivated by the promising results obtained in [64] and in Sec. 7.3 for downlink flows, in this section we describe a cross layer approach for TCP flows from the UE towards a remote destination in the Internet. The UE can use the information gathered from different layers in the cellular protocol stack to directly change the value of the TCP congestion window. The basic algorithm is described with the pseudocode in Algorithm 7.1. In particular, as for milliProxy, we set the congestion window to the optimal bandwidth delay product (BDP) that is estimated from the round trip time (RTT) (measured as in Sec. 7.3.2) of the end-to-end connection and the data rate provided by the mmWave link.

As for milliProxy, the available datarate at the physical layer is estimated using cross-layer information from the MAC and PHY layers. Moreover, we scale the PHY data rate by the overhead introduced by the MAC, RLC, PDCP, IP and TCP headers. Finally, if the congestion on the network increases or the SINR $\Gamma$ is below a certain threshold (i.e., $\Gamma < 0$ dB), a scaling factor $\lambda \leq 1$ is used in the computation of the congestion window, in order to decrease the aggressiveness of the protocol and account for the additional retransmissions performed by the RLC and MAC layers, or for the possible congestion in an intermediate link on the end-to-end path. The algorithm assumes that the path is congested when the estimated round trip time exceeds $RTT_{\min}$ by a certain threshold $\epsilon$, which is set to 10 ms following the approach in [64]. The parameter $\lambda$ depends on the scenario and on the configuration of the cellular network. In our implementation, we used a value $\lambda = 0.85$, which was observed to be a good tradeoff between the throughput loss and the reduction of the queueing delay and, in turn, of the RTT. A context-based optimization of this parameter is left for future work.

In order to test the performance of the proposed approach, we conducted an extensive simulation campaign that will be described in Sec. 4.6.

---

**Algorithm 7.1** Cross layer congestion window update

---

**initialization**
$rtt_{\min} \leftarrow \infty$
cwnd $\leftarrow$ Maximum Segment Size (MSS)

**for** every received ACK
   estimate RTT $RTT_e$
   from the mmWave stack:
      estimate mmWave data rate $R_e$
      get SINR value $\Gamma$
   **if** $RTT_e < RTT_{\min}$
     $RTT_{\min} \leftarrow RTT_e$
   **end if**

   **if** $\Gamma \geq 0$ and $RTT_e \leq RTT_{\min} + \epsilon$
     cwnd$\leftarrow R_e\, RTT_{\min}$
   **else**
     cwnd$\leftarrow \lambda\, R_e\, RTT_{\min}$
   **end if**
**end for**

---

| Parameter | Value |
|---|---|
| mmWave TX power | 30 dBm |
| mmWave carrier frequency | 28 GHz |
| mmWave bandwidth | 1 GHz |
| Number of subframes in one frame | 10 |
| Length of one subframe in $\mu$s | 100 |
| Number of OFDM symbols per slot | 24 |
| Length of one OFDM symbol in $\mu$s | 4.16 |
| Number of sub-bands | 72 |
| UE speed $v$ | [1.75, 5] m/s |
| $R_{\mathrm{app,max}}$ | [1, 2] Gbps |
| $L_{\mathrm{pck}}$ | 1400 byte |
| $\lambda$ | 0.85 |
| $\epsilon$ | 10 ms |
| RLC AM buffer size | 10 MB |
| Core network latency | 1 ms |
| Remote host latency | 10 ms |

**(a)** Simulation parameters



**(b)** First simulation scenario. The grey rectangles are randomly deployed non-overlapping obstacles (e.g., cars, buildings, people, trees).

**Figure 7.7:** Simulation parameters and a random scenario.

## 7.6 X-TCP Performance Evaluation

### 7.6.1 Simulation Setup

The performance evaluation campaign was conducted using the mmWave ns–3 module described in Chapter 2. The mmWave channel model is based on the statistical channel model presented in [37]. We implemented the algorithm described in Sec. 7.5 as a TCP congestion control module, and it can be tested against several other flavors, like BIC, CUBIC, Illinois, NewReno. In all the simulations, we consider uplink traffic from the UE to a remote host, connected to the core network gateway with a high-capacity wired link. The traffic model is full buffer [15], i.e., it always fills the transmission capacity of the TCP pipe with packets of Maximum Segment Size (MSS) equal to $L_{\mathrm{pck}}$ bytes, but we add the possibility of limiting the maximum application data rate to $R_{\mathrm{app}}$. The main parameters of the simulations are summarized in Table 7.7a.

### 7.6.2 Evaluation in Random Scenarios

The scenario we consider consists of a rectangular area, with a mmWave eNB in the center and some objects (buildings, cars, people) randomly deployed over the area, without overlapping. We consider two UEs, named UE1 and UE2, moving at constant speed $v = 1.75$ m/s along straight horizontal lines. Both trajectories cross the area from left to right, but UE1 moves along the lower border of the rectangle, while UE2 is placed on the upper part of the area (see Fig. 7.7b). While moving, the links between the UEs and the eNB alternate between LOS and NLOS conditions, depending on the blockage due to the objects distributed in the area.

The first set of results is reported in Fig. 7.8 and shows the time evolutions of the TCP congestion window size, the RTT, and the application throughput of UE1 when using TCP CUBIC[2] (red lines) and X-TCP (blue lines). It can be immediately seen where the cross layer approach gains with respect to the CUBIC algorithm. As shown in Fig. 7.8a, since the buffers and the retransmissions in the mmWave cellular stack mask most of the losses on the channel,

---

[2]The implementation of TCP CUBIC for ns–3 can be found at: `https://github.com/kronat/ns-3-dev-git/tree/tcp-versions-updated`

**(a)** Window　　　　　　　　　**(b)** RTT　　　　　　　　　**(c)** Throughput

**Figure 7.8:** Evolution over time of the congestion window, RTT and throughput of a X-TCP flow and of a TCP CUBIC flow for UE1, on the path shown in Fig. 7.7b (i.e., always NLOS).

TCP CUBIC is unaware of the actual rate provided by the lower layers and keeps increasing its congestion window, thus injecting more and more packets in the buffers at the lower layers. This yields an increase of the RTT, as shown in Fig. 7.8b. Conversely, X-TCP keeps a small congestion window, proportional to the actual rate supported by the channel. The throughput of both approaches is instead comparable, as shown in Fig. 7.8c. In this scenario, a limited data rate $R_{app,max}$ and large buffers were used. However, when the data rate is higher or the buffer size is smaller, TCP CUBIC may end up filling the whole buffer, thus causing an overflow and triggering an RTO. With X-TCP, instead, this would not happen, since the congestion window is adapted to the actual rate provided by the link.

To gain more insights on the performance of X-TCP and TCP CUBIC and on their mutual interactions, we ran a number of simulations by randomly changing the position and the number of obstacles in the area. To avoid any bias, each simulation has been run two times (i.e., resetting the pseudo-generator seed to the same value), but swapping the trajectories of the two UEs. Furthermore, we fixed the upper bound of the application-layer rate to $R_{app,max} = 2$ Gbps and $R_{app,max} = 1$ Gbps, which are larger and lower than the average bitrate of the mmWave links, respectively.

Fig. 7.9 shows the average RTT, buffer occupancy and throughput for three different combinations of TCP flavors: (i) both UEs use X-TCP, (ii) both UEs use TCP CUBIC, and (iii) one uses X-TCP and the other TCP CUBIC. For the first two cases, we show the time average of the mean performance of the two UEs, while for the latter case we present separately the results for each TCP flavor.

**RTT and RLC buffer occupancy:** as shown in Figs. 7.9a and 7.9b, the main advantage of the cross layer approach is the reduced latency (i.e., smaller RTT), which on average is half of that of TCP CUBIC. In particular, it can be seen in Fig. 7.9a that this behavior is consistent for the two different $R_{app,max}$ values, with a slightly higher RTT for $R_{app,max} = 2$ Gbps, and also across the three different simulated scenarios. TCP CUBIC experiences a higher RTT on average, because the congestion window continues to grow also in NLOS (as shown in Fig. 7.8a), and packets are queued in the RLC layer buffers because the capacity offered by the mmWave link is not enough to transmit all of them. This behavior is consistent with the results in Fig. 7.9b, where the average RLC buffer occupancy in the three scenarios is reported. The trend of this metric is similar to that of the RTT, suggesting that the RLC queueing is what causes the increase in the RTT of TCP CUBIC. X-TCP, instead, is able to adapt its congestion window to the actual rate available at the mmWave physical layer, therefore it limits the buffering at the RLC layer, only in the case of retransmissions triggered by unavoidable packet losses in the

165

**(a)** Average RTT  **(b)** RLC buffer occupancy  **(c)** Average throughput

**Figure 7.9:** RTT, buffer occupancy and throughput for different TCP configurations. The two leftmost bars represent scenarios with two TCP flows of the same flavor, i.e., TCP CUBIC or X-TCP. The two rightmost bars report the values for a scenario with TCP flows with different flavors, i.e., one is TCP CUBIC and the other is X-TCP. The errorbars represent the 95% confidence interval.

channel.

**Throughput and fairness:** it can be seen in Fig. 7.9c that, when the mmWave eNB is not saturated, then there is no difference in throughput among the three scenarios. In particular, in the third scenario the aggressiveness of X-TCP does not harm TCP CUBIC, since there are enough resources available to both. When the sum of the application rates exceeds the mmWave capacity, the results in the third scenario (two concurrent flows with different TCP) suggest that X-TCP may be unfair to TCP CUBIC. When flows use the same congestion control algorithm, then the average achievable throughput is around 600 Mbps, independently of the TCP flavor. However, when one UE uses X-TCP and the other uses TCP CUBIC, then the two flows do not split the available resources fairly, but the cross layer congestion control algorithm achieves a 23% higher throughput than TCP CUBIC. In particular, the TCP CUBIC flow in the third scenario loses 70 Mbps (i.e., 11%) with respect to the average throughput in the second scenario (TCP CUBIC only), while the UE using X-TCP gains 50 Mbps (i.e., 8.5%) with respect to the first scenario, when both UEs use X-TCP. This can be explained by the fact that, when a transition from NLOS to LOS happens, the cross layer approach restores full bandwidth utilization much more quickly than TCP CUBIC and this extra capacity is not easily released to the other flow. This has a negative impact also on the sum-rate (i.e., the sum of the throughput of the two flows), which decreases by 20 Mbps (i.e., 1.6%).

## 7.7 Conclusions

In this chapter we introduced milliProxy and X-TCP, two novel cross-layer schemes designed to enhance the performance of TCP in mmWave cellular networks. MilliProxy splits the TCP control loop in two segments, while keeping the end-to-end semantics of TCP. It has a modular design, which enables the use of different MSS values and flow window management algorithms in the two portions of the connection (i.e., wired and wireless). The window control policy can benefit from the interaction of milliProxy with the protocol stack of the mmWave networks, which enables cross-layer approaches. We showed how a FW policy based on the BDP of the end-to-end connection allows a reduction in latency of up to 10 times or an increase in goodput of up to 2 times with respect to traditional TCP NewReno, as well as a robustness with respect to where milliProxy is placed in the network. Similarly, X-TCP introduces a performance enhancement for

uplink flows in terms of buffer occupancy and, consequently, end-to-end latency, but introduces a decrease in fairness with respect to legacy flows

Cross-layer designs are an option to reach high goodput with low latency with TCP in mmWave networks. As part of our future work we will test the milliProxy and X-TCP performance in a wide variety of scenarios in ns-3, analyzing the performance with multiple users, and with different flow window policies, and will consider the implementation in a real setup.

<div style="text-align: right; font-size: 3em; color: #8B0000;">**8**</div>

# Alternative Multi-Connectivity-Based Transport Layer Solutions

A trend that has recently emerged in modern mobile networks is the exploitation of multiple network interfaces at the transport layer. Given the availability of different communications technologies in the same device (e.g., LTE, NR, Wi-Fi), multi-connectivity-based transport protocols can make the most out of the different characteristics of these connections and relay the application data over the different interfaces to decrease latency, and to increase the reliability and the throughput [298, 299]. This approach is different from the one considered in Chapter 3, where the different RATs are integrated in the RAN, because the integration happens at the two endpoints of the communication at a higher layer (i.e., the transport or application layer). Nonetheless, when it comes to mmWaves, it is beneficial to combine multiple links with different propagation characteristics (e.g., the good propagation at sub-6 GHz frequencies and the bandwidth available at mmWaves) [300]. Moreover, solutions based on multi connectivity at higher layers of the protocol stack make the deployment of multi connectivity independent of the choices of the network operator, and, as long as a final user can use multiple independent radios in its smartphone, an over-the-top content provider can exploit multi connectivity in its application (e.g., for video streaming).

In this chapter, which is based on [389, 404, 413], we investigate the performance on mmWave networks of two transport layer protocols which are an alternative to or an extension of TCP and are based on multi connectivity. The first is the multipath version of TCP, i.e., MPTCP, originally proposed in [139, 301, 302] and that we will evaluate for a combined mmWave (at 28 and 73 GHz) and LTE usage in Sec. 8.1. The second is a proposal of a network-coding-based application layer transport scheme that relies on UDP to provide high-quality, low-latency video streaming, and will be presented in Sec. 8.2.

## 8.1 Multipath TCP for mmWave and Sub-6 GHz Networks

### 8.1.1 Introduction on MPTCP

MPTCP has been proposed as a way of allowing vertical and seamless handovers between cellular networks and Wi-Fi hotspots and is currently under discussion for standardization at the Internet Engineering Task Force (IETF). It may also be used to provide path diversity in mmWave cellular

| Parameter | Value | Parameter | Value |
|---|---|---|---|
| mmWave carrier frequency | 28 GHz, 73 GHz | LTE carrier frequency (UL) | 1.9 GHz |
| mmWave bandwidth | 1 GHz | LTE bandwidth | 20 MHz |
| mmWave TX power | 30 dBm | LTE downlink TX power | 30 dBm |
| LTE carrier frequency (DL) | 2.1 GHz | LTE uplink TX power | 25 dBm |

**Table 8.1:** Simulation parameters

networks. The three main design goals of MPTCP are [302]:

1. Improve throughput: an MPTCP flow should perform at least as well as a traditional Single Path TCP (SPTCP) flow on the best path available.

2. On shared links, MPTCP should not get more resources than standard TCP flows.

3. MPTCP should prefer less congested paths, subject to the previous two conditions.

There are three Request for Commentss (RFCs) that describe MPTCP [139, 301, 302]. They discuss the signaling and setup procedures [301], the architectural choices for the deployment of MPTCP [139], and a CC algorithm [302]. Finally the document in [303] discusses the impact on the application layer.

There are several studies that propose coupled congestion control algorithms for MPTCP connections. By coupling over the different subflows, the authors of [302] claim that it is possible to reach goals 2 and 3 above. In particular they propose a first coupled CC, that is however criticized in [304] and in [145], because it (i) transmits too much traffic on congested paths and (ii) is unfriendly with respect to SPTCP. Therefore two more coupled CC were proposed:

- In [304] the Opportunistic Linked Increase Algorithm (OLIA) is designed to overcome these two issues, but presents non-responsiveness problems with respect to congestion changes in the subflows;

- In [145] the BALIA addresses both the problems of the original CC and those of OLIA. In particular, the parameters of the protocol are derived through a theoretical analysis of the performance of multipath congestion control algorithms.

However, these schemes are based on the legacy design of Reno and New Reno congestion control algorithms (Additive Increase - Multiplicative Decrease, AIMD), which are shown to suffer from the highly dynamic behavior of mmWave links more than the newer TCP CUBIC congestion control algorithm, as discussed in Chapter 6. Other CC algorithms have been recently proposed for MPTCP, we refer the interested reader to our survey in [396].

### 8.1.2 Performance Evaluation on mmWaves

MPTCP could be used as an end-to-end solution for multi-connectivity, i.e., next generation mobile devices may connect both to an LTE and to a mmWave eNB, or to two or more mmWave eNBs with no need for coordination at the lower layers. However, there are some issues with its performance in mmWave networks, as we will show in the following paragraphs. In this performance evaluation campaign, based on the DCE framework described in Sec. 2.11, we used the real Linux implementation of MPTCP (v0.90), which includes several CC algorithms, namely the original coupled CC, OLIA, BALIA, uncoupled (with any desired TCP flavor, e.g., CUBIC), and others. We co-deploy an LTE eNB and a mmWave eNB, or a mmWave eNB capable of transmissions at different frequencies (28 and 73 GHz, with the same bandwidth and

**(a)** MPTCP throughput for different distances $d$ and different MPTCP options. The black dotted line shows the performance of a SPTCP connection with TCP CUBIC, as a reference.

**(b)** Contribution of the two subflows as a function of the distance $d$, for MPTCP with CUBIC CC.

**Figure 8.1:** Throughput for MPTCP.

the maximum number of antennas available in the ns–3 mmWave module), and vary the distance of the multi-connected UE from the eNBs, using the channel model described in [37]. The remote host is a multi-homed server, supporting MPTCP connections. The UE uses IPERF, and starts the connection on the 28 GHz mmWave link. Then another subflow is added on the LTE link, or on the 73 GHz mmWave link. The main parameters for the simulations are reported in Table 8.1.

Fig. 8.1a shows the performance in terms of throughput of different MPTCP congestion control algorithms over different connections, with respect to the baseline of a SPTCP connection with TCP CUBIC. The dashed lines represent a scenario with paths on LTE and on mmWave (28 GHz), while the solid ones refer to paths on mmWave links with 28 GHz and 73 GHz as carrier frequencies.

**LTE as mmWave secondary path:** When the UE is close to the eNB and has a LOS link most of the time on both the 28 and the 73 GHz connections (e.g., for $d = 50$ m), then the solution with multipath TCP on mmWave-only links outperforms SPTCP, with a gain that ranges from 800 Mbps (28%) to 1 Gbps (36%). Instead, due to the limit of the LTE uplink, the performance of a multipath on LTE and mmWave is close to that of SPTCP (when CUBIC is used, because BALIA has much worse performance, as will be discussed later).

However, it can be seen from Fig. 8.1a that MPTCP with LTE and mmWave links performs better than with only mmWave connections for $d \geq 100$ m, and with the CUBIC uncoupled CC algorithm also for $d = 75$ m. Indeed, the 73 GHz link offers a potentially larger throughput than an LTE uplink connection, but it has a lossy behavior that penalizes the overall throughput, except for small distances. In particular, for $d = 150$ m, MPTCP with LTE and 28 GHz mmWave offers a gain of more than 450 Mbps (i.e., 100%) with respect to the SPTCP (i.e., more than the LTE uplink throughput), showing that the presence of the secondary and reliable LTE path improves the throughput on the mmWave link. This can be seen also in Fig. 8.1b, where we plot the contribution of the two subflows of MPTCP connections at $d \in \{100, 150\}$ m when the second subflow is LTE or mmWave. It can be seen that the contribution given by the reliable LTE uplink subflow is smaller than that of the 73 GHz mmWave subflow, but the primary 28 GHz mmWave subflow reaches a higher throughput when coupled with the LTE secondary subflow.

171

**Figure 8.2:** Download time as a function of the file size and of the distance, for MPTCP with CUBIC CC and secondary subflow on LTE or mmWave at 73 GHz.

For short-lived TCP sessions, instead, using a secondary subflow on mmWave links improves the system performance. This can be seen in Fig. 8.2, which shows the download time of a file using wget. However, the performance gain, especially for smaller files, is minimal, showing that the LTE link makes up for its smaller capacity with a higher reliability that benefits the performance of TCP.

**Coupled vs uncoupled CC**: Another important observation is that MPTCP with the BALIA CC algorithm fails to meet target 1, since in many cases its throughput is lower than that of SPTCP, as shown in Fig. 8.1a. The most striking cases are those with MPTCP on LTE and mmWave, and $d \in \{50, 75\}$ m. Here the congestion control algorithm sees the losses on the 28 GHz mmWave link as congestion, and, according to design goal 3, it steers the whole traffic to the LTE subflow, degrading the performance of the end-to-end connection. Instead, the uncoupled congestion control algorithm is not affected by this issue, since each path behaves independently. However, in this case design goal 2 is not met.

When considering short-lived TCP sessions and file download times, there are two different outcomes according to the file size. As shown in Fig. 8.3a, when the file is smaller than 1 MB the BALIA coupled congestion control algorithm exhibits a slightly smaller download time than the CUBIC uncoupled CC. Instead, when the file is larger than 5 MB, as in Fig. 8.3b, the MPTCP



**(a)** Small file size.



**(b)** Large file size.

**Figure 8.3:** Download time as a function of the file size and of the distance, for MPTCP with secondary subflow on LTE and CUBIC or BALIA CC.

172

solution with CUBIC as CC mechanism manages to download the file in less than a fifth of the time required by MPTCP with BALIA. This behavior can be explained by considering the shape of the window growth function of CUBIC, which recalls a cubic function, i.e., flat at the beginning and then rapidly increasing.

## 8.2    Network Coding and Multi Connectivity with UDP on mmWaves

In this section, we describe the protocol we introduce in [413], which aims at providing high quality (live) video streaming by combining the high data rates at mmWaves with reliability, low packet loss, a stable data rate, and low latency. The proposed solution exploits (i) multi-connectivity between the LTE and mmWave RANs, to provide continuous coverage with LTE and high capacity with mmWave, and (ii) network coding, in order to simplify the management of the transmission on multiple links and provide additional robustness. We evaluate the performance in terms of packet loss, latency and video quality using a novel framework that combines for the first time the ns-3 mmWave simulator described in Chapter 2 with real video traces and a network coding library [305]. The results verify that the proposed solution provides a high level of video quality with low delays.

### 8.2.1    State of the Art

In the literature, there are both research results and commercial products for indoor applications of video streaming at mmWaves, with a limited range and based on either proprietary technologies [306] or IEEE 802.11ad [185]. In [307], the 60 GHz band is shown as a candidate for the transmission of uncompressed, high quality video up to 3 Gbps. In an indoor environment, mmWave links can also be configured to stream virtual reality content from a local server to the headset [308]. However, the evaluation of the end-to-end performance of applications in mmWave cellular networks is a research area still in its infancy, given the lack of large mmWave cellular network testbeds or deployments.

In conventional LTE cellular networks, network coding has been studied as an enabler of high quality video streaming. The authors of [309] propose to use it as an error correction technique in the LTE RAN, and re-design the MAC layer in order to use network coding instead of the traditional HARQ mechanism for multimedia traffic. A similar proposal for WiMAX can be found in [310]. In [311], network coding was shown to increase the efficiency of resource allocation when used for video broadcasting in LTE.

To the best of our knowledge, the combination of multi connectivity and network coding at mmWaves for video streaming has not been studied yet. Multi connectivity was studied in [312] to satisfy the quality of service constraints of video streaming, but without network coding and at sub-6 GHz frequencies. On the other hand, a packet-level encoding technique similar to network coding (i.e., the Luby Transform codes) was used on top of UDP on a mmWave link in [313], to increase the connection goodput with respect to TCP, but not for video streaming and without multi connectivity.

### 8.2.2    Video Streaming Framework

In this section, we present the proposed framework for video streaming in a mmWave cellular network scenario. As shown in Fig. 8.4, the proposed protocol that handles the intelligent distribution policy operates from the Video Streaming Server (VSS), which can be deployed either in the operator's core network as a caching server, or in the public internet. A middle layer (called video distribution layer) that manages network coding, any retransmissions, and the

**Figure 8.4:** Video Streaming Framework.

multiple interfaces to the different RANs is placed between the encoding layer, which generates video frames, and the transport layer. Both UDP and TCP have been used as transport protocols for video distribution: for example the DASH protocol [54] relies on TCP, while the Real Time Streaming Protocol (RSTP) [314] can operate on both. In this framework we consider UDP for two reasons. The first is that the reliability typically offered by TCP is provided in our architecture by network coding at the middle layer, and the second is related to the limitations of the TCP performance on mmWave links, as discussed in Chapter 6.

In the next paragraphs, we describe how each of the three components of our solution (i.e., multi connectivity, network coding and the video transmission policy) are engineered to yield the best performance for the final user.

### Multi Connectivity

The possibility of using multiple network interfaces at the same time is an emerging paradigm in wireless and data center communications [299]. In particular, the modern smartphones are generally equipped with multiple radios and network interfaces. In this context, we use multi connectivity at the application layer to communicate using different RATs, such as LTE at sub-6 GHz frequencies and NR at mmWave frequencies, in order to benefit from (i) a more reliable end-to-end packet transmission on LTE when the mmWave link is not available, and (ii) the very high data rate of the mmWave connection when the link quality is high enough. Moreover, the LTE connection is also used to send feedback messages from the UE to the VSS to signal the availability and the quality of the mmWave link.

### Network Coding

Network coding is a packet-level encoding technique that combines source packets using algebraic operations in order to increase the resilience with respect to packet loss in an efficient way [315]. In our architecture, we rely on the rateless version of Random Linear Network Coding (RLNC) [316], as it provides a good trade-off between bandwidth efficiency, complexity and delay, compared to other network coding or forward error correction strategies [317]. The network coding library chosen for this section is Kodo [305].

With RLNC, the packets generated by the video encoding layer are grouped into generations, i.e., sets of $K$ packets meant to be encoded together, where $K$ is the generation size. For each generation, coded packets are obtained as independent random linear combinations of the $K$ packets, where symbols, coefficients and all operations are defined in the Galois field with $q$ elements, $\mathbb{F}_q$. As a result, every encoded packet is an equally useful representation of the packets from the generation, such that the decoder is able to decode the original information using any combination of (slightly more than) $K$ encoded packets. The number of encoded packets that

174

can be generated from $K$ packets, i.e., the RLNC code rate is not fixed. If some packets are lost on the mmWave link, it is possible to produce newly encoded packets without re-encoding and retransmitting the whole generation. This is the rateless property of the encoding scheme.

When an encoded packet is produced, it can be immediately transmitted. The decoder collects encoded packets, and needs to receive at least $K$ packets to attempt a successful decoding. At the decoder side, the original packets are retrieved through Gaussian elimination, by constructing a decoding matrix with the linearly independent encoded packets that have been successfully received. Since the encoding coefficients are randomly chosen, it is not guaranteed that each encoded packet will be linearly independent of the others, and thus that the original payloads will be re-constructed given $K$ encoded packets. In order to increase the decoding probability, in our design we send $N \geq K$ encoded packets and start decoding on-the-fly as soon as $K$ packets are received.

There are trade-offs between (i) the latency and the decoding probability, which both increase with the generation size $K$, and (ii) the decoding complexity and the decoding probability, which both increase with the field size $q$ [318]. We test two different configurations: configuration $LC$ ($K = 40$, $q = 4$), which offers low latency and decoding overhead at the cost of a lower probability of successful decoding; and configuration $HC$ ($K = 100$, $q = 8$) where the latency, overhead, and also the probability of successful decoding are increased.

Finally, network coding simplifies the management of multi connectivity, since the retransmissions do not need to be performed on the path in which the lost packet was originally transmitted, but the best available one can be used when needed. In order to protect from both unsuccessful decoding and packet loss on the wireless link, $N$ is set to respectively $1.2K$ or $1.1K$ when the mmWave or the LTE link is used. Transmissions of additional, newly encoded packets can be triggered, up to a maximum number of 5 attempts.[1]

Video Encoding Policy

The video is encoded using the H.264/Advanced Video Coding (AVC) standard [319] with the Scalable Video Coding (SVC) extension [320]. This framework provides the possibility of avoiding the transmission of some portions of the video bit stream in order to adapt the source rate to the channel capacity or to the needs of the end users: this property has been referred to as *scalability*. The source content can be divided into subsets with a reduced picture size (spatial scalability) or lower frame rate (temporal scalability). In the time domain, it is possible to identify key frames that will carry most of the content, and enhancement frames that are placed between two key frames and can be discarded (with a loss of quality). The different kinds of frames belong to different temporal layers. The key frames are part of the temporal base layer, and two of these frames together with a set of enhancement frames form a Group of Pictures (GOP). In the spatial domain, the scalability makes it possible to code two or more versions of the same video at different resolutions in a unique bit stream, which is therefore composed of different layers corresponding to different spatial resolutions (i.e., a spatial layer). According to the H.264/AVC standard, the bit stream generated by the encoder is divided into Network Abstraction Layer Units (NALUs), each with a payload containing a portion of the encoded video frame. Each NALU is then split into packets of size $P = 1000$ bytes, which are forwarded to the network coding layer. According to the Non Overlapping Window (NOW) policy [321], the network coding layer maps packets of different NALUs into different generations, so that the encoding is independent for each NALU.

---

[1]Strictly speaking, the proposed scheme is not fully rateless, as the number of generated encoded packets has an upper limit due to latency constraints.

In this chapter, we consider a 50 Hz video with GOP of 16 frames [320], 5 temporal layers, and 2 spatial layers at a resolution of 720p (base layer) and 1080p (enhancement layer).

### 8.2.3  Performance Evaluation

#### Simulation Setup

The performance evaluation is carried out using the ns-3 mmWave module, which was integrated with the Kodo network coding library [305] and several tools to process the video traces.

We implemented the elements of the framework of Fig. 8.4 (e.g., the Video Streaming Server), as well as the protocol to manage multi connectivity and network coding. The scenario contains a single cell with radius equal to 100 m and 5 users, 2 in LOS and 3 in NLOS, which move at a random speed between 2 and 4 m/s around fixed positions in the cell. The main parameters of the simulations are given in Table 8.2.

In order to provide a realistic video streaming model, the chosen video sample is first encoded in the format specified in Sec. 8.2.2, from which a bit stream is then generated using the JSVM software [322]. Using our extension of the tool provided in [323], the bit stream is adapted to the processing in ns-3. The NALUs are then handled by the video distribution layer and transmitted in the simulation, and, at the UE, the correctly received frames are first buffered and then played-out. The play-out action in the simulation corresponds to writing the frame-related information in an output trace, which is then processed with the tool in [323] in order to allow the video reconstruction and quality evaluation with FFmpeg [324]. The video buffer considered in the simulation has a memory of 25 frames, i.e., 500 ms of video.

#### Results

The metrics considered in this performance evaluation, obtained via Monte Carlo simulations with 90 independent runs, are (i) the NALU loss ratio; (ii) the application layer latency, i.e., the delay between the time at which the video frame is generated at the VSS and when it is consumed by the application at the UE; and (iii) the average frame Peak Signal to Noise Ratio (PSNR). The PSNR is a measure of the quality of reconstructed video that is inversely proportional to the Mean Square Error (MSE) of the received frame $R$ with respect to the original frame $I$. Given the frame width $W$ and height $H$ in pixels, the PSNR for frame $n$ is given by [322]:

$$PSNR(n) = 10 \log_{10} \frac{WH(2^8 - 1)^2}{\sum_{w=1}^{W} \sum_{h=1}^{H} [I_n(w, h) - R_n(w, h)]^2}. \tag{8.1}$$

| Parameters | Value | Parameters | Value |
|---|---|---|---|
| LTE carrier frequency (DL) | 2.1 GHz | 3GPP Channel Scenario | Urban Micro |
| LTE carrier frequency (UL) | 1.9 GHz | mmWave SNR Outage Threshold | -5 dB |
| LTE bandwitdh | 20 MHz | RLC buffer size $B_{RLC}$ | 20 MB |
| LTE downlink TX power $P_{TX}$ | 43 dBm | RLC reordering timer | 1 ms |
| mmWave carrier frequency | 28 GHz | RLC Buffer Status Report timer | 2 ms |
| mmWave bandwidth | 1 GHz | Number of UEs | 5 (2 LOS, 3 NLOS) |
| mmWave $P_{TX}$ | 30 dBm | VSS-UE latency | 10 ms |

**Table 8.2:** Main simulation parameters

**(a)** Latency for different configurations of the packet transmission policy. The y-axis starts from 10 ms, i.e., the value of the VSS-UE latency.

**(b)** NALU loss ratio for different configurations of the packet transmission policy.

**Figure 8.5:** Performance evaluation for the video streaming framework. NC FEC/RAN retx means that network coding error correction/RLC and HARQ retransmission are used, no NC FEC/no RAN retx otherwise.

When there are no differences between the reconstructed video frames and the original ones, the PSNR as in JSVM is assigned the maximum value of 99.99 dB.

Fig. 8.5a and Fig. 8.5b compare the end-to-end latency and the NALU loss ratio of different configurations. In particular, two different kinds of error correction to compensate for the packet losses are considered, via link-level retransmissions in the RAN and/or by transmissions of additional NC packets (henceforth denoted by NC FEC).

Consider the blue bars in Fig. 8.5a and Fig. 8.5b, i.e., the cases in which multi connectivity is not used. When comparing the schemes without NC FEC, and without or with RAN retransmission, i.e., the first and the third group of bars, it can be seen that the RAN retransmissions do not significantly increase the latency (Fig. 8.5a), while enabling a 25% reduction in the NALU loss (Fig. 8.5b) for the HC configuration. Instead, when comparing the options without RAN retransmissions, and without or with NC FEC, i.e., the first and second group of bars, it can be seen that NC FEC reduces the NALU loss by 38% (HC option) or 36% (LC option): it is more efficient for reducing NALU loss than the RAN retransmissions only. However, the latency increases much more when applying NC FEC than RAN retransmissions (up to 50% for the HC configuration, as shown in Fig. 8.5a). The link-level retransmissions are indeed more efficient with respect to single-packet losses in the channel, while NC FEC protects larger chunks of packets and can yield a lower NALU loss in case of more extended bursty errors.

Nonetheless, the best performance in terms of NALU loss when multi connectivity is not used is achieved when combining both RAN retransmissions and NC FEC (fourth group of bars in Fig. 8.5b), at the price of a modest increase in latency (comparable to that achieved by applying only NC FEC for the HC configuration, and 25% higher for the LC option, as shown by the fourth group of bars in Fig. 8.5a).

Multi connectivity, however, is the configuration that performs best both for the latency and the NALU loss ratio, as shown by the comparison between the blue bars (no multi connectivity) and the red bars for all the configurations. In particular, multi connectivity makes it possible to continuously transmit packets even when the mmWave link is in outage, thanks to the LTE fall back and to the seamless switch enabled by the fact that the UE is already connected to both RANs. Therefore, as shown by the red bars in Fig. 8.5a, the average value of the latency when RAN retransmissions and/or NC FEC are introduced does not increase significantly with

177

**Figure 8.6:** PSNR of the spatial base layer vs. latency for the configuration with no RAN retransmission and no NC FEC (1) and the one with both RAN and NC FEC (2).

respect to the case in which neither of them is used (i.e., the first group of bars), and is in general less than 2 ms higher than the 10 ms delay introduced by the fixed backhaul network. Additionally, the NALU loss with multi connectivity, RAN retransmissions and NC FEC has very small values (in the order of $10^{-5}$, as shown by the fourth group of bars in Fig. 8.5b), given that a more reliable LTE link is used when the mmWave one is in outage.

Finally, Fig. 8.6 shows the trade-off between the average end-to-end latency and the average PSNR of the spatial base layer.[2] It can be seen that the PSNR without error control (points marked with "1") is limited to about 26 dB because of the relatively high NALU loss ratio, whereas the combination of RAN retransmissions and network coding error correction (points marked with "2") is able to guarantee almost perfect reconstruction (note that 100 dB is the conventional value given by JSVM to error-free packet delivery). The figure also shows that multi connectivity, while not necessarily needed for high PSNR (which can be achieved even by a stand-alone mmWave network), can be very effective in reducing latency, especially in the presence of error control (points marked with "2"), with a reduction in delay of up to 40%, as also shown in the fourth group of bars in Fig. 8.5a.

## 8.3 Conclusions

This chapter has investigated two different multi-connectivity-based transport layer options at mmWave frequencies, with the goal of understanding which alternatives to TCP can yield a good performance on mmWave links.

The performance analysis of MPTCP for mmWave networks showed that at larger distances and for long-lived TCP sessions it is preferable to use a more stable LTE-like link, and that the deployment of MPTCP coupled congestion control algorithms on mmWave links is not able to satisfy the original design goals of [302]. A possible improvement of MPTCP CC algorithms should adapt the TCP CUBIC scheme to a coupled scenario, so that the reactiveness and stability of TCP CUBIC enhance the performance of the transport protocol while not harming other legacy TCP flows.

Moreover, in the second part of this chapter, we highlighted how it is possible to deploy reliable and low-latency video streaming on mmWave links. The proposed video streaming framework is based on a combination of network coding and multi connectivity with LTE, which are managed by a middle layer between the application and the transport layers. The performance evaluation

---

[2]Note that FFmpeg does not support spatial scalability for video decoding, therefore it is possible to reconstruct only the spatial base layer at 720p. The evaluation of the PSNR for the combined spatial base and enhancement layers is left for future work.

of the proposed solution (based on real video traces) confirms the benefit introduced by multi connectivity, and shows that network coding can help reduce the NALU loss and increase the PSNR, especially when the mmWave-only solution is used.

As future work, we will consider and investigate the performance on mmWave links of other transport protocols that have recently emerged as alternatives to TCP, such as, for example, QUIC. We provide and describe an implementation of a QUIC module for ns-3 in [423].

**Part IV**

# The Intelligence: Data-Driven 5G Networks Optimization

<div align="right">

# 9

</div>

# Machine Learning at the Edge

## 9.1 Introduction

As mentioned in Chapter 1, the fifth generation of cellular networks is being designed to satisfy the massive growth in capacity demand, number of connections and the evolving use cases of a connected society for 2020 and beyond [1]. In particular, 5G networks target the following KPIs: (i) very high throughput, in the order of 1 Gbps or more, to enable virtual reality applications and high-quality video streaming; (ii) ultra-low latency, possibly smaller than 1 ms on the wireless link, to support autonomous control applications; (iii) ultra-high reliability; (iv) low energy consumption; and (v) high availability of robust connections [6, 325].

In order to meet these requirements, a new approach in the design of the network is required, and new paradigms have recently emerged [6]. First, the densification of the network will increase the spatial reuse and, combined with the usage of mmWave frequencies, the available throughput. On the other hand, this will introduce new challenges related to mobility management [42]. Second, with MEC, the content will be brought closer to the final users, in order to decrease the end-to-end latency [6]. Third, a higher level of automation will be introduced in cellular networks, relying on Machine Learning (ML) techniques and SDN, in order to manage the increased complexity of 5G networks.

The usage of ML and Artificial Intelligence (AI) techniques to perform autonomous operations in cellular networks has been widely studied in recent years, with use cases that range from optimization of video flows [326] to energy-efficient networks [327] and resource allocation [328]. This trend is coupled with the application of big-data analytics that leverage the huge amount of monitoring data generated in mobile networks to provide more insights on the behavior of networks at scale [329]. In the domain of mobile networks, these two technological components can empower costs savings, but also new applications, as we will show in this and the following chapters. In particular, in this chapter we will focus on the proposal of an architecture to practically deploy intelligent and machine-learning-based algorithms in a 5G networks, and on how data-driven techniques can enable self-organizing approaches for 5G.

### 9.1.1 Contributions

Despite the interest of the industrial and research communities towards the deployment of machine learning in networks, the state of the art lacks considerations on how it is possible to

effectively deploy intelligence in cellular networks, and an evaluation of the gains of a data-driven approach with real large-scale network datasets.

To address these limitations, in this chapter, which is based on [392, 398, 420], we propose a data-driven architecture for the practical implementation of ML techniques in 5G cellular networks, and evaluate the gains that this architecture can introduce in some data-driven applications, using real data collected from hundreds of base stations of a major U.S. carrier in the San Francisco and Mountain View areas for more than a month. In particular, the main contributions related to this topic are:

- the design of a scalable and efficient multi-layer edge-based architecture to deploy big-data and ML applications in 5G systems. We propose to exploit controllers implemented in MEC and cloud facilities to collect the data generated by the network, run analytics and extract relevant metrics, that can be fed to intelligent algorithms to control the network itself and provide new services to the users. The RAN controllers, deployed at the edge, are associated with a cluster of base stations, and are thus responsible not only for RAN control, as proposed in [330], but also for running the data collection and ML infrastructure. The network controller, placed in the operator's cloud, orchestrates the operations of the RAN controllers. We characterize this architecture with respect to the latest 5G RAN specifications for 3GPP NR, the 5G standard for cellular networks [7], and provide insights on how the controllers can interface with an NR deployment, following the approach of an emerging open RAN initiative contributed by multiple operators and vendors [330].

- the demonstration of the gains that data-driven techniques enabled by the proposed architecture can yield in network applications, leveraging a real world dataset on two use cases. In the first, big data analytics are used to control the association between the base stations and the RAN controllers. We propose a *dynamic clustering* method where base stations and controllers are grouped according to the day-to-day user mobility patterns, which are collected and processed by the ML infrastructure. With respect to a static algorithm, based on the position of the base stations, the data-driven algorithm manages to decrease the number of inter-controller interactions and thus reduce the control plane latency. In the second example, we test different machine learning techniques (i.e., the Bayesian Ridge Regressor, the Gaussian Process Regressor and the Random Forest Regressor) for the *prediction* of the number of users in the base stations of the network. We show that, thanks to the proposed ML edge-based architecture, which makes it possible to exploit the spatial correlation of the users, it is possible to increase the prediction accuracy with respect to that of decentralized schemes, with a reduction of the prediction error by up to 53%.

To the best of our knowledge, this is the first exhaustive contribution in which a practical ML architecture, that can be applied on top of 5G NR cellular networks, is evaluated using a real network dataset, showing promising results that indicate that new user services and optimization techniques based on machine learning in cellular networks are possible.

Finally, we also evaluate how data from external sources (e.g., sensors deployed in smart cities) can help optimize the network itself, following a data-driven paradigm. Building upon the "SymbioCity" concept proposed in [331], in the final part of this chapter (based on [392]), we exploit the traffic data from the Transport for London (TfL) Urban Traffic Control (UTC) network [332] in order to dynamically optimize network parameters such as the number of virtualized MMEs deployed in a certain market (e.g., in London). Since handovers will be one of the major issues in 5G ultra-dense networks, the techniques we propose will decrease the operating costs for an operator without compromising the handover completion time. The ability to choose the point in the tradeoff between cost and performance is going to be a key element in the design of self-organizing 5G networks, and data-driven techniques will play a fundamental role in this.

**Table 9.1:** Relevant literature on machine learning, MEC and edge controllers in cellular networks and novel contributions of this study.

| Topic | Relevant References | Contribution of this study |
|---|---|---|
| Application of ML in cellular networks | [13, 329, 333–337] | Novel network-level architecture, integrated with 3GPP 5G specifications, and evaluation of its performance gains based on real network dataset. |
| Mobility prediction in cellular networks | [338–340] | Cluster-based approach to capture spatial correlation |
| Mobile Edge Cloud | [6, 341–344] | MEC-based architecture used for ML for network control and applications |
| SDN in cellular networks | [12, 330, 345–348] | ML-driven edge-SDN controllers integrated in the ML architecture |
| NFV in cellular networks | [252, 349–353] | Data-driven optimization of the NFV function deployment |

## 9.1.2 Related Work

In the following paragraphs we will discuss the literature relevant to the scope of this study, which is also summarized in Table 9.1, and highlight the main differences we introduce with respect to the state of the art.

**ML in cellular networks** The application of ML techniques to cellular networks is a topic that has gained a lot of attention recently, thanks to the revived importance of ML and AI throughout all facets of the industry. The surveys in [13, 333] present some recent results on how it is possible to apply regression techniques to mobile and cellular scenarios in order to optimize the network performance. The paper [334] gives an overview of how machine learning can play a role in next-generation 5G cellular networks, and lists relevant ML techniques and algorithms. The usage of big-data-driven analytics for 5G is considered in [329,335], with a discussion of how data-driven approaches can empower self-organizing networks. However, none of these papers provides results based on real operators datasets at large scale that show the actual gains of data-driven and machine learning based approaches. Moreover, while practical implementations of machine learning algorithms for networks indeed exist for host-based applications (e.g., TCP [336], video streaming [337]), or base-station-based use cases (e.g., scheduling [354]), the literature still lacks a discussion and an analysis of how it is possible to practically deploy the algorithms, collect real-time data and process it to enable new services in large-scale commercial networks.

Furthermore, several papers report results on the prediction of mobility patterns of users in cellular networks. The authors of [338,339] use network traces to study human mobility patterns, with the goal to infer large-scale patterns and understand city dynamics. The paper [340] proposes to use a leap graph to model the mobility pattern of single users. Other works focus on the prediction of the traffic generated by single base stations [355,356], or by groups of base stations [357], and do not consider the mobility patterns. With respect to the state of the art, in this study we focus on the prediction of the number of users associated to a base station, in order to provide innovative services to the users themselves, and propose a novel cluster-based approach to improve the prediction accuracy, evaluating the performance of different algorithms on a real large-scale dataset.

**MEC and controllers in cellular networks**   The role of MEC has also been discussed in the context of 5G networks, e.g., to perform coordination [341] and caching [342], and to offer low-latency content and control applications to the end users [6]. MEC is indeed considered a key element in the deployment of future autonomous driving vehicles, for which very short control loops will be needed [358]. A few papers consider specific cases for the application of machine learning and big data techniques at the edge, for example for intelligent transportation systems [343], or the processing of data collected by internet-of-things devices [344], but, to the best of our knowledge, the usage of MEC to run data collection and machine learning algorithms for the prediction and optimization in 5G cellular networks has not been discussed in detail yet.

The edge has also been proposed for hosting controllers in cellular networks [330, 345, 346]. As the SDN paradigm has become popular in wired networks [359], several software-defined approaches for the RAN have been described in the literature [12, 347, 348], and the telecom industry is moving towards open-controllers-based architectures for the deployment of 5G networks [330]. With respect to existing studies, in this work we propose to exploit the RAN controllers as proxies for the data collection in the RAN and the enforcement of machine learning algorithm-based policies. This approach has been considered in a wired-network context [360], but this is the first study that applies it in a 5G cellular network.

**Network Function Virtualization and Virtual MME**   Finally, the other main architectural trend in the evolution towards 5G is NFV: instead of using specialized and costly hardware in both the core and the access network, most of the processing is virtualized and run on general-purpose machines in the cloud [349]. This allows a larger flexibility and adaptability to the instantaneous load of the control and user planes. The initialization cost of a new Virtual Machine (VM) is orders of magnitude smaller than the cost of the equivalent worst-case dimensioned hardware. A broad overview of the issues and other potential benefits of NFV is presented in [252]. Although this research is still ongoing, preliminary results [352] show that it is possible to increase the energy efficiency of the network without significant performance losses.

In the last part of this chapter, we focus on handover management in virtualized MMEs. A first model of the performance of the different virtualized CN functions is presented in [351], and the MME is identified as a critical element for scalability of control plane functionalities. Virtualization can also enable distributed MME designs [350]. An optimized design of a virtualized MME is given in [353], where the number of vMME instances is adapted to the traffic load in an Machine to Machine (M2M) scenario, using a traffic model for CN-related events. With respect to the state of the art, this study exploits a data-driven approach to drive the optimization.

### 9.1.3   Chapter Structure

The remainder of the chapter is organized as follows. In Sec. 9.2 we present the real network dataset used in the first part of the chapter, and in Sec. 9.3 we describe the proposed architecture. In Sec. 9.4 we provide details on the first application, i.e., the autonomous data-driven clustering of base stations. Results on the second application, i.e., the prediction accuracy for the number of users in the cells, are given in Sec. 9.5, together with possible use cases. The data-driven NFV optimization approach based on vehicular traffic data is discussed in Sec. 9.6. Finally, in Sec. 9.7 we conclude the chapter.

**(a)** Utilization (averaged over a 15-minute interval).



**(b)** Number of active UEs (summed over a 15-minute interval).

**Figure 9.1:** Example of timeseries from the traces collected for 4 eNBs in the Palo Alto dataset over 5 days.

## 9.2 The Dataset

This section describes the data that will be used in the evaluations in the first part of the chapter. The traces we exploit are based on the monitoring logs generated by 650 base stations of a national U.S. operator in two different areas, i.e., San Francisco and Palo Alto/Mountain View, for more than 600000 UEs per day, properly anonymized during the collection phase. The base stations in the dataset belongs to a 4G LTE-A deployment, which represents the most advanced cellular technology commercially deployed at a large scale. Even if 5G NR networks will have more advanced characteristics than LTE ones, this dataset can be seen as representative of an initial 5G deployment at sub-6 GHz frequencies in a dense urban scenario [361]. We consider two separate measurement campaigns, conducted in February 2017 in the San Francisco area and in June and July 2018 in the Palo Alto and Mountain View areas. Table 9.2 summarizes the most relevant details of each measurement campaign.

Given the sensitivity of this kind of data, we adopted standard procedures to ensure that individuals' privacy was not compromised during the data collection and the analysis. In particular, the records were anonymized by hashing the UEs' International Mobile Subscriber Identitys (IMSIs), which is the unique identifier that can be associated to a single customer in these traces. Moreover, for our analysis, we only used anonymized metrics that are based on aggregated usage at multiple layers: first, we consider users' data for each single cell (a cell is mapped to a sector and carrier frequency), and, then, aggregate the cells associated to the same base station (i.e., with the RF equipment in the same physical location). In this way, no user can be singled out by the results we present.

The traces used for this study record a set of standardized events in LTE eNBs, mainly related to the mobility of users. The raw data is further processed to construct time series of

| | Location | Time interval | Number of eNBs |
|---|---|---|---|
| Campaign 1 | San Francisco | 01/31/2017 − 02/26/2017, every day from 3 P.M. to 8 P.M. | 472 |
| Campaign 2 | Palo Alto, Mountain View | 06/22/2018 − 07/15/2018, whole day | 178 |

**Table 9.2:** Anonymized datasets used in this study.

different quantities of interest in each eNB at different time scales (from minutes to weeks): (i) the utilization of the eNB, which is represented by the ratio of used and available Physical Resource Blocks (PRBs); (ii) the number of incoming and outgoing handovers, for both X2 and S1 handover events [259]; and (iii) the number of active UEs, obtained from context setup and release events. The measurement framework we used also offered the possibility of logging other events and extract other metrics, for example related to the latency experienced by the users, link statistics (e.g., error probability), or different estimates of the user and cell throughput. The events associated to these quantities, however, are reported less regularly and less frequently than those we consider, therefore they do not represent a reliable source for the estimation of the network performance. With respect to other publicly available datasets [362], this presents a more precise characterization of the mobility dynamics in the network and a finer granularity in the collected data.

Fig. 9.1 shows an example of different timeseries for 4 eNBs in the Mountain View/Palo Alto area, with a time step of 15 minutes. It can be seen that, even though daily patterns can be identified, each eNB presents characteristic differences with the others.

## 9.3   RAN Controllers as Enablers of Machine-Learning Applications at the Edge

The past and current generations of cellular networks were not designed to deploy machine learning and artificial intelligence algorithms at scale. The main reason is that there are no standardized interfaces that network operators can exploit to collect data from the base stations and the equipments of different vendors, and/or to modify the behavior of the network according to custom policies. Indeed, despite the Self-Organizing Network (SON) capabilities embedded in the LTE standard [259], the deployment of autonomous networks is not widespread, and LTE eNBs are usually self-contained appliances to which the telecom operators have restricted access. Therefore, the control plane is usually decentralized, and the exchange of information among eNBs is limited [330]. Accordingly, practical machine learning solutions that can deployed in a 4G LTE network are generally limited to SON parameters optimization for a few eNBs, generally with offline training and/or optimization, thus without real-time insights, or to the application of intelligent algorithms to the data that is collected in each single eNB, for example to predict the channel gain [403], perform smart handovers [363] or scheduling [328, 354].

In order to make network management and operation more efficient, new design paradigms have emerged in the 5G domain. The main trend is related to the disaggregation of the base station (which in 3GPP NR networks is the gNB). The 3GPP has proposed different splits of the gNB protocol stack [7], so that it will be possible to deploy a different RAN architecture, with the lower layers in DUs on poles and towers, and the higher layers in CUs which can be hosted in a datacenter. The pooling of CUs can enable more sophisticated orchestration operations, and energy savings [12]. On the other hand, the DUs that are deployed in the RAN are simpler and possibly smaller than 4G full-fledged base stations.

The second trend is related to the deployment in the wireless RAN of SDN solutions based on open and smart network controllers [364], which have already been adopted with success in large wired backbone networks [359]. Along this line, the O-RAN Alliance, a consortium of network operators and equipment vendors, is standardizing controller interfaces between the CUs and new custom RAN controllers that can be implemented and deployed by the telecom operators themselves. As mentioned in [330], an architecture with a split between the distributed hardware that performs data-plane-related functions and a more centralized software-based control plane can enable more advanced control procedures, thanks to the centralized view and the context awareness, and thus this approach is quickly becoming a de facto standard for the deployment

**Figure 9.2:** Proposed controller architecture for RAN control and machine learning at the edge.

of 5G cellular networks.

### 9.3.1 Proposed Architecture

In this chapter, and in [398], we propose to exploit the new design paradigms for the 5G RAN to make it possible to practically deploy intelligence in cellular networks, without the constraints and limitations previously described for 4G LTE deployments. As shown in Fig. 9.2, our architecture leverages the different layers of controllers to aggregate and process the network data using machine learning and AI techniques, with a multi-layer semi-distributed point of view that strikes a balance between the decentralized 4G approach and a completely centralized approach, which would be infeasible due to the amount of data to be processed.

In the following paragraphs, we will introduce the proposed architecture and describe how it can be integrated in the NR and O-RAN Alliance designs, following MEC paradigm. Moreover, we will discuss the costs and the technical challenges associated to the deployment of the proposed architecture. In Sec. 9.4 and Sec. 9.5 we will describe two ML-based applications for networks, showing that the usage of the proposed architecture makes it possible to improve the performance with respect to decentralized, 4G-based approaches.

**Integration with 3GPP networks**

The proposed architecture exploits a multi-layer overlay that is compliant with 3GPP NR networks, as reported in Fig. 9.2. The overlay is composed by three main elements:

- the *RAN*, which is deployed to provide cellular service to the users, and includes the 3GPP NR CUs and DUs. The RAN handles the data plane of the users, i.e., the user traffic is forwarded from or to the core network and the public Internet from the CUs [7].

- the *RAN controllers*, which control and coordinate the RAN elements, as proposed in [330]. Each RAN controller is associated to a cluster of gNBs, and is deployed in MEC, to minimize the communication latency with the RAN. Some of the control-plane processes are assigned to the RAN controllers, which can benefit from the cluster-based overview.

For example, as proposed in [330], the RAN controllers can manage UE-level connectivity, by coordinating handover decisions and performing load balancing, or can enforce QoS policies.

- the *Cloud Network Controller*, that orchestrates the RAN controllers (e.g., to establish the RAN controllers/gNBs association) and provides application-layer services, and can be deployed in a remote cloud facility.

A multi-layer controller architecture combines the benefits of the scalability of a distributed approach with the performance gain given by a partially-centralized view of the network. Each layer implements control functionalities with different latency constraints, allowing the network to scale: the DUs schedule over-the-air transmissions on a sub-ms basis, the RAN controllers may decide upon users' association on a time scale of tens of milliseconds, and, finally, the Cloud Network Controller can operate on multiple-second (or even longer) intervals, for example to update the association between gNBs and RAN controllers. At each additional layer, it is possible to support a larger number of devices (e.g., a DU controls tens of UEs at most, while the RAN controller can be designed to handle hundreds of UEs), and, given the more relaxed constraints on the decision time scale, it is possible to implement more refined and complex decision policies, based on machine learning algorithms enabled by the larger amount of data given by the clustered and/or centralized views.

**RAN Controllers, Machine Learning and Data Collection**

While the RAN controllers are generally deployed to perform the aforementioned control plane task, we propose to leverage them to implement machine learning techniques at the edge of the network. A network operator can indeed use the proposed overlay to manage the data collection from the distributed gNBs and enforce policies based on the learning applied to this data. Notice that, for some metrics, the controllers would not need explicit signaling for the data collection: for example, if a controller manages the UEs sessions, as proposed in [330], then it is already aware of the number of users connected to each gNB it controls.

The position of the RAN controllers in the overlay network strikes a balance between the breadth of their point of view, the amount of data they need to collect and process, and the number of the user sessions they can handle. In general, as the number of base stations associated to a controller grows (and, consequently, the number of controllers decreases, up to a single controller), it is possible to perform more refined optimizations, given that the knowledge of the state of the network is more complete. However, there is a limit to how much the data collection can be centralized. Indeed, if the operator is interested in running *real-time* data-driven algorithms, for example to decide upon the association of UEs and gNBs, then we argue that a completely centralized architecture does not scale because of (i) the amount of data (for example, related to channel measurements) that needs to be collected and (ii) the collection and processing delay. In this regards, we observed that it is not possible to perform a real-time collection and processing of a subset of the monitoring data streamed from the Palo Alto/Mountain View network (178 base stations) in a single virtual machine with 8 x86 CPUs at 2.1 GHz. On the other hand, a completely distributed approach (as in a 4G LTE network) cannot exploit *any* centralized view and/or enforce coordinated policies, as previously mentioned, and, as we will show in Sec. 9.5 with real network data, does not perform as well as the controller-based architecture for the accurate prediction of the number of users in the network.

190

**Technical Challenges**

The usage of RAN controllers, however, introduces new technical challenges. First, new standard interfaces and signaling between the gNBs and the controllers will need to be defined.[1] For example, in a completely distributed architecture (e.g., LTE), for a handover there is a message exchange between neighboring base stations, and, then, the core network [259], while, if controllers are used, the gNBs can interface directly with their controller to exploit its global view. Once the actual specifications for RAN controllers will be completed, it will be possible to also evaluate the signaling difference among these different architectures.

Another interesting problem is related to the association of controllers and gNBs. This issue has already been studied for SDN controllers in wired networks [365], but wireless cellular networks have characteristics that introduce new dimensions to this problem, mainly related to the higher level of mobility of the endpoints of such networks, i.e., the UEs.[2] If the RAN controllers are used to manage user sessions and mobility events, then they will need to maintain a consistent state for each user associated to the gNBs they control. Given that cellular users often move through the area covered by the cellular networks, it becomes of paramount importance to minimize the number of times a user performs a handover between two base stations controlled by different controllers. In this case, indeed, the two controllers would need to synchronize and share the user's state, and this would increase the control plane latency, as also observed in case of inter-controller communications in wired SDN networks [367]. Therefore, in the following section, we will describe a practical data-driven method to perform the association between gNBs and controllers, testing the proposed algorithm on the San Francisco and the Mountain View/Palo Alto datasets.

## 9.4  Big-data Driven RAN Controller Association

In the remainder of this chapter we introduce our second major contribution, i.e., we describe two applications related to network control and optimization that show the advantages of using the proposed controller-based architecture described in Fig. 9.2. In particular, in this section, we illustrate a data-driven approach for the control-plane association of RAN controllers and gNBs. The algorithm we designed aims at minimizing the number of interactions between gNBs belonging to different RAN controllers (since any controller that is added in the control loop severely impacts the control plane latency), and enables a dynamic allocation of the base stations to the different controllers. Moreover, it is based on the real data that the network itself can collect, thus it represents another example of how it is possible to exploit real-time analytics to self-optimize the performance.

### 9.4.1  Proposed Algorithm

Our method is based on a semi-supervised constrained clustering on a graph weighted according to the transition probabilities among base stations. The algorithm is summarized with the pseudocode in Alg. 9.2. The input is represented by the timeseries of X2 and S1 handovers for all the $N_g$ gNBs in the set $\mathcal{B}$, each tagged with the timestamp of the event and the pair $< source, destination >$ gNBs, and by the time step $T_c$ to be considered for the computation of the transition probability matrices (e.g., fifteen minutes or a day). Moreover, the network operator

---

[1] This effort is being pursued, among others, by the O-RAN Alliance [330]

[2] Notice that in this study we consider a control-plane gNB-controller association, i.e., the controller is not involved in the processing of data-plane packets and low-level scheduling, which is what is instead usually considered in the design of controllers for interference coordination problems [366].

---

**Algorithm 9.2** Network-data-driven RAN Controller Association Algorithm

---

1: **for** every time step $T_c$
2:     **distributed data collection step (performed in each RAN controller):**
3:       **for** every RAN controller $p \in \{0, \ldots, N_c - 1\}$ with associated gNBs set $\mathcal{B}_p$
4:         **for** every gNB $i \in \mathcal{B}_p$
5:           compute the number of handovers $N_{i,j}^{\text{ho}} \forall j \in \mathcal{B}$
6:         **end for**
7:       report the statistics on the number of handovers to the Cloud Network Controller
8:     **end for**
9:     **clustering and association step (performed in the Cloud Network Controller):**
10:       compute the transition probability matrix $H$ based on the handovers between every pair of gNBs
11:       define weighted graph $G = (V, E)$ with weight $W(G)_{i,j} = H_{i,j} + H_{j,i}$
12:       perform spectral clustering with constrained K means on $G$ to identify $N_c$ clusters
13:       apply the new association policy for the next time step
14: **end for**

---

**Algorithm 9.3** Graph spectral clustering algorithm with constrained K means

---

1: **input:** graph $G = (V, E)$ with weights $W(G)$
2: compute the degree matrix $D_{i,i} = \sum_{j=1}^{N_g} W(G)_{i,j}$
3: compute the normalized Laplacian of $G$ as $L = I - D^{-1} W(G)$
4: create the matrix $U \in \mathbb{R}^{N_g \times N_c}$ with the eigenvectors of $L$ associated to the $N_c$ smallest eigenvalues as columns

5: apply constrained K means on the rows of $U$ to get $N_c$ clusters

---

can tune the number of RAN controllers $N_c$ according to the availability of computational resources and the number of base stations and related UEs that each controller can support.

Every $T_c$, each RAN controller $p \in \{0, \ldots, N_c - 1\}$, which has collected the timeseries of events for its gNB $i$ in the set of controlled gNBs $\mathcal{B}_p$, will process this data to extract the number of handovers $N_{i,j}^{\text{ho}}, \forall i \in \mathcal{B}_p, \forall j \in \mathcal{B}$, and will report this information to the Cloud Network Controller described in Sec. 9.3.1. The Cloud Network Controller then aggregates the statistics from each RAN controller and builds a complete transition probability matrix $H$, where entry $(i, j)$ is

$$H_{i,j} = \begin{cases} \dfrac{N_{i,j}^{\text{ho}}}{\sum_{j=1}^{N_g} N_{i,j}^{\text{ho}}} & \text{if } \sum_{j=1}^{N_g} N_{i,j}^{\text{ho}} \neq 0, \\ 0 & \text{otherwise.} \end{cases} \tag{9.1}$$

Then, consider the fully-connected undirected graph $G = (V, E)$, where $V = \mathcal{B}$ is the set of $N_g$ vertices, and $E$ is the set of edges that represent possible transitions among the gNBs. Each edge $e_{i,j}$ is weighted by the sum of the transition probabilities between gNBs $i$ and $j$, i.e., $W(G)_{i,j} = H_{i,j} + H_{j,i}$, with $W(G)$ the weight matrix, to account for all the possible transitions (and thus interactions, and, possibly, message exchanges and state synchronizations) between the two gNBs. In order to identify the set of gNB-to-controllers associations that minimize the inter-controller communications, the proposed algorithm clusters the undirected graph $G$ to identify the groups of gNBs in which the intra-cluster interactions (i.e., handovers and transfer of user sessions) are more frequent than inter-cluster ones.

We tested and considered different approaches for the clustering [368,369], which, in this case, has to satisfy two constraints: (i) the number of clusters should be an input of the algorithm, to match the number of available controllers[3]; and (ii) the size of the clusters (i.e., number of gNBs per cluster) should be balanced, to avoid overloading certain controllers while under-utilizing

---

[3]Notice that in this case finding the optimal solution to the clustering problem is NP-hard, thus identifying the optimal solution is not feasible in large scale networks [370].

others. The first constraint rules out popular unsupervised graph clustering techniques based on community detection algorithms, which are also generally applied to directed graphs [371]. Therefore, we propose to use a variant of standard spectral clustering techniques for graphs [372], which relies on a constrained version of K-means to balance the size of the clusters. Alg. 9.3 lists the main steps of the procedure.

Consider the degree matrix $D \in \mathbb{R}^{N_g \times N_g}$, i.e., a diagonal matrix with an entry $D_{i,i} = \sum_{j=1}^{N_g} W(G)_{i,j}$ for each gNB $i \in 1, \dots, N_g$. Then, it is possible to compute the normalized graph Laplacian as $L = I - D^{-1}W(G)$ and extract the eigenvectors associated to the $N_c$ smallest eigenvalues, i.e., as many eigenvalues as the number of clusters to identify. The result is a matrix $U \in \mathbb{R}^{N_g \times N_c}$ with the eigenvectors as columns. Each row of this matrix, which corresponds to a specific gNB, can be considered as a point in $\mathbb{R}^{N_c}$, and can be clustered using K means [372]. Standard K means, however, does not generate balanced clusters. Therefore, we replace this last step with a constrained K means algorithm, which modifies the standard K means by adding constraints on the minimum and maximum size of the clusters during the cluster assignment step. In this way, the cluster assignment problem can be formulated as a linear programming problem [373]. The final result is a set of $N_c$ clusters, and the Cloud Network Controller can apply the clustering policy to assign the gNBs to the respective RAN controllers.

### 9.4.2 Evaluation with Real Data

We compare the proposed network-data-based strategy (whose results are reported in Fig. 9.3a for the San Francisco area and Fig. 9.3b for the Mountain View area) with a baseline, in which the constrained K means is directly applied to the latitude and longitude of the gNBs (shown in Figs. 9.3c and 9.3d, respectively). Indeed, several approaches have been proposed in the literature to cluster, for example, remote radio heads and Base Band Units (BBUs) into BBU pools, according to different targets [374–376]. However, none of these focuses on the minimization of the control plane latency, but rather on data-plane issues, such as the minimization of interference or coordinated multipoint transmissions. Therefore, as a baseline, we consider the basic clustering approach based on the geographical position of the base stations. This method is static, and can be applied in networks that do not rely on data-driven approaches for configuration purposes, for example because the operator does not collect and/or make use of real-time network analytics. In the absence of this kind of data, we argue that geographic clustering is an approach in line with the goal to minimize inter-controller interactions, given that users are expected to move among neighboring base stations, which the geographical clustering will group under the same RAN controller.

Fig. 9.3a reports an example of the clustering applied to the $N_g = 472$ San Francisco base stations, with $N_c = 22$ clusters and $T_c = 24$ hours, i.e., with one clustering update per day, using the data collected in the previous day. The size of the clusters is constrained in $\{0.8N_g/N_c, \dots, 1.2N_g/N_c\}$. By comparing Figs. 9.3a and 9.3c, it can be seen that network-based clustering maintains a proximity criterion (i.e., base stations which are close together are generally clustered together), but this is not as strict as in the geographical one. Consider for example the base station at the bottom right of the figures: it serves an area close to U.S. Route 101, and public transportation stations, thus there are a lot of handovers happening directly from base stations in the downtown area to that gNB. Consequently, the network-based approach clusters it with the purple cluster in the city center, while the position-based strategy associates it to the other base stations at the bottom of the map. In general, it can be seen that in Fig. 9.3c there are more large black lines connecting the gNBs, meaning that base stations with a high level of interactions are placed under different controllers in different clusters. Another example of this can be seen in the comparison between Figs. 9.3b and 9.3d for the transitions

**(a)** Clustering with Alg. 9.2 in San Francisco.



**(b)** Clustering with Alg. 9.2 in Mountain View.



**(c)** Clustering with the positions of the gNBs in San Francisco.



**(d)** Clustering with the positions of the gNBs in Mountain View.

**Figure 9.3:** Network-data- and position-based clusters in San Francisco, using data from 2017/02/01 with $T_c = 24$ hours and $N_c = 22$, and Mountain View/Palo Alto, with data from 2018/06/28 with $T_c = 24$ hours and $N_c = 10$. The colored dots represent the base stations, with different colors associated to different clusters. The lines connecting the dots represent the weights in the graph $G$ of the edge between the two gNBs, with a thicker line representing a larger weight, i.e., sum of transition probabilities between the gNBs. Finally, lines with the same color as the dots represent edges between vertices in the same cluster, and vice versa for black lines.

along the Caltrain railway line that crosses the map on the diagonal. In Fig. 9.3b, most of the lines along the railway are colored, showing that intra-cluster handovers happen between the interested base stations, and vice versa in Fig. 9.3d.

In order to further compare the location-based, static clustering and that obtained from the network data, we compare the number of intra- and inter-controller handovers as a function of the number of controllers[4] (and thus clusters) $N_c$ and the time interval between two consecutive updates $T_c$. As mentioned in Sec. 9.3, intra-controller handovers can be managed locally, by the controller which is in common to the source and target base stations. Inter-controller handoffs, instead, require the coordination and synchronization of the two controllers, thus increasing the

---

[4]The number of controllers an operator will need to deploy on a network will depend on the capacity of the controllers themselves and the signaling they will need to support.

**(a)** San Francisco scenario, 2017/02/02.

**(b)** Mountain View/Palo Alto scenario, 2018/06/28.

**Figure 9.4:** Ratio $R$ between intra- and inter-cluster handovers as a function of the number of clusters $N_c$, with clustering based on daily updates.

control plane latency to at least twice that of handovers related to a single controller. The actual overhead on the latency introduced by inter-controller communications will depend on signaling specifications that have not been developed yet, and on the controller implementation and processing capabilities, as mentioned in Sec. 9.3, but the need to avoid inter-controller synchronization is valid in any case. Therefore, we report as metrics the number of intra- and inter-controller handovers and their ratio.

In Fig. 9.4 we present the ratio $R$ between intra- and inter-cluster handovers by considering $T_c = 24$ hours as fixed, and changing the number of clusters $N_c$. For each value of $N_c$, we run multiple times the clustering algorithms, to average the behavior of K means and provide confidence intervals. It can be seen that the gain of the network-data-based solution over the position-based one is almost constant, especially as the number of clusters grows, with an average increase of the ratio $R$ of 45.38% for the San Francisco case and 42.62% for the Mountain View/Palo Alto scenario. The behavior in the two scenarios with $N_c = 2$, however, is different: while in the San Francisco case $N_c = 2$ yields the largest difference for the value of $R$ between the network-data- and the location-based clustering, in the Mountain View context it corresponds to the minimum difference. This is due to the difference in the geography of the two areas, as shown in Fig. 9.3: the San Francisco dataset covers a much larger number of base stations than the other one, and the mobility patterns of the users are less regular, thus the clustering based on the network data can find a better solution than that based on location.

Finally, in Figs. 9.5a and 9.5b, we report the number of handovers for the two configurations shown in Fig. 9.3, with $T_c = 24$ hours, and for a more dynamic solution based on more frequent updates (i.e., $T_c = 15$ minutes). Moreover, Figs. 9.5c and 9.5d also plot the ratio between the intra- and inter-cluster handovers. Notice that the number of handovers reported in Fig. 9.5a refers to the events happened on February 2nd, while the clustering is based on the data from the previous day. For the 15-minute update case, the clustering is updated every 15 minutes to reflect the statistics from the previous 15 minutes. However, as Fig. 9.5a shows, updating the clusters with a daily periodicity, using data from the previous day, does not result in significantly degraded performance with respect to the 15-minute updates case. Notice also that a cluster update has some cost in terms of control signaling between the gNBs and the controllers. Moreover, the daily-based update builds the graph and the clustering according to a more robust statistics, i.e., based on the transitions for the whole day. This is particularly evident if we consider the example in Figs. 9.5b and 9.5d, which report the same metrics but for a whole day in the Mountain View/Palo Alto area and $N_c = 10$ clusters. As it can be seen, at night, when the number of handovers is low, the clustering with update step $T_c = 15$ minutes exhibits a very high variation in the ratio between intra- and inter-cluster handovers, and in some cases has a performance which is similar to that of the geographic case, while the curve for the daily-based

195

**(a)** Number of intra- and inter-cluster handovers for 2017/02/02 in San Francisco, $N_c = 22$.

**(b)** Number of intra- and inter-cluster handovers for 2018/06/28 in Mountain View, $N_c = 10$.

**(c)** Ratio between intra- and inter-cluster handovers for 2017/02/02 in San Francisco, $N_c = 22$.

**(d)** Ratio between intra- and inter-cluster handovers for 2018/06/28 in Mountain View, $N_c = 10$.

**Figure 9.5:** Number of intra- and inter-cluster handovers (and relative ratio $R$) with different clustering strategies, in different deployments (i.e., San Francisco, with 472 base stations, and Mountain View/Palo Alto, with 178).

update shows a more stable behavior and better performance.

To summarize, we showed that the data-driven clustering based on the proposed architecture (i) adapts to the mobility of users, in different scenarios, thus reducing the inter-controller interactions and, consequently, the control plane latency, and (ii) can be updated on a daily basis without significant performance loss with respect to a more dynamic solution.

## 9.5 Predicting Network KPIs Using Controllers

In this section, we present an additional application of the ML architecture presented in Sec. 9.3, in which the point of view of the RAN controllers is exploited to predict the number of users attached to each base station of the cellular network. This metric can be used to forecast useful KPIs such as the user throughput, the outage duration and the overall network load. In the following paragraphs, we will first discuss the quality of the prediction with several machine learning algorithms by considering a single cluster among those presented in Fig. 9.3a for San Francisco. The main comparison will be between the accuracy of the prediction with (i) methods that only use local information, i.e., in which each base station is a separate entity (as in 4G) and has available only its own data for the training of the machine learning algorithm, and (ii) techniques that exploit the architecture described in Sec. 9.3 to collect and process data, and thus for which it is possible to perform predictions based on the joint history of multiple base stations associated to each controller. Then, we will extend the analysis to all the clusters, using the most promising approaches identified for the test-cluster, showing how a cluster-based

| Regression method | Hyperparameters |
|---|---|
| Bayesian Ridge Regressor [377, 378] | $\alpha \in \{10^{-6}, 10^{-3}, 1, 10, 100\}$, $\lambda \in \{10^{-6}, 10^{-3}, 1, 10, 100\}$ |
| Random Forest Regressor [379, 380] | Number of trees $N_{rf} \in \{1000, 5000, 10000\}$ |
| Gaussian Process Regressor [381] | $\alpha \in \{10^{-6}, 10^{-4}, 10^{-2}, 0.1\}$, $\sigma_k \in \{0.001, 0.01\}$ |

**Table 9.3:** Values of the hyperparameters of the different regressors for the k-fold cross-validation.

approach reduces the prediction error with respect to a local-based approach. Finally, we will describe some prediction-based applications for network automation and new user services.

### 9.5.1 Data Preprocessing

The performance analysis presented in this section is based on the San Francisco dataset. We sampled the number of users in each base station with a time step $T_s = 5$ minutes, and divided the dataset into a training set (which will be used for k-fold cross validation) and a test set. The training set is based on the interval from January 31st to February 20th, while the test set goes from February 21st to February 26th.

For base station $i \in \mathcal{B}$, with $\mathcal{B}$ the set of base stations in San Francisco, consider a multi-step ahead prediction of the number of users $N_u^i(t + L)$ at times $t + 1, \ldots, t + L$ (where $L \geq 1$ is the *look-ahead* step of the prediction), given the real-time data before time $t$. The features we identified are (i) the past $W$ samples of the number of users (where $W$ is the window of the history used for the prediction), i.e., $N_u^i(t + \tau), \tau \in [-W + 1, 0]$; (ii) an integer $h(t) \in \{0, \ldots, 4\}$ that represents the hour of the day (from 3 P.M. to 8 P.M.); and (iii) a boolean $b(t)$ that indicates whether the selected day is a weekday. We also tested the cell utilization and the number of handovers as possible features, however they showed small correlation with the prediction target. For each day, given the discontinuities of the collected data, we discard the first $W$ samples, thus the actual size of the training ($N_{tr}$) and test ($N_{te}$) sets depends on the value of $W$.

For the local-based prediction, in which each base station predicts the future number of users based on the knowledge of its own data, the training and test set are composed by the feature matrix $\mathbf{X} \in \mathbb{R}^{N_i, 3W}, i \in \{tr, te\}$, in which each row is a vector $[N_u^i(t - W + 1), h(t - W + 1), b(t - W + 1) \ldots, N_u^i(t), h(t), b(t)]$, and by the target vector $\mathbf{y} \in \mathbb{R}^{N_i, 1}, i \in \{tr, te\}$. For the cluster-based method, instead, the goal is to predict the vector of the numbers of users for all the base stations in the cluster. Therefore, for the set $\mathcal{C}_d = \{i_d, \ldots, j_d\} \subset \mathcal{B}$ with the $N_b^d$ base stations of cluster $d$, each row of the target matrix $\mathbf{Y} \in \mathbb{R}^{N_i, N_b^d}, i \in \{tr, te\}$ is a vector $[N_u^{i_d}(t + L), \ldots, N_u^{j_d}(t + L)]$. The feature matrix $\mathbf{X} \in \mathbb{R}^{N_i, W(N_b^d + 2)}, i \in \{tr, te\}$ is composed in each row by a vector with the form $[N_u^{i_d}(t - W + 1), \ldots, N_u^{j_d}(t - W + 1), h(t - W + 1), b(t - W + 1), \ldots, N_u^{i_d}(t), \ldots, N_u^{j_d}(t), h(t), b(t)]$.

The values of the numbers of users in the training and test sets are transformed with the function $\log(1 + x)$ and scaled so that each feature assumes values between 0 and 1. The scaling is fitted on the training set, and then applied also to the test set. For the evaluation of the performance of the different methods and prediction algorithms, we use the Root Mean Squared Error (RMSE), defined for a single base station $i$ as $\sigma_i = \sqrt{1/N_{te} \sum_{t=1}^{N_{te}} (y_i(t) - \hat{y}_i(t))^2}$, with $y_i$ the time series of the real values for the number of users for base station $i$, and $\hat{y}_i$ the predicted one.

### 9.5.2 Algorithm Comparison

We tested several machine learning algorithms tailored for prediction, i.e., the Bayesian Ridge Regressor (BRR) for the local-based prediction, and the Gaussian Process Regressor (GPR) and

| Look-ahead step $L$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| BRR | 6 | 6 | 4 | 4 | 3 | 3 | 3 | 2 | 2 |
| cluster-GPR | 3 | 2 | 2 | 2 | 2 | 1 | 6 | 5 | 4 |

**Table 9.4:** Values of $W$ for the plot in Fig. 9.6b for the BRR and the cluster-based GPR

Random Forest Regressor (RFR) for both the local- and the cluster-based predictions, using the implementations from the popular open-source library scikit-learn [382].[5] For each of these methods, we considered different values of $W \in \{1, \ldots, 10\}$ and predicted at different future steps $L \in \{1, \ldots, 9\}$, i.e., over a time horizon of 45 minutes. 3-fold cross-validation was performed for each method and value of $L$ and $W$ to identify the best hyperparameters, among those summarized in Table 9.3. The split in each fold is done using the `TimeSeriesSplit` of scikit-learn, i.e., without shuffling, and with increasing indices in each split, to maintain the temporal relation among consecutive samples.

The BRR (which is used for urban traffic prediction in [378]) combines the Bayesian probabilistic approach and the ridge $L_2$ regularization [377]. The Bayesian framework makes it possible to adapt to the data, and only needs the tuning of the parameters $\alpha$ and $\lambda$ of the Gamma priors. However, it does not generalize to multi-output prediction, thus we applied this method only to the local-based scenario.

The RFR (used in [380] for population prediction) is a classic ensemble method that trains $N_{rf}$ regression trees from bootstrap samples of the training set and averages their output for the prediction [379]. The only hyperparameters to be tuned are (i) the number of trees $N_{rf}$, for which a higher value implies better generalization properties, but also longer training time; and (ii) the number of random features to sample when splitting the nodes to build additional tree branches, which is set to be equal to the number of features for regression problems. It supports prediction of scalars and vectors, thus we tested it with both the local- and the cluster-based approaches.

Finally, the GPR is a regressor that fits a Gaussian Process to the observed data [381]. The prior has a zero mean, and the covariance matrix described by a kernel. In this case, we chose a kernel in the form

$$k(x_i, x_j) = \sigma_k^2 + x_i \cdot x_j + \left(1 + \frac{d(x_i, x_j)^2}{2\alpha l^2}\right)^{-\alpha} + \delta_{x_i x_j}, \tag{9.2}$$

i.e., the sum of a dot product kernel, that can model non-stationary trends, a rational quadratic kernel with $l = 1$ and $\alpha = 1$, and a white kernel, that explains the noisy part of the signal. The GPR can be used for both single-output and multi-output regressions.

### 9.5.3 Performance analysis for a sample cluster

For the comparison between the aforementioned regressors, we consider the cluster $d = 0$ with $N_d^0 = 22$ base stations in the San Francisco area. We assume that the cluster is stable throughout the training and testing period. In a real deployment, when the base station association to the available controllers changes, a re-training will be needed, together with additional signaling between the controllers, to share the data related to the base stations whose association was updated.

---

[5] An approach based on neural networks was also considered, but, due to the reduced size of the training set, underperformed with respect to the other regression methods.

**(a)** $W = 1$

**(b)** A different window $W$ is selected for each method and look-ahead step $L$ to minimize the RMSE $\hat{\sigma}$. The values of $W$ are reported in Table 9.4.

**Figure 9.6:** RMSE $\hat{\sigma}$ for different local- and cluster-based prediction methods, as a function of the look-ahead step $L$, and for different windows $W$.



**(a)** High number of users

**(b)** Low number of users

**Figure 9.7:** Example of predicted vs true time series, for $L = 3$ (i.e., 15 minutes ahead), $W = 3$ and the cluster-based GPR on two base stations for cluster 0.

In order to compare the local- and the cluster-based methods, we report in Fig. 9.6 the average RMSE $\hat{\sigma} = \mathbb{E}_{i \in \mathcal{C}_0}[\sigma_i]$ of the base stations in the set $\mathcal{C}_0$ associated to cluster 0. As expected, the RMSE increases with the look-ahead step $L$. Among the local-based methods, the BRR gives the best results for all the values of the look-ahead step $L$, with a gain of up to 18% and 55% with respect to the GPR and RFR for $L = 9$. The GPR, instead, is the best among the cluster-based techniques, with an improvement up to 50% from the RFR (for $L = 1$). When comparing the local- and the cluster-based methods, the latter performs better, especially as the look-ahead step increases, since the curve of the RMSE for the cluster-based GPR flattens around $\hat{\sigma} = 14.8$, while that for both the BRR and the local-based GPR continues to increase. In this case, instead, for small values of $L$ the performance of local- and cluster-based methods is similar.

Table 9.4 reports the values of the window $W$ used in Fig. 9.6b for the two best performing methods, the BRR and the GPR. By comparing Figs. 9.6a, in which the window $W$ is fixed, and 9.6b, where $W$ is selected for each step $L$ to yield the smallest RMSE $\hat{\sigma}$, it can be seen that the difference is minimal for the best performing methods (i.e., below 5%), while it is more significant for the local-based RFR. Moreover, the spatial dimension has more impact on the quality of the prediction than the temporal one. Indeed, while by changing $W$ the RMSE for the GPR and BRR improves by up to 5%, when introducing the multi-output prediction with the GPR the RMSE decreases by up to 50%. Differently from prior works in which the single user mobility is predicted [340], we are indeed considering the number of users at a cell level, and, in this case, the possible transitions between neighboring cells are limited by the geography of the scenario, and by the available means of transport. Therefore, there exists a spatial correlation

**(a)** RMSE $\hat{\sigma}$ of the cluster-based GPR on cluster 0 when varying the amount of data used for training, at different future time steps $L$.

**(b)** Residual error $N_u(t) - \hat{N}_u(t)$, where $N_u(t)$ is the true value of the number of users at time $t$, and $\hat{N}_u(t)$ is the predicted one, as a function of the true value of the number of users $N_u(t-1)$ at time $t-1$. $L = 2$.

**Figure 9.8:** Additional results on the prediction accuracy for cluster 0 with the cluster-based GPR, $W = 2$.

between the number of users in the neighboring base stations and the number of users in the considered base station at some time in the future, given that the mobility flows are constrained by the aforementioned factors.

Nonetheless, there exist still some limitations to the accuracy of the prediction of the number of users. Fig. 9.7 reports an example of the predicted (for $L = 3$, i.e., 15 minutes) and the true time series for two different base stations, with a high and low number of users. As it can be seen, the true time series have some daily patterns, but are also quite noisy. As a consequence, the predicted time series manage to track the daily pattern, but cannot predict the exact value of the number of users. This is more evident when the number of UEs is low, as in Fig. 9.7b, which also exhibits smaller daily variations.

Finally, Fig. 9.8 reports additional results on the prediction performance of the cluster-based GPR. In Fig. 9.8a, we compare the RMSE $\hat{\sigma}$ obtained on the testing dataset when using partial training datasets of different sizes, i.e., with 25, 50, 75 hours, or the complete training dataset (i.e., 100 hours). The RMSE monotonically decreases as the size of the training dataset increases, showing that there is room for improvement with a richer past history. Moreover, the difference is more marked when considering a higher prediction lag $L$, i.e., the full training dataset yields an RMSE which is 25% smaller than the 25-hours dataset for $L = 1$ and 40% for $L = 5$.

Fig. 9.8b shows an example of residual analysis, which can help understand the limits of the cluster-based GPR on the available San Francisco dataset. The y-axis reports the residual error $N_u(t) - \hat{N}_u(t)$, with $N_u(t)$ and $\hat{N}_u(t)$ the true and predicted number of users at time $t$, and the x-axis one of the features used in the prediction, i.e., the true number of users $N_u(t-1)$ at the previous time step $t-1$. Notice that the x-axis is quantized into 100 bins in order to improve the visualization of the residuals. It can be seen that the largest errors happen (infrequently) on the left part of the plot, i.e., when there is a sudden increase in the number of users in the base station, transitioning from a small $N_u(t-1)$ to a large $N_u(t)$.

### 9.5.4 Performance analysis for the other clusters

Given the promising results of the cluster-based approach on the first cluster, we selected the best performing local- and cluster-based methods, i.e., respectively, the BRR and the GPR, and performed the prediction on all the clusters reported in Fig. 9.3a. The results are reported in Fig. 9.9 for each single cluster. The cluster-based method always outperforms the local-based one, and, in most cases, also exhibits a smaller RMSE for small values of the look-ahead step $L$,

**Figure 9.9:** Cluster-based GPR vs local-based BRR for the other clusters.

contrary to what happens for cluster 0. The reduction in the average RMSE over all the clusters $\mathbb{E}_{clusters}[\hat{\sigma}]$ is 18.3% for $L = 1$ (from $\mathbb{E}_{clusters}[\hat{\sigma}] = 7.24$ to $\mathbb{E}_{clusters}[\hat{\sigma}] = 6.11$) and increases up to 53% for $L = 9$ (from $\mathbb{E}_{clusters}[\hat{\sigma}] = 17.42$ to $\mathbb{E}_{clusters}[\hat{\sigma}] = 11.34$).

### 9.5.5 Possible Applications

The results presented in Figs. 9.6 and 9.9 show that the cluster-based method is more capable than local-based ones to capture the user dynamics in the cellular network. The prediction of the number of users in a base station can be used to optimize the performance of the network in a number of different ways: for example, it can enable predictive load-balancing, bearer pre-configuration, scaling of RAN resources, sleeping periods for base stations, and so on. We believe that the increase in the prediction accuracy that the cluster-based method yields can be beneficial to practically enable these anticipatory and prediction-based optimizations.

Moreover, network operators can exploit the prediction to offer novel services to the end users. For example, consider a vehicle that has to travel from point $A$ to point $B$ in an area covered by cellular service. While on the journey, the passengers may want to participate in a conference call, or, if not driving, surf the web or stream multimedia content. Therefore, given the choice of multiple routes with similar Estimated Times of Arrival (ETAs), the passengers may prefer to choose an itinerary with a slightly higher ETA but with a better network performance, because, for example, it crosses an area with a better coverage, or with fewer users. This becomes

| | Feb. 23rd, 19:00 | | | | Feb. 24th, 19:00 | | | | Feb. 24th, 19:20 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Route | R1 | R2 | R3 | R4 | R1 | R2 | R3 | R4 | R1 | R2 | R3 | R4 |
| $\hat{S}$ [Mbps] | 1.93 | 2.51 | 2.36 | 2.74 | 1.72 | 2.00 | 2.28 | 2.89 | 2.05 | 2.49 | 1.98 | 2.86 |
| $D_{o,\max}$ [s] | 133.47 | 157.8 | 172.5 | 171.2 | 152.4 | 157 | 148.8 | 169.1 | 152.1 | 123.7 | 172.5 | 116.7 |

**Table 9.5:** Average throughput $\hat{S}$ and maximum outage duration $D_{o,\max}$ on the four itineraries from Fig. 9.10, for different departure times in February 2017. For the three routes with a similar duration, the colored cells represent the best route for the metric of interest.

particularly relevant in view of the envisioned transition to an autonomous driving future, in which active driving might not be required and working or getting entertained in the car will become a common trend. In order to address this need, cellular network operators can exploit the architecture described in Sec. 9.3 and the prediction of the number of active users in the cells to offer anticipatory services to the end users and inform them on which is the best route for their journey.

As mentioned in Sec. 9.2, the throughput cannot be directly and reliably collected from the measurement framework we used, which provides instead network KPIs and exact counters for mobility-related quantities such as the number of active users. Therefore, we estimate the user throughput as inversely proportional to the number of active users. In particular, we express the user throughput at base station $i$, time $t$ and user's position $p$ as

$$S_i(t,p) = \frac{\hat{U}(t)}{\frac{N_u^i(t)+1}{N_s^i}} B^i \rho^i(p),$$  (9.3)

where $N_u^i(t)$ is the number of users, $N_s^i$ the number of sectors, $B^i$ is the bandwidth and $\rho^i(t,p)$ is the spectral efficiency. $\hat{U}(t) \in [0,1]$ is the maximum PRB utilization, defined as the median over the considered dataset of the maximum daily PRB utilization of all the base stations, and in this case it is equal to 0.91. Both $N_s^i$ and $B^i$ are known, given the network configuration. The spectral efficiency $\rho^i(p)$, instead, depends on the mapping of the estimated SINR of the user in position $p$ to the CQI, using the map in [383], and then of the CQI to the spectral efficiency, according to 3GPP mapping from [384, Table 7.2.3-1]. The SINR is computed as

$$\Gamma^i(p) = \frac{P_{tx}^i L^i(p)}{I(p) + B^i N_0},$$  (9.4)

where $P_{tx}^i$ is the transmitted power of base station $i$, $L^i(p)$ the pathloss, computed as a function of distance and frequency using the equations in [385], $I(p)$ the interference, and $N_0 = -174$ dBm/Hz the thermal noise. For the interference, we consider the set of all the base stations except $i$, i.e., $\mathcal{B} \setminus \{i\}$, and, for each of them, compute the received power in position $p$.[6] Then, if the power is above a certain threshold (e.g., 10 dB below the thermal noise), it is added to the total count for $I(p)$.

Table 9.5 reports the value of different throughput-related metrics for the three itineraries with similar travel time, and a longer one, shown in Fig. 9.10, and identifies the best route according to each metric. The average throughput is measured as the average of the user throughput over

---

[6]This is a worst case scenario, since the base station may not be always transmitting, or may be using beamforming to steer the power towards its users and not omnidirectionally

the drive time for each itinerary, i.e.,

$$\hat{S} = \frac{1}{D} \sum_{d=1}^{D} S_{i(p_d)}(t_d, p_d), \tag{9.5}$$

where $D$ it the number of points sampled along the itinerary (e.g., provided by Google Maps), each at time $t_d$ and with position $p_d$, and $i(p_d)$ is the index of the closest base station to the position $p_d$. The maximum outage duration is given by the maximum time interval on the journey in which the user is offered a zero throughput, for example, because it is too far from the base stations, or the interference from the neighbors is too strong, and thus CQI 0 is selected. A high average throughput is desirable for web browsing, video and audio streaming, while a short maximum outage duration is preferable, for example, to attend conference calls.

As it can be seen from Table 9.5, the fastest route (i.e., route 1, in blue), is not always the one offering the best service in the three departure times considered. For the first three routes, which have a similar travel time, the best route changes at different departure times: for the throughput, on Feb. 23rd, 19:00, route 2 (red) is better than the others, while in the next day at the same time the best itinerary is route 3 (green). When considering also the longest route, which still leads from the origin to the destination, but takes 50% more time than the shortest, it can be seen that it always offers the highest average throughput, but, in some cases, is one of the worst in terms of maximum outage duration.

This example shows that, according to the users' needs, it is possible to identify and select different routes that have different performance in terms of throughput and outage. Moreover, the routes are ranked differently according to various departure times. Therefore, simply apply-
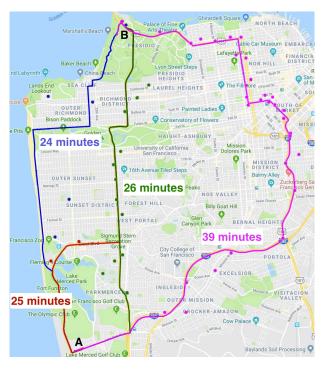


**Figure 9.10:** Map of the routes. The dots represent the visited base stations. Notice that, for route 2 (the red one), several base stations are shared with either the blue or the green routes.

ing the analytics given by the average statistics from the previous days may not yield reliable results in terms of routes ranking. This makes the case for adopting the medium-term prediction techniques described in this section to forecast the expected value of the metrics in the time interval in which the user will travel, based on the actual network conditions for the same day.

## 9.6  Data-driven vMME Allocation

This section, based on [392], discusses the data-driven optimization of the number of virtual MMEs instances in a 5G core network based on NFV. We consider a dataset based not on network data, but on Smart City sensors (described in Sec. 9.6.1), and show how it can be integrated to perform 5G network optimizations in Sec. 9.6.2.

### 9.6.1  Smart City dataset

The TfL UTC network is composed of more than 10000 road sensors, placed at all critical crossings around the city. The Split Cycle Offset Optimization Technique (SCOOT) optimizer uses the traffic flow data from the sensors to adapt the traffic light times to the traffic situation in real time. TfL released the raw sensor data of the first three months of 2015 for the North and Central regions of London, and we use those data in our optimization.

The sensors are actually very basic presence-detectors: every $T_s = 250$ ms, each sensor returns a 1 if it detects a vehicle in close proximity, and a 0 otherwise. The resulting binary signal is packetized and sent to a central collector through different types of technologies.

In this study, we extract from the TfL dataset the number of handovers between Macro eNBs over the whole city to perform data-driven cellular network optimizations.

These values are not directly provided by TfL. However, they can be roughly estimated using the binary signals generated by the detectors. In particular, we assume that the Macro eNBs in a ultra-dense scenario are placed using a standard regular hexagonal tiling, with sides of 100 m, and associate the detection of a car by a sensor in a cell with a handover. Then, given a time interval $T_{per}$ equal to 1 hour, we estimate the number of handovers $H_m$ as the total number of detections from the different sensors in cell $m$ during $T_{per}$. Since the timescale is long and each vehicle is likely detected only once when crossing the area (because of the relatively low density of sensors), the number of vehicles counted in the area in the period $T_{per}$ is roughly equal to the number of cell handovers performed by the vehicles crossing that area in the considered time interval. This assumption is not necessarily realistic for a single cell, but is a valid approximation on the city-wide scale and for timescales of minutes or hours. Moreover, we assume that on average each vehicle carries an LTE device. This is a working assumption based on the available data, and the integration of additional data such as bus position and usage can be easily accommodated by the framework.

After computing $H_m$ for all eNBs, the cells are partitioned into $N$ areas, with $N \in \{1, 2, 3, 4\}$, each controlled by a different vMME; given the estimated number of handovers at peak hours, 4 vMMEs should be enough to maintain network stability. The results in Sec. 9.6.2 confirm this hypothesis. These groups are obtained using a clustering algorithm that divides the cells among $N$ vMMEs so that each vMME handles approximately the same number of handovers. An example of this is shown in Fig. 9.12, which reports the partitions for $N \in \{2, 3, 4\}$.

We define $I_i$ as the total number of handovers for vMME $i$, and $S_{i,j}$ as the number of handovers from vMME $i$ to vMME $j$. $I_i$ is given by

$$I_i = \sum_{m \in A_i} H_m \tag{9.6}$$

**Figure 9.11:** Map of traffic in London from 12 PM to 1 PM of January 23, 2015. Free intersections are shown in green, heavily congested ones in red.



**(a)** $N = 2$            **(b)** $N = 3$            **(c)** $N = 4$

**Figure 9.12:** Partition for a different number $N$ of vMMEs. The colors indicate the areas controlled by each vMMEs.

where $A_i$ is the set of cells controlled by vMME $i$. $S_{i,j}$ can be approximated with this formula:

$$S_{i,j} = \sum_{m \in A_i} \sum_{n \in A_j} \frac{H_m}{6} e_{m,n} \tag{9.7}$$

where the variable $e_{m,n} \in \{0, 1\}$ indicates the number of sides that cells $m$ and $n$ have in common.

### 9.6.2 Dynamic Allocation of vMMEs

As already mentionded in Sec. 9.1.2, NFV allows to dynamically allocate the resources needed by a cellular network. In traditional mobile networks a single dedicated MME is typically used to manage millions of end users, such as those in the London metropolitan area [353]. With the NFV approach, instead, it is possible to change the number of vMME instances on the fly, adapting to the number of handovers that are expected to happen in a certain interval.

In this application, we use data processed as in Sec. 9.6.1 to determine the number of handovers that happen in the London area during a typical day. We distinguish between the two kinds of handovers that may happen in LTE networks [119], i.e., intra MME (X2–based) and inter MME (S1–based) handovers, since they require different procedures and different interactions with the MMEs. The X2–based handover happens when the UE remains in an area managed by the same MME and changes the eNB to which it is attached. The S1–based procedure, instead, is used when the UE performs a handover between two eNBs managed by different MMEs. The two procedures are described in detail in [119]. In this paper, we consider the duration of a handover procedure as the interval from the instant in which the source eNB (SeNB) triggers the handover to the instant in which SeNB receives the `RELEASE_RESOURCES` command. During this period the UE first experiences a degraded channel, and then receives packets with an increased latency, thus the Quality of Service perceived by the final user decreases. The goal of this application is to minimize the duration of these intervals, while using as few vMME instances as possible.

In particular, we model the duration of an X2–based handover handled by vMME $i$ as a function of the number of vMMEs $N$ and of the total number of handovers $I_i$ that involve that vMME during an interval $T_{per}$:

$$t_{HO}^{X2}(N, I_i) = 3t_{Se-Te} + 2t_{Te-SM}(N) + t_{HR} + \tau(I_i) \tag{9.8}$$

while the time required to complete an S1–based handover that involves vMMEs $i$ and $j$ also depends on the number of handovers $I_j$ that are served by the target vMME $j$:

$$\begin{aligned} t_{HO}^{S1}(N, I_i, I_j) =& \tau_1(I_i) + 3\tau_2(I_j) + \\ & 4t_{Se-SM}(N) + 4t_{Te-TM}(N) + \\ & 2t_{SM-TM}(N) + \max\{t_{TM-SM}(N) + \\ & t_{SM-Se}(N) + t_{HR}, t_{TM-Se}(N)\} + \\ & \max\{t_{TM-SM}(N) + \tau_1(I_i), t_{TM-Se}\} \end{aligned} \tag{9.9}$$

In Eqs. (9.8) and (9.9), $t_{A-B}(N)$ with $A, B \in \{Te, Se, SM, TM\}$[7] is the latency between element $A$ and element $B$ of the network. Unless both $A$ and $B$ represent eNBs, we have

$$t_{A-B}(N) = t_{tx} + \frac{d_N(A, B)}{v_f}, \tag{9.10}$$

where $t_{tx} = 5$ ms is a factor that models the time spent in middleboxes and $t_{PROP} = d_N(A, B)/v_f$ is the propagation delay, given by the ratio of the distance between the two devices and the speed of light inside optical fibers[8] (i.e., $v_f = 2 \cdot 10^8$ m/s). The dependence on the number of vMMEs $N$ is in the distance $d_N(A, B)$ between two network elements, that changes according to the allocation of eNBs to the vMMEs. Instead, $t_{Te-Se}$ is the latency between two adjacent eNBs and does not depend on the relative position between the eNBs and the MMEs, therefore, as in [182], it is modeled as a constant latency $t_{Te-Se} = 2.5$ ms. $t_{HR}$ is the duration of the interval from when the UE actually disconnects from the SeNB to when it connects to the TeNB. In [386], $t_{HR}$ is estimated to be in the order of 50 ms.

Finally, $\tau(I_i)$ is the time that a vMME takes to process the received command. In [353] the process of handover requests is modeled as a Markov process. We adopt the same approach and in particular we model the vMME as an M/D/1 queue, assuming a Poisson arrival process with *arrival rate* $\lambda = I_i/T_{per}$ and a deterministic service time $T_s$. Given these assumptions, it is possible to compute the value of $\tau$ as the *system time* of an M/D/1 queue:

$$\tau = \frac{1}{\mu} + \frac{\rho}{2 \cdot \mu \cdot (1 - \rho)}, \tag{9.11}$$

where $\mu = 1/T_s$ and $\rho = \lambda T$ are the *service rate* and the *loading factor* of the vMME. The study in [351] uses the value $T_s = 110$ $\mu s$ as *service time* of a vMME, requiring considerable computational resources. Since our work only considers vehicular UEs, and the adaptive nature of our system, overdimensioning each vMME would be a waste of resources: a number of slow vMMEs can provide the same performance as a single powerful vMME during rush hour, and the additional vMMEs can be turned off at less congested times, with a substantial reduction in server management costs and energy requirements. For this reason, we limit the processing

---

[7]$Te$ stands for Target eNB, $Se$ stands for Source eNB, $TM$ stands for Target MME and $SM$ stands for Source MME

[8]We assume that the backhaul network uses fiber-optic links.

**(a)** Average number of X2–based and S1–based handovers per vMMEs instance, for a different $N$ and different time slots, during January 23, 2015.

**(b)** Average service time $\tau$ for different $N$, during January 23, 2015.

power of our vMMEs dedicated to vehicular handovers to the value of $\mu = 1000$ handovers per second.

Since our goal is to find the optimal number of vMMEs $N$ that minimizes the total duration of the handovers, we consider the objective function

$$
J_{T_{per}}(N) = \sum_{i=1}^{N}(I_i - \sum_{\substack{j=1 \\ j \neq i}}^{N} S_{i,j})t_{HO}^{X2}(N, I_i)
$$
$$
+ \sum_{i=1}^{N}\sum_{\substack{j=1 \\ j \neq i}}^{N} S_{i,j}t_{HO}^{S1}(N, I_i, I_j) + C(N), \tag{9.12}
$$

where the sums consider all the handovers in a time slot $T_{per}$ of one hour, and $C(N)$ is a penalty function representing the operational cost of $N$ vMMEs. We consider it to be a linear function of the number of vMMEs $N$, i.e., $C(N) = kN$.

The optimization problem uses the vehicular traffic data processed as in Sec. 9.6.1 to compute the value of $I_i$, $S_{i,j}$ and $\lambda(I_i) = I_i/T_{per}$ for each vMME $i, j \in \{1, \cdots, N\}$ and computes

$$
N_{opt} = \min_{N} J_{T_{per}}(N) \tag{9.13}
$$

for each interval $T_{per}$ during a certain day.

In the following results we consider the data of January 23, 2015. Fig. 9.13a shows the average number of handovers inside a single vMME in different time slots. Notice that since we consider only the inter MME handovers for the London area MMEs, then $S_{ij}$ is zero for $N = 1$. The number of handovers in different time slots changes greatly, from $1.5 \cdot 10^6$ per hour during the night to more than $7 \cdot 10^6$ at midday. This justifies a dynamic allocation of resources; a single and dedicated MME that targets the worst case scenario at midday would be wasted during the night. Instead the adaptive approach allows the use of less powerful vMMEs, which are able to serve a smaller number of handover requests, and have lower operational expenses than dedicated hardware [252], but can be instantiated on the fly according to the control traffic intensity.

In Fig. 9.13b, the average service time of the vMME instances is shown for different values of $N$. It can be seen that during the night the values have a small difference, but one or two vMME instances are not enough to handle the load during the day. Fig. 9.14a, instead, shows the value

**(a)** Objective function $J(N)$, for $N \in \{1, 2, 3, 4\}$, during January 23, 2015.

**(b)** $N_{opt}$ for different costs $C(N)$, during January 23, 2015.

of the objective function $J(N)$ throughout the whole day, assuming a cost factor $k = 0$. In this case, one vMME instance is enough only from midnight to 5 AM, and more instances (up to 3) must be allocated during the day to meet the vehicular handover traffic load.

If we increase the value of $k$, as shown in Fig. 9.14b, the optimal number of vMMEs changes. At certain times using a lower number of vMMEs becomes more convenient, because of the operational cost which is now accounted for.

The adaptation of the number of vMMEs significantly improves the efficiency of the system: while a worst-case dimensioned system would need 3 vMMEs at all times, the average number of active vMME instances for the most aggressive adaptive system ($k = 0$) is 2.42, while a more conservative system ($k = 100000$) only uses an average of 2.17 vMMEs. This translates into a lower operating cost for the network provider because of a reduced energy consumption and of the need of using fewer virtual functions.

## 9.7 Conclusions

Machine learning, software-defined networks, network function virtualization and edge cloud will be key components of the next generation of cellular networks. In this chapter we investigated how these three elements can be jointly used in the data-driven design and optimization of 5G networks, providing insights and results based on (i) a dataset collected from hundreds of base stations of a major U.S. cellular network in two different cities for more than a month; and (ii) a dataset of vehicular traffic in London.

After reviewing the relevant state of the art, we investigated how it is possible to practically introduce machine learning and big-data-based policies in 5G cellular networks. We proposed an overlay architecture on top of 3GPP NR, in which multiple layers of controllers with different functionalities are used to collect the data from the RAN, process it and use it to infer intelligent policies that can be applied to the cellular network.

Next, we discussed a first application of the proposed architecture, i.e., a data-driven association algorithm between the gNBs and the RAN controllers themselves. We described a clustering solution that limits the interactions among different controllers to minimize the need for inter-controller synchronization and reduce the control plane latency, and evaluated the performance of the proposed approach using data from a real network.

Then, we outlined a second possible application enabled by our architecture, providing an extensive set of results related to the prediction accuracy of the number of users in base stations, using one month of data collected from the San Francisco base stations. In particular, we showed how the usage of the cluster-based architecture proposed in this chapter can reduce the prediction error. With respect to a solution in which each base station tries to perform the regression based solely on its own data, as realized by a completely distributed architecture (e.g., in LTE), the controller-based design makes it possible to aggregate data from multiple neighboring base stations, and to predict a vector with the number of users in the nodes associated to the controller. This captures the spatial correlation given by the mobility of users, and, especially when increasing the temporal horizon of the prediction, reduces the RMSE by up to 53%. Finally, we also described some prediction-enabled use cases, either to control the network itself, or to offer innovative predictive services to network users, for example by recommending different driving itineraries to improve the user experience in the network. We illustrated a real example in the San Francisco area, showing how the fastest route does not necessarily yield the best throughput, or the minimum outage, and that the best itinerary according to these metrics (which we derive from the number of users in each base station) may differ according to the departure time, so that a prediction-based approach is useful.

Finally, we used the vehicular traffic flow data to adaptively provision virtual resources and add or remove virtual MMEs, reducing operating costs without impacting the performance with respect to a worst-case dimensioned system. The performance benefits of this data-driven scheme can only increase as the integration of smart city (and network) data in the optimization of 5G deployment progresses, for example, by considering public transportation. Moreover, periodic or forecastable events (i.e., holidays and changes in the weather conditions) that impact mobility patterns can be added to the model in order to improve its accuracy.

We believe that this study addresses for the first time several issues related to the practical deployment of machine learning and data-driven techniques in 5G cellular networks, providing results and conclusions based on real-world datasets. As future work, we will test different prediction algorithms (e.g., neural networks) to understand if it is possible to improve even more the prediction accuracy, and will extend the regression to other relevant metrics in the network (e.g., the number of handovers, the utilization), to verify the limits of what can be actually predicted in a cellular network.

# 10
## Conclusions

This thesis has investigated the design and performance of mmWave and 5G cellular networks from an end-to-end and system level perspective. In particular, we focused on how to efficiently deploy mmWave networking architectures, on which are the end-to-end protocols that provide the best performance in a complex network with a mmWave RAN, and on some options to deploy data-driven intelligent techniques in 5G cellular networks.

We first discussed the main tool that has been used for the end-to-end performance evaluation, i.e., the ns-3 mmWave module that was developed as part of this thesis, with a 3GPP channel model for mmWave frequencies, and a wireless protocol stack that adapts and extend an LTE implementation towards an NR and mmWave support, also providing features such as carrier aggregation and dual connectivity. The development of this simulator allowed us to study for the first time with an open source tool the interactions that emerge among the different parts of the network, and to understand the interplay of the mmWave RAN and channel with the protocols of the higher layers of the networking stack.

In the second part of the thesis, we described the architectural solution that can be deployed to make mmWave networks more reliable, robust and with improved performance. Our first proposal was the usage of multi connectivity in the RAN to combine the benefits of different carrier frequencies, i.e., sub-6 GHz (e.g., with LTE), to provide a reliable coverage layer, and mmWaves, for high capacity in the hotspots where the signal is available. We proposed to tightly integrate the two systems at the PDCP layer, to implement an efficient mobility management framework which allows network operators to provide a seamless service to the end users. The second contribution was the analysis of the performance of beam management frameworks in 5G mmWave cellular networks, with simulation and analytical results to characterize the tradeoffs among the different parameters that can be tuned in a 3GPP NR deployment. In particular, we highlighted some design choices (e.g., high density of base stations with fewer directions to scan) which could strike a good compromise between the accuracy and the reactiveness of initial access and tracking processes. Moreover, we proposed mmBAC, a context-based beam management framework for highly mobile UAVs, and characterized its performance with a real prototype based on a commercial drone and 60 GHz radios. Finally, we investigated the benefits and challenges that IAB introduces in 3GPP NR deployments at mmWave frequencies, with the first end-to-end evaluation of this relaying technology. Our results show that IAB is a feasible solution to relay cell-edge traffic, even if the performance degrages in more congested scenarios. We have also introduced and studied the performance of path selection policies for IAB relays that do not require a central coordinator.

The third part studied how the interplay with the mmWave channel affects the performance of TCP, and how it is possible to improve it. Thanks to the ns-3 mmWave module, we were able to evaluate the performance of TCP in 3GPP scenarios, with mmWave links operating at 28 GHz. This comprehensive evaluation has highlighted the major pitfalls and limitations of TCP, which fails at tracking the capacity that the mmWave physical layer offers, especially after the LOS to NLOS transitions, and introduces bufferbloat, high latency spikes and low resource utilization. Consequently, we proposed a number of different optimizations at the transport layer and with in-network solutions that can boost the end-to-end performance. The first, milliProxy, is a transparent proxy that can be installed at the base stations, and coordinates with the lower layers of the wireless stack, to control the rate at which TCP at a remote server injects data in the network, without the need for any modification of the TCP/IP stack in the two endpoints of the communication. Similarly, X-TCP controls the congestion window of an uplink TCP flow following a similar cross layer approach. Finally, we studied two solutions that exploit multi-connectivity at the transport layer. MPTCP, i.e., the multipath extension of TCP, can improve the performance of users at the cell-edge, especially when one of the flows operates on a reliable sub-6 GHz link. Alternatively, we proposed a protocol that piggybacks on UDP and uses network coding to efficiently distribute packets over sub-6 GHz and mmWave connections. We showed that this solution manages to stream video with low latency and high quality.

The fourth and final part was instead related to the usage of data-driven and machine-learning techniques in 5G networks. We proposed to use controllers at the edge of the network to aggregate data and statistics on the behavior of the network itself, that can then be exploited to train machine learning algorithms and perform different optimizations. We highlighted how our proposal can be integrated in 3GPP NR networks, following the O-RAN approach, which advocates the deployment of controllers at the edge of 5G networks. The method we introduced was tested using a dataset with hundreds of base stations of a major U.S. operator, with two different use cases: (i) a clustering problem, where the controllers use the mobility data of the network to determine the association between the base stations and the controllers themselves; and (ii) a prediction problem, in which the controllers forecast the number of users in the base stations they supervise. We showed that, thanks to the proposed strategy, the results in both scenarios improve with respect to static or uncoordinated solutions. We also discussed an optimization problem to dynamically allocate the number of vMMEs in a cellular network that covers the area of London, using a dataset of vehicular traffic traces to determine the load of handovers at different times of the day.

## 10.1 Future Directions

While the research on mmWave communications has been particularly active in the last few years, there are several research directions that are relatively unexplored. In particular, most of the studies in this area are either based on analysis or simulations, and just a few are based on real-world experiments [387]. Analysis and simulation are valid tools, especially when the model is accurate and, as discussed in this thesis, captures the complexity and the interactions of the channel and the different elements of the protocol stack. Nonetheless, especially when it comes to mmWaves, an experimental validation of the design and proposals is important, given the impact that the channel behavior has on the overall performance. However, existing research testbeds for mmWave communications are either closed – because of commercial and intellectual property reasons or because they are based on commercial off-the-shelf devices – or extremely expensive. As part of our future research, we will study the feasibility of low-cost, open source and open hardware testbeds at mmWave and also terahertz frequencies, following the vision we

**Figure 10.1:** Architectural innovations introduced in 6G networks.

described in [425].

Besides the investigation of novel research methodologies in the mmWave domain, we will also focus on understanding how mobile networks can be designed to address the ever increasing connectivity requirements of future generations. In [400] we discuss possible scenarios and use cases for beyond 5G and 6G networks, and propose our vision for the technologies that can enable them. Fig. 10.1 illustrates the major system-level and full-stack research directions that we envision, which can be organized in three main areas:

- *Novel disruptive communication technologies:* 6G networks could very much benefit from even higher spectrum technologies, e.g., through terahertz and visible light communications. Moreover, 6G will also transform wireless networks by leveraging a set of technologies that have been enabled by recent physical layer and circuits research, but are not part of 3GPP NR 5G specifications, e.g., with full duplex communications in the radio access, and the support of simultaneous communications, sensing and localization.

- *Innovative network architectures:* the heterogeneity of the requirements of future network applications calls for new radical paradigms in the design of mobile network architectures. In order to combine the different communication technologies that will be available, a tight integration based on advanced multi connectivity techniques at different layer of the stack will be required. Moreover, future wireless networks will relax the boundaries of traditional cells, towards a cell-less architectural paradigm and a support of connectivity in the 3D space (e.g., with a mixed use of terrestrial and flying infrastructure). Finally, the advances in computing capabilities will bring the virtualization and disaggregation of the networking equipment to an extreme, with the 6G PHY and MAC layers fully virtualized, and simple and low-cost distributed units with just the antennas and minimal processing units.

- *Towards Intelligent End Devices:* we expect 6G to bring intelligence from centralized or edge computing facilities (as those discussed in this thesis) to the end terminals, thereby giving concrete form to distributed learning models that have been studied from a theoretical point of view in a 5G context. Unsupervised learning and inter-user inter-operator knowledge sharing will also promote data-driven real-time network decisions.

# References

[1] Cisco, "Cisco Visual Networking Index: Global Mobile Data Traffic Forecast Update, 2016–2021," *White Paper*, March 2017.

[2] ITU-R, "IMT Vision - Framework and overall objectives of the future development of IMT for 2020 and beyond," Recommendation ITU-R M.2083, September 2015.

[3] M. Khoshnevisan, V. Joseph, P. Gupta, F. Meshkati, R. Prakash, and P. Tinnakornsrisuphap, "5G Industrial Networks With CoMP for URLLC and Time Sensitive Network Architecture," *IEEE Journal on Selected Areas in Communications*, vol. 37, no. 4, pp. 947–959, April 2019.

[4] K. Antonakoglou, X. Xu, E. Steinbach, T. Mahmoodi, and M. Dohler, "Toward Haptic Communications Over the 5G Tactile Internet," *IEEE Communications Surveys & Tutorials*, vol. 20, no. 4, pp. 3034–3059, Fourth quarter 2018.

[5] A. Osseiran, F. Boccardi, V. Braun, K. Kusume, P. Marsch, M. Maternia, O. Queseth, M. Schellmann, H. Schotten, H. Taoka, H. Tullberg, M. A. Uusitalo, B. Timus, and M. Fallgren, "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.

[6] F. Boccardi, R. W. Heath, A. Lozano, T. L. Marzetta, and P. Popovski, "Five disruptive technology directions for 5G," *IEEE Communications Magazine*, vol. 52, no. 2, pp. 74–80, February 2014.

[7] 3GPP, "NR and NG-RAN Overall Description," TS 38.300, V15.0.0, 2018.

[8] J. Andrews, S. Buzzi, W. Choi, S. Hanly, A. Lozano, A. Soong, and J. Zhang, "What will 5G be?" *IEEE J. Sel. Areas Commun.*, vol. 32, no. 6, pp. 1065–1082, June 2014.

[9] M. Shafi, A. F. Molisch, P. J. Smith, P. Z. T. Haustein, P. D. Silva, F. Tufvesson, A. Benjebbour, and G. Wunder, "5G: A Tutorial Overview of Standards, Trials, Challenges, Deployment, and Practice," *IEEE Journal on Selected Areas in Communications*, vol. 35, no. 6, pp. 1201–1221, June 2017.

[10] European Commission - Radio Spectrum Policy Group, "Strategic Spectrum Roadmap Towards 5G For Europe - RSPG Opinion on 5G implementation challenges (RSPG 3rd opinion on 5G)," RSPG19 - 007 FINAL, Jan. 2019. [Online]. Available: https://rspg-spectrum.eu/wp-content/uploads/2013/05/RSPG19-007final-3rd_opinion_on_5G.pdf

[11] T. L. Marzetta, "Noncooperative cellular wireless with unlimited numbers of base station antennas," *IEEE Transactions on Wireless Communications*, vol. 9, no. 11, pp. 3590–3600, November 2010.

[12] A. Checko, H. L. Christiansen, Y. Yan, L. Scolari, G. Kardaras, M. S. Berger, and L. Dittmann, "Cloud RAN for Mobile Networks – Technology Overview," *IEEE Communications Surveys & Tutorials*, vol. 17, no. 1, pp. 405–426, First quarter 2015.

[13] N. Bui, M. Cesana, S. A. Hosseini, Q. Liao, I. Malanchini, and J. Widmer, "A survey of anticipatory mobile networking: Context-based classification, prediction methodologies, and optimization techniques," *IEEE Communications Surveys & Tutorials*, vol. 19, no. 3, pp. 1790–1821, Third quarter 2017.

[14] E. Dahlman, S. Parkvall, and J. Skold, *5G NR: The next generation wireless access technology*. Academic Press, 2018.

[15] 3GPP, "Study on Scenarios and Requirements for Next Generation Access Technologies," TR 38.913, V14.1.0, 2017.

[16] ——, "WF on Work plan of Self Evaluation SI," RP-172101, 3GPP TSG RAN77, Sapporo, Japan, September 2017.

[17] A. A. Zaidi, R. Baldemair, H. Tullberg, H. Bjorkegren, L. Sundstrom, J. Medbo, C. Kilinc, and I. D. Silva, "Waveform and Numerology to Support 5G Services and Requirements," *IEEE Communications Magazine*, vol. 54, no. 11, pp. 90–98, November 2016.

[18] 3GPP, "NR - Physical channels and modulation," TR 38.211, V15.0.0, 2017.

[19] N. Patriciello, S. Lagen, L. Giupponi, and B. Bojovic, "5G New Radio Numerologies and their Impact on the End-To-End Latency," in *IEEE 23rd International Workshop on Computer Aided Modeling and Design of Communication Links and Networks (CAMAD)*, Sept 2018, pp. 1–6.

[20] R. Ford, M. Zhang, M. Mezzavilla, S. Dutta, S. Rangan, and M. Zorzi, "Achieving Ultra-Low Latency in 5G Millimeter Wave Cellular Networks," *IEEE Communications Magazine*, vol. 55, no. 3, pp. 196–203, March 2017.

[21] K. Takeda, L. H. Wang, and S. Nagata, "Latency Reduction toward 5G," *IEEE Wireless Communications*, vol. 24, no. 3, pp. 2–4, 2017.

[22] 3GPP, "Study on New Radio (NR) access technology," TR 38.912, V15.0.0, 2018.

[23] S. Parkvall, E. Dahlman, A. Furuskar, and M. Frenne, "NR: The New 5G Radio Access Technology," *IEEE Communications Standards Magazine*, vol. 1, no. 4, pp. 24–30, Dec 2017.

[24] E. Björnson, E. G. Larsson, and T. L. Marzetta, "Massive MIMO: ten myths and one critical question," *IEEE Communications Magazine*, vol. 54, no. 2, pp. 114–123, February 2016.

[25] 3GPP, "NG-RAN; Architecture description - Release 15," TS 38.401, V15.3.0, 2018.

[26] ——, "System Architecture for the 5G System; Stage 2 - Release 15," TS 23.501, V15.3.0, 2018.

[27] N. Makris, C. Zarafetas, P. Basaras, T. Korakis, N. Nikaein, and L. Tassiulas, "Cloud-based Convergence of Heterogeneous RANs in 5G Disaggregated Architectures," in *IEEE International Conference on Communications (ICC)*. IEEE, 2018.

[28] A. Maeder, A. Ali, A. Bedekar, A. F. Cattoni, D. Chandramouli, S. Chandrashekar, L. Du, M. Hesse, C. Sartori, and S. Turtinen, "A scalable and flexible radio access network architecture for fifth generation mobile networks," *IEEE Communications Magazine*, vol. 54, no. 11, pp. 16–23, November 2016.

[29] 3GPP, "Evolved Universal Terrestrial Radio Access (E-UTRA) and NR; Multi-connectivity," TS 37.340 (Rel. 15), 2018.

[30] X. Foukas, G. Patounas, A. Elmokashfi, and M. K. Marina, "Network Slicing in 5G: Survey and Challenges," *IEEE Communications Magazine*, vol. 55, no. 5, pp. 94–100, May 2017.

[31] 3GPP, "Study on requirements for NR beyond 52.6 GHz," TR 38.807, V0.2.0, 2019.

[32] Z. Pi and F. Khan, "An introduction to millimeter-wave mobile broadband systems," *IEEE Communications Magazine*, vol. 49, no. 6, June 2011.

[33] S. Akoum, O. El Ayach, and R. W. Heath, "Coverage and capacity in mmwave cellular systems," in *Conference Record of the Forty Sixth Asilomar Conference on Signals, Systems and Computers (ASILOMAR)*. IEEE, 2012, pp. 688–692.

[34] IEEE, "IEEE Standard for Information technology—Telecommunications and information exchange between systems Local and metropolitan area networks—Specific requirements - Part 11: Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications," *IEEE Std 802.11-2016 (Revision of IEEE Std 802.11-2012)*, pp. 1–3534, Dec 2016.

[35] ——, "IEEE Draft Standard for Information Technology – Telecommunications and Information Exchange Between Systems Local and Metropolitan Area Networks – Specific Requirements Part 11: Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications Amendment Enhancements for High Efficiency WLAN," *IEEE Draft Std P802.11ax (Amendment to IEEE Std 802.11-2012, as amended by IEEE Std 802.11ae-2012 and IEEE Std 802.11aa-2012)*, 2019.

[36] Y. Niu, Y. Li, D. Jin, L. Su, and A. V. Vasilakos, "A survey of millimeter wave communications (mmWave) for 5G: opportunities and challenges," *Wireless Networks*, vol. 21, no. 8, pp. 2657–2676, Nov. 2015.

[37] M. R. Akdeniz, Y. Liu, M. K. Samimi, S. Sun, S. Rangan, T. S. Rappaport, and E. Erkip, "Millimeter wave channel modeling and cellular capacity evaluation," *IEEE Journal on Selected Areas in Communications*, vol. 32, no. 6, pp. 1164–1179, June 2014.

[38] T. S. Rappaport, *Wireless Communications: Principles and Practice*, 2nd ed. Upper Saddle River, NJ: Prentice Hall, 2002.

[39] T. S. Rappaport, S. Sun, R. Mayzus, H. Zhao, Y. Azar, K. Wang, G. N. Wong, J. K. Schulz, M. Samimi, and F. Gutierrez, "Millimeter wave mobile communications for 5G cellular: It will work!" *IEEE Access*, vol. 1, pp. 335–349, May 2013.

[40] J. S. Lu, D. Steinbach, P. Cabrol, and P. Pietraski, "Modeling human blockers in millimeter wave radio links," *ZTE Communications*, vol. 10, no. 4, pp. 23–28, 2012.

[41] V. Raghavan, L. Akhoondzadeh-Asl, V. Podshivalov, J. Hulten, M. A. Tassoudji, O. H. Koymen, A. Sampath, and J. Li, "Statistical blockage modeling and robustness of beamforming in millimeter-wave systems," *IEEE Transactions on Microwave Theory and Techniques*, vol. 67, no. 7, pp. 3010–3024, July 2019.

[42] S. Rangan, T. S. Rappaport, and E. Erkip, "Millimeter-Wave Cellular Wireless Networks: Potentials and Challenges," *Proceedings of the IEEE*, vol. 102, no. 3, pp. 366–385, March 2014.

[43] S. Sun, T. S. Rappaport, R. W. Heath, A. Nix, and S. Rangan, "MIMO for millimeter-wave wireless communications: beamforming, spatial multiplexing, or both?" *IEEE Communications Magazine*, vol. 52, no. 12, pp. 110–121, December 2014.

[44] S. Dutta, C. N. Barati, A. Dhananjay, and S. Rangan, "5G Millimeter Wave Cellular System Capacity with Fully Digital Beamforming," in *51st Asilomar Conference on Signals, Systems and Computers*, Nov 2017.

[45] A. Alkhateeb, O. El Ayach, G. Leus, and R. W. Heath, "Channel estimation and hybrid precoding for millimeter wave cellular systems," *IEEE Journal of Selected Topics in Signal Processing*, vol. 8, no. 5, pp. 831–846, Oct 2014.

[46] H. Shokri-Ghadikolaei, C. Fischione, G. Fodor, P. Popovski, and M. Zorzi, "Millimeter wave cellular networks: A MAC layer perspective," *IEEE Trans. Comm.*, vol. 63, no. 10, pp. 3437–3458, Oct. 2015.

[47] M. Rebato, M. Mezzavilla, S. Rangan, F. Boccardi, and M. Zorzi, "Understanding Noise and Interference Regimes in 5G Millimeter-Wave Cellular Networks," in *22th European Wireless Conference*, 2016.

[48] W. Feng, Y. Wang, D. Lin, N. Ge, J. Lu, and S. Li, "When mmwave communications meet network densification: A scalable interference coordination perspective," *IEEE Journal on Selected Areas in Communications*, vol. 35, no. 7, pp. 1459–1471, July 2017.

[49] N. Lee and R. W. Heath Jr, "Advanced interference management technique: potentials and limitations," *IEEE Wireless Communications*, vol. 23, no. 3, pp. 30–38, June 2016.

[50] C. Pielli, T. Ropitault, and M. Zorzi, "The Potential of mmWaves in Smart Industry: Manufacturing at 60 GHz," in *Ad-hoc, Mobile, and Wireless Networks*. Springer International Publishing, 2018, pp. 64–76.

[51] M. N. Islam, S. Subramanian, and A. Sampath, "Integrated access backhaul in millimeter wave networks," in *IEEE Wireless Communications and Networking Conference (WCNC)*, March 2017, pp. 1–6.

[52] S. H. Ali Shah, S. Aditya, S. Dutta, C. Slezak, and S. Rangan, "Power Efficient Discontinuous Reception in THz and mmWave Wireless Systems," in *IEEE 20th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, July 2019, pp. 1–5.

[53] M. Allman, V. Paxson, and E. Blanton, "TCP congestion control," IETF, RFC 5681, Sep. 2009. [Online]. Available: https://rfc-editor.org/rfc/rfc5681.txt

[54] T. Stockhammer, "Dynamic Adaptive Streaming over HTTP: Standards and Design Principles," in *Proceedings of the Second Annual ACM Conference on Multimedia Systems (MMSys)*, 2011.

[55] T. Bai and R. W. Heath, "Coverage and rate analysis for millimeter-wave cellular networks," *IEEE Transactions on Wireless Communications*, vol. 14, no. 2, pp. 1100–1114, February 2015.

[56] J. G. Andrews, T. Bai, M. N. Kulkarni, A. Alkhateeb, A. K. Gupta, and R. W. Heath, "Modeling and analyzing millimeter wave cellular systems," *IEEE Transactions on Communications*, vol. 65, no. 1, pp. 403–430, Jan 2017.

[57] M. Zhang, M. Mezzavilla, R. Ford, S. Rangan, S. Panwar, E. Mellios, D. Kong, A. Nix, and M. Zorzi, "Transport layer performance in 5G mmWave cellular," in *IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS)*, April 2016, pp. 730–735.

[58] T. S. Rappaport, R. W. Heath Jr., R. C. Daniels, and J. N. Murdock, *Millimeter Wave Wireless Communications*. Pearson Education, 2014.

[59] K. Allen *et al.*, *Building penetration loss measurements at 900 MHz, 11.4 GHz, and 28.8 MHz*, ser. NTIA report – 94-306. Boulder, CO: U.S. Dept. of Commerce, National Telecommunications and Information Administration, 1994.

[60] S. Singh, F. Ziliotto, U. Madhow, E. M. Belding, and M. J. W. Rodwell, "Millimeter Wave WPAN: Cross-Layer Modeling and Multi-Hop Architecture," in *26th IEEE International Conference on Computer Communications (INFOCOM)*, May 2007, pp. 2336–2340.

[61] A. Ghosh, T. A. Thomas, M. C. Cudak, R. Ratasuk, P. Moorut, F. W. Vook, T. S. Rappaport, G. MacCartney, S. Sun, and S. Nie, "Millimeter wave enhanced local area systems: A high data rate approach for future wireless networks," *IEEE J. Sel. Areas Commun.*, vol. 32, no. 6, pp. 1152–1163, June 2014.

[62] M. Giordani, M. Mezzavilla, S. Rangan, and M. Zorzi, "Multi-Connectivity in 5G mmwave cellular networks," in *Ad Hoc Mediterranean Networking Workshop (Med-Hoc-Net)*. IEEE, 2016.

[63] F. B. Tesema, A. Awada, I. Viering, M. Simsek, and G. P. Fettweis, "Mobility modeling and performance evaluation of multi-connectivity in 5G intra-frequency networks," in *IEEE Globecom Workshops*. IEEE, 2015.

[64] M. Zhang, M. Mezzavilla, J. Zhu, S. Rangan, and S. Panwar, "TCP dynamics over mmwave links," in *IEEE 18th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, July 2017, pp. 1–6.

[65] R. Ford, A. Sridharan, R. Margolies, R. Jana, and S. Rangan, "Provisioning Low Latency, Resilient Mobile Edge Clouds for 5G," in *IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS)*. IEEE, 2017.

[66] M. Mezzavilla, S. Dutta, M. Zhang, M. R. Akdeniz, and S. Rangan, "5G mmWave Module for the ns-3 Network Simulator," in *Proceedings of the 18th ACM International Conference on Modeling, Analysis and Simulation of Wireless and Mobile Systems*, 2015, pp. 283–290. [Online]. Available: http://doi.acm.org/10.1145/2811587.2811619

[67] R. Ford, M. Zhang, S. Dutta, M. Mezzavilla, S. Rangan, and M. Zorzi, "A framework for end-to-end evaluation of 5G mmwave cellular networks in ns-3," in *Proceedings of the Workshop on ns-3*. ACM, 2016, pp. 85–92.

[68] T. R. Henderson, M. Lacage, G. F. Riley, C. Dowell, and J. Kopena, "Network simulations with the ns-3 simulator," *SIGCOMM demonstration*, vol. 14, no. 14, p. 527, 2008.

[69] N. Baldo, M. Miozzo, M. Requena-Esteso, and J. Nin-Guerrero, "An open source product-oriented lte network simulator based on ns-3," in *Proceedings of the 14th ACM International Conference on Modeling, Analysis and Simulation of Wireless and Mobile Systems*, 2011, pp. 293–298. [Online]. Available: http://doi.acm.org/10.1145/2068897.2068948

[70] "LTE-EPC Network Simulator," Available at `http://iptechwiki.cttc.es/ LTE-EPC_Network_Simulator_(LENA)`, Feb. 2012.

[71] 3GPP, "Study on channel model for frequency spectrum above 6 GHz," TR 38.900, V14.2.0, 2017.

[72] H. Assasa and J. Widmer, "Implementation and Evaluation of a WLAN IEEE 802.11ad Model in ns-3," in *Proceedings of the Workshop on ns-3*. ACM, 2016, pp. 57–64.

[73] ——, "Extending the IEEE 802.11ad Model: Scheduled Access, Spatial Reuse, Clustering, and Relaying," in *Proceedings of the Workshop on ns-3*. ACM, 2017, pp. 39–46.

[74] H. Assasa, J. Widmer, T. Ropitault, and N. Golmie, "Enhancing the ns-3 IEEE 802.11ad Model Fidelity: Beam Codebooks, Multi-antenna Beamforming Training, and Quasi-deterministic mmWave Channel," in *Proceedings of the 2019 Workshop on ns-3*, ser. WNS3 2019. New York, NY, USA: ACM, 2019, pp. 33–40. [Online]. Available: http://doi.acm.org/10.1145/3321349.3321354

[75] H. Assasa, J. Widmer, T. Ropitault, A. Bodi, and N. Golmie, "High Fidelity Simulation of IEEE 802.11ad in ns-3 Using a Quasi-deterministic Channel Model," in *Proceedings of the 2019 Workshop on Next-Generation Wireless with ns-3*, ser. WNGW 2019. New York, NY, USA: ACM, 2019, pp. 22–25. [Online]. Available: http://doi.acm.org/10.1145/3337941.3337946

[76] H. Assasa, J. Widmer, J. Wang, T. Ropitault, and N. Golmie, "An Implementation Proposal for IEEE 802.11ay SU/MU-MIMO Communication in ns-3," in *Proceedings of the 2019 Workshop on Next-Generation Wireless with ns-3*, ser. WNGW 2019. New York, NY, USA: ACM, 2019, pp. 26–29. [Online]. Available: http://doi.acm.org/10.1145/3337941.3337947

[77] T. Kim, J. Park, J.-Y. Seol, S. Jeong, J. Cho, and W. Roh, "Tens of Gbps support with mmWave beamforming systems for next generation communications," in *IEEE Global Communications Conference (GLOBECOM)*, Dec. 2013, pp. 3685–3690.

[78] K. Zheng, L. Zhao, J. Mei, M. Dohler, W. Xiang, and Y. Peng, "10 Gb/s HetSNets with Millimeter-Wave Communications: Access and Networking - Challenges and Protocols," *IEEE Communications Magazine*, vol. 53, no. 1, pp. 222–231, Jan. 2015.

[79] C. Dehos, J. L. González, A. D. Domenico, D. Kténas, and L. Dussopt, "Millimeter-wave access and backhauling: the solution to the exponential data traffic increase in 5G mobile communications systems?" *IEEE Communications Magazine*, vol. 52, no. 9, pp. 88–95, September 2014.

[80] R. Taori and A. Sridharan, "Point-to-multipoint in-band mmwave backhaul for 5G networks," *IEEE Communications Magazine*, vol. 53, no. 1, pp. 195–201, Jan. 2015.

[81] S. Choi, J. Song, J. Kim, S. Lim, S. Choi, T. T. Kwon, and S. Bahk, "5G K-SimNet: End-to-End Performance Evaluation of 5G Cellular Systems," in *16th IEEE Annual Consumer Communications Networking Conference (CCNC)*, Jan 2019, pp. 1–6.

[82] N. Patriciello, S. Lagen, B. Bojovic, and L. Giupponi, "An E2E simulator for 5G NR networks," *Simulation Modelling Practice and Theory*, vol. 96, p. 101933, 2019. [Online]. Available: http://www.sciencedirect.com/science/article/pii/S1569190X19300589

[83] N. Patriciello, S. Lagen, L. Giupponi, and B. Bojovic, "An Improved MAC Layer for the 5G NR ns-3 Module," in *Proceedings of the 2019 Workshop on ns-3*, ser. WNS3 2019. New York, NY, USA: ACM, 2019, pp. 41–48. [Online]. Available: http://doi.acm.org/10.1145/3321349.3321350

[84] T. R. Henderson, S. Roy, S. Floyd, and G. F. Riley, "ns-3 project goals," in *Proceeding of the 2006 workshop on ns-2: the IP network simulator*. ACM, 2006, p. 13.

[85] M. Casoni, C. A. Grazia, M. Klapez, and N. Patriciello, "Implementation and validation of TCP options and congestion control algorithms for ns-3," in *Proceedings of the 2015 Workshop on ns-3*. ACM, 2015, pp. 112–119.

[86] G. Pei and T. Henderson, "Validation of ns-3 802.11 b PHY model," Technical Report, 2009. [Online]. Available: https://www.nsnam.org/~pei/80211b.pdf

[87] J. Farooq and T. Turletti, "An IEEE 802.16 WiMAX Module for the ns-3 Simulator," in *Proceedings of the 2nd International Conference on Simulation Tools and Techniques*, 2009, pp. 8:1–8:11. [Online]. Available: http://dx.doi.org/10.4108/ICST.SIMUTOOLS2009.5644

[88] H. Narra, Y. Cheng, E. K. Çetinkaya, J. P. Rohrer, and J. P. G. Sterbenz, "Destination-sequenced Distance Vector (DSDV) Routing Protocol Implementation in ns-3," in *Proceedings of the 4th International ICST Conference on Simulation Tools and Techniques*, 2011, pp. 439–446. [Online]. Available: http://dl.acm.org/citation.cfm?id=2151054.2151132

[89] H. Tazaki, F. Uarbani, E. Mancini, M. Lacage, D. Camara, T. Turletti, and W. Dabbous, "Direct Code Execution: Revisiting Library OS Architecture for Reproducible Network Experiments," in *Proceedings of the Ninth ACM Conference on Emerging Networking Experiments and Technologies*, ser. CoNEXT '13. New York, NY, USA: ACM, 2013, pp. 217–228. [Online]. Available: http://doi.acm.org/10.1145/2535372.2535374

[90] Centre Tecnologic de Telecomunicacions de Catalunya (CTTC), "The LENA ns-3 LTE Module Documentation," Available at `http://iptechwiki.cttc.es/ LTE-EPC_Network_Simulator_(LENA)`, Jan 2014.

[91] 3GPP, "Spatial channel model for multiple input multiple output (MIMO) simulations," TR 25.996, V6.0.0, 2003.

[92] "Winprop software," Available at *http://www.awe-communications.com/Products/*.

[93] R. N. Almesaeed, A. S. Ameen, E. Mellios, A. Doufexi, and A. Nix, "3D Channel Models: Principles, Characteristics, and System Implications," *IEEE Communications Magazine*, vol. 55, no. 4, pp. 152–159, April 2017.

[94] S. Jaeckel, L. Raschkowski, K. Borner, and L. Thiele, "QuaDRiGa: A 3-D Multi-Cell Channel Model With Time Evolution for Enabling Virtual Field Trials," *IEEE Transactions on Antennas and Propagation*, vol. 62, no. 6, pp. 3242–3256, Mar. 2014.

[95] S. Sun, G. R. MacCartney, Jr., and T. S. Rappaport, "A novel millimeter-wave channel simulator and applications for 5G wireless communications," in *IEEE International Conference on Communications (ICC)*, May 2017, pp. 1–7.

[96] M. K. Samimi and T. S. Rappaport, "3-D millimeter-wave statistical channel model for 5G wireless system design," *IEEE Transactions on Microwave Theory and Techniques*, vol. 64, no. 7, pp. 2207–2225, July 2016.

[97] NYU Wireless, "NYUSIM: The Open Source 5G Channel Model Simulator software," 2016. [Online]. Available: http://wireless.engineering.nyu.edu/5g-millimeter-wave-channel-modeling-software/

[98] M. Giordani, M. Mezzavilla, A. Dhananjay, S. Rangan, and M. Zorzi, "Channel dynamics and SNR tracking in millimeter wave cellular systems," in *22th European Wireless Conference*. VDE, 2016.

[99] P. A. Eliasi, S. Rangan, and T. S. Rappaport, "Low-Rank Spatial Channel Estimation for Millimeter Wave Cellular Systems," *IEEE Transactions on Wireless Communications*, vol. 16, no. 5, pp. 2748–2759, May 2017.

[100] D. J. Love and R. W. Heath, "Equal gain transmission in multiple-input multiple-output wireless systems," *IEEE Transactions on Communications*, vol. 51, no. 7, pp. 1102–1110, July 2003.

[101] J. H. Wilkinson, *The Algebraic Eigenvalue Problem*. Clarendon Press Oxford, 1965, vol. 87.

[102] J. Wang, "Beam codebook based beamforming protocol for multi-Gbps millimeter-wave WPAN systems," *IEEE Journal on Selected Areas in Communications*, vol. 27, no. 8, pp. 1390–1399, October 2009.

[103] R. W. Heath, N. Gonzalez-Prelcic, S. Rangan, W. Roh, and A. M. Sayeed, "An Overview of Signal Processing Techniques for Millimeter Wave MIMO Systems," *IEEE Journal of Selected Topics in Signal Processing*, vol. 10, no. 3, pp. 436–453, Apr. 2016.

[104] M. Rebato, J. Park, P. Popovski, E. D. Carvalho, and M. Zorzi, "Stochastic Geometric Coverage Analysis in mmWave Cellular Networks with a Realistic Channel Model," in *GLOBECOM 2017 - IEEE Global Communications Conference*, Dec 2017.

219

[105] M. Mezzavilla, M. Miozzo, M. Rossi, N. Baldo, and M. Zorzi, "A Lightweight and Accurate Link Abstraction Model for the Simulation of LTE Networks in ns-3," in *Proceedings of the 15th ACM International Conference on Modeling, Analysis and Simulation of Wireless and Mobile Systems*, ser. MSWiM '12, 2012, pp. 55–60.

[106] "The Vienna LTE Simulators." [Online]. Available: https://www.nt.tuwien.ac.at/research/mobile-communications/vccs/vienna-lte-a-simulators/lte-advanced-link-level-simulator/

[107] Korea Telecom (KT) 5G-SIG, "5G.201, KT 5th Generation Radio Access; Physical Layer; General description (5G pre-specifications)," 2016.

[108] M. Cudak, T. Kovarik, T. A. Thomas, A. Ghosh, Y. Kishiyama, and T. Nakamura, "Experimental mmWave 5G cellular system," in *Globecom Workshops (GC Wkshps), 2014*. IEEE, 2014, pp. 377–381.

[109] T. Levanen, J. Pirskanen, and M. Valkama, "Radio interface design for ultra-low latency millimeter-wave communications in 5G era," in *Proc. IEEE Globecom Workshops (Gc Wkshps)*, Dec. 2014, pp. 1420–1426.

[110] S. Dutta, M. Mezzavilla, R. Ford, M. Zhang, S. Rangan, and M. Zorzi, "Frame structure design and analysis for millimeter wave cellular systems," *IEEE Transactions on Wireless Communications*, vol. 16, no. 3, pp. 1508–1522, March 2017.

[111] ——, "MAC layer frame design for millimeter wave cellular system," in *European Conference on Networks and Communications (EuCNC)*, June 2016, pp. 117–121.

[112] D. Astély, E. Dahlman, A. Furuskär, Y. Jading, M. Lindström, and S. Parkvall, "LTE: the evolution of mobile broadband," *IEEE Communications magazine*, vol. 47, no. 4, May 2009.

[113] Verizon, "5G TF; Air Interface Working Group; Verizon 5th Generation Radio Access; Physical channels and modulation (Release 1)," October 2016. [Online]. Available: http://www.5gtf.net/V5G_211_v1p7.pdf

[114] 3GPP, "Study on New Radio (NR) Access Technology - Physical Layer Aspects," TR 38.802, V14.0.0, 2017.

[115] P. Popovski, V. Brau, H.-P. Mayer, P. Fertl, Z. Ren, D. Gonzales-Serrano, E. G. Ström, T. Svensson, H. Taoka, P. Agyapong *et al.*, "EU FP7 INFSO-ICT-317669 METIS, D1. 1 scenarios, requirements and KPIs for 5G mobile and wireless system," 2013.

[116] P. Kela, M. Costa, J. Salmi, K. Leppanen, J. Turkka, T. Hiltunen, and M. Hronec, "A novel radio frame structure for 5G dense outdoor radio access networks," in *Proc. IEEE 81st Vehicular Technology Conference (VTC Spring)*, May 2015, pp. 1–6.

[117] S. Choi and K. G. Shin, "A class of adaptive hybrid ARQ schemes for wireless links," *IEEE Transactions on Vehicular Technology*, vol. 50, no. 3, pp. 777–790, May 2001.

[118] 3GPP, "Medium Access Control (MAC) protocol specification - Release 14," TS 36.321, 2016.

[119] S. Sesia, M. Baker, and I. Toufik, *LTE-the UMTS long term evolution: from theory to practice*. John Wiley & Sons, 2011.

[120] E. Dahlman, S. Parkvall, J. Sköld, and P. Beming, *4G LTE/LTE-Advanced for Mobile Broadband*. Oxford, UK: Academic Press, 201.

[121] 3GPP, "Technical Specification Group Services and System Aspects; Policy and charging control architecture," TR 23.203, V14.0.0, 2017.

[122] B. Bojovic, M. D. Abrignani, M. Miozzo, L. Giupponi, and N. Baldo, "Towards LTE-Advanced and LTE-A Pro Network Simulations: Implementing Carrier Aggregation in LTE Module of ns-3," in *Proceedings of the Workshop on ns-3*, ser. WNS3 '17, Porto, Portugal, 2017, pp. 63–70.

[123] ITU-R World Radiocommunication Conference - Geneva, "Studies on Frequency-related Matters for International Mobile Telecommunications Identification Including Possible Additional Allocations to the Mobile Services on a Primary basis in Portion(s) of the Frequency Range Between 24.25 and 86 GHz for the Future Development of International Mobile Telecommunications for 2020 and Beyond," Resolution 238 (WRC-15), 2015. [Online]. Available: https://www.itu.int/dms_pub/itu-r/oth/0c/0a/R0C0A00000C0014PDFE.pdf

[124] G. Abbas, Z. Halim, and Z. H. Abbas, "Fairness-driven queue management: A survey and taxonomy," *IEEE Communications Surveys & Tutorials*, vol. 18, no. 1, pp. 324–367, Jan. 2016.

[125] F. Baker and G. Fairhurst, "IETF recommendations regarding Active Queue Management," RFC 7567, 2015.

[126] P. Imputato and S. Avallone, "Design and implementation of the traffic control module in ns-3," in *Proceedings of the Workshop on ns-3*. ACM, 2016, pp. 1–8.

[127] ——, "Traffic differentiation and multiqueue networking in ns-3," in *Proceedings of the Workshop on ns-3*. ACM, 2017, pp. 79–86.

[128] A. Deepak, K. Shravya, and M. P. Tahiliani, "Design and Implementation of AQM Evaluation Suite for ns-3," in *Proceedings of the Workshop on ns-3*. ACM, 2017, pp. 87–94.

[129] J. Gettys and K. Nichols, "Bufferbloat: Dark buffers in the internet," *ACM Queue*, vol. 9, no. 11, pp. 40:40–40:54, Nov 2011.

[130] S. Floyd and V. Jacobson, "Random early detection gateways for congestion avoidance," *IEEE/ACM Transactions on networking*, vol. 1, no. 4, pp. 397–413, Aug. 1993.

[131] T. J. Ott, T. Lakshman, and L. H. Wong, "SRED: stabilized RED," in *INFOCOM'99. Eighteenth Annual Joint Conference of the IEEE Computer and Communications Societies*, vol. 3. IEEE, 1999, pp. 1346–1355.

[132] K. Nichols and V. Jacobson, "Controlling queue delay," *Communications of the ACM*, vol. 55, no. 7, pp. 42–50, July 2012.

[133] S. Chandrashekar, A. Maeder, C. Sartori, T. Höhne, B. Vejlgaard, and D. Chandramouli, "5G multi-RAT multi-connectivity architecture," in *IEEE International Conference on Communications Workshops (ICC)*, May 2016, pp. 180–186.

[134] J. G. Rois, B. Lorenzo, F. J. Gonzalez-Castano, and J. C. Burguillo, "Heterogeneous millimeter-wave/micro-wave architecture for 5G wireless access and backhauling," in *European Conference on Networks and Communications (EuCNC)*, June 2016, pp. 179–184.

[135] AT&T, "Migration and Interworking Aspects - SA WG2 Temporary Document S2-163348," July 2016. [Online]. Available: http://www.3gpp.org/ftp/tsg_sa/WG2_Arch/TSGS2_116_Vienna/Docs/S2-163348.zip

[136] 3GPP, "Study on small cell enhancements for E-UTRA and E-UTRAN," TR 36.842, V12.0.0, 2013.

[137] I. D. Silva, G. Mildh, J. Rune, P. Wallentin, J. Vikberg, P. Schliwa-Bertling, and R. Fan, "Tight Integration of New 5G Air Interface and LTE to Fulfill 5G Requirements," in *IEEE 81st Vehicular Technology Conference (VTC Spring)*, May 2015.

[138] B. Nguyen, A. Banerjee, V. Gopalakrishnan, S. Kasera, S. Lee, A. Shaikh, and J. Van der Merwe, "Towards Understanding TCP Performance on LTE/EPC Mobile Networks," in *Proceedings of the 4th Workshop on All Things Cellular: Operations, Applications, and Challenges*, ser. AllThingsCellular '14. Chicago, Illinois, USA: ACM, 2014, pp. 41–46. [Online]. Available: http://doi.acm.org/10.1145/2627585.2627594

[139] A. Ford, C. Raiciu, M. Handley, S. Barre, and J. Iyengar, "Architectural Guidelines for Multipath TCP Development," RFC 6182, 2011.

[140] B. Chihani and C. Denis, "A Multipath TCP model for ns-3 simulator," in *Workshop on ns-3 held in conjunction with SIMUTools 2011*, 2011.

[141] M. Coudron and S. Secci, "An implementation of multipath TCP in ns-3," *Computer Networks*, vol. 116, pp. 1–11, 2017.

[142] C. Paasch, S. Barre, and al., "Multipath TCP in the Linux Kernel," available at http://www.multipath-tcp.org.

[143] "Iperf 2.0," available at https://iperf.fr.

[144] H. Tazaki, R. Nakamura, and Y. Sekiya, "Library operating system with mainline linux network stack," in *Proceedings of Netdev 0.1*, 2015.

[145] Q. Peng, A. Walid, J. Hwang, and S. H. Low, "Multipath TCP: Analysis, Design, and Implementation," *IEEE/ACM Transactions on Networking*, vol. 24, no. 1, pp. 596–609, Feb 2016.

[146] A. Kassler and M. Pieskä, "TCP Performance over 5G mmWave Links-Tradeoff between Capacity and Latency," in *The 13th IEEE International Conference on Wireless and Mobile Computing, Networking and Communications, 9.11 th October 2017 Rome Italy*, 2017, pp. 385–394.

[147] M. Kim, S. W. Ko, and S. L. Kim, "Enhancing TCP end-to-end performance in millimeter-wave communications," in *IEEE 28th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC)*, Oct 2017.

[148] Q. Hu, Y. Liu, Y. Yan, and D. M. Blough, "End-to-end Simulation of mmWave Out-of-band Backhaul Networks in ns-3," in *Proceedings of the 2019 Workshop on Next-Generation Wireless with ns-3*, ser. WNGW 2019. New York, NY, USA: ACM, 2019, pp. 1–4. [Online]. Available: http://doi.acm.org/10.1145/3337941.3337943

[149] M. Giordani, A. Zanella, and M. Zorzi, "LTE and Millimeter Waves for V2I Communications: An End-to-End Performance Comparison," in *IEEE 89th Vehicular Technology Conference (VTC2019-Spring)*, April 2019, pp. 1–7.

[150] P. J. Mateo, C. Fiandrino, and J. Widmer, "Analysis of TCP Performance in 5G mm-Wave Mobile Networks," in *IEEE International Conference on Communications (ICC)*, May 2019, pp. 1–7.

[151] Y. Hou, Z. Wen'an, L. Song, and M. Gao, "A QoE Estimation Model for Video Streaming over 5G Millimeter Wave Network," in *International Conference on Broadband and Wireless Computing, Communication and Applications.* Springer, 2016, pp. 93–104.

[152] O. Semiari, W. Saad, and M. Bennis, "Joint millimeter wave and microwave resources allocation in cellular networks with dual-mode base stations," *IEEE Transactions on Wireless Communications*, vol. 16, no. 7, pp. 4802–4816, July 2017.

[153] ——, "Downlink cell association and load balancing for joint millimeter wave-microwave cellular networks," in *IEEE Global Communications Conference (GLOBECOM)*, Dec. 2016, pp. 1–6.

[154] M. E. Rasekh, Z. Marzi, Y. Zhu, U. Madhow, and H. Zheng, "Noncoherent mmwave path tracking," in *Proceedings of the 18th International Workshop on Mobile Computing Systems and Applications.* New York, NY, USA: ACM, 2017, pp. 13–18. [Online]. Available: http://doi.acm.org/10.1145/3032970.3032974

[155] O. Semiari, W. Saad, M. Bennis, and B. Maham, "Caching Meets Millimeter Wave Communications for Enhanced Mobility Management in 5G Networks," *IEEE Transactions on Wireless Communications*, vol. 17, no. 2, pp. 779–793, Feb 2018.

[156] C. Herranz, M. Zhang, M. Mezzavilla, D. Martin-Sacristán, S. Rangan, and J. F. Monserrat, "A 3GPP NR Compliant Beam Management Framework to Simulate End-to-End mmWave Networks," in *Proceedings of the 21st ACM International Conference on Modeling, Analysis and Simulation of Wireless and Mobile Systems*, ser. MSWIM '18. New York, NY, USA: ACM, 2018, pp. 119–125. [Online]. Available: http://doi.acm.org/10.1145/3242102.3242117

[157] J. Palacios, D. De Donno, and J. Widmer, "Tracking mm-Wave channel dynamics: Fast beam training strategies under mobility," in *IEEE Conference on Computer Communications (INFOCOM)*. IEEE, 2017.

[158] A. Tassi, M. Egan, R. J. Piechocki, and A. Nix, "Modeling and design of millimeter-wave networks for highway vehicular communication," *IEEE Transactions on Vehicular Technology*, vol. 66, no. 12, pp. 10 676–10 691, Dec. 2017.

[159] I. Mavromatis, A. Tassi, R. J. Piechocki, and A. Nix, "mmWave System for Future ITS: A MAC-Layer Approach for V2X Beam Steering," in *IEEE 86th Vehicular Technology Conference (VTC-Fall)*, Sept 2017, pp. 1–6.

[160] I. Mavromatis, A. Tassi, R. J. Piechocki, and A. R. Nix, "Beam alignment for millimetre wave links with motion prediction of autonomous vehicles," *Proc. of IET Colloquium on Antennas, Propagation & RF Technology for Transport and Autonomous Platforms*, 2017.

[161] M. C. Gonzalez, C. Hidalgo, and A. Barabasi, "Understanding individual human mobility patterns," *Nature*, vol. 453, pp. 779–782, 07 2008.

[162] J. Harri, F. Filali, and C. Bonnet, "Mobility models for vehicular ad hoc networks: a survey and taxonomy," *IEEE Communications Surveys & Tutorials*, vol. 11, no. 4, pp. 19–41, Fourth quarter 2009.

[163] M. Giordani, M. Mezzavilla, S. Rangan, and M. Zorzi, "An Efficient Uplink Multi-Connectivity Scheme for 5G Millimeter-Wave Control Plane Applications," *IEEE Transactions on Wireless Communications*, vol. 17, no. 10, pp. 6806–6821, Oct 2018.

[164] C. N. Barati, S. A. Hosseini, S. Rangan, P. Liu, T. Korakis, S. S. Panwar, and T. S. Rappaport, "Directional cell discovery in millimeter wave cellular networks," *IEEE Transactions on Wireless Communications*, vol. 14, no. 12, pp. 6664–6678, December 2015.

[165] C. N. Barati, S. A. Hosseini, M. Mezzavilla, P. Amiri-Eliasi, S. Rangan, T. Korakis, S. S. Panwar, and M. Zorzi, "Directional initial access for millimeter wave cellular systems," in *49th Asilomar Conference on Signals, Systems and Computers.* IEEE, 2015, pp. 307–311.

[166] A. Zakrzewska, D. Lopez-Perez, S. Kucera, and H. Claussen, "Dual connectivity in LTE HetNets with split control- and user-plane," in *IEEE Globecom Workshops (GC Wkshps)*, Dec 2013, pp. 391–396.

[167] Z. He, S. Mao, and T. S. Rappaport, "Minimum time length link scheduling under blockage and interference in 60 GHz networks," in *IEEE Wireless Communications and Networking Conference (WCNC)*, March 2015, pp. 837–842.

[168] X. Yan, Y. A. Sekercioglu, and S. Narayanan, "A survey of vertical handover decision algorithms in Fourth Generation heterogeneous wireless networks," *Computer Networks*, vol. 54, no. 11, pp. 1848 – 1863, 2010. [Online]. Available: http://www.sciencedirect.com/science/article/pii/S1389128610000502

[169] F. Guidolin, I. Pappalardo, A. Zanella, and M. Zorzi, "Context-aware handover policies in HetNets," *IEEE Transactions on Wireless Communications*, vol. 15, no. 3, pp. 1895–1906, March 2016.

[170] A. Talukdar, M. Cudak, and A. Ghosh, "Handoff rates for millimeterwave 5G systems," in *IEEE 79th Vehicular Technology Conference (VTC Spring)*, May 2014.

[171] H. Song, X. Fang, and L. Yan, "Handover scheme for 5G C/U plane split heterogeneous network in high-speed railway," *IEEE Transactions on Vehicular Technology*, vol. 63, no. 9, pp. 4633–4646, Nov 2014.

[172] S. Sadr and R. S. Adve, "Handoff rate and coverage analysis in multi-tier heterogeneous networks," *IEEE Transactions on Wireless Communications*, vol. 14, no. 5, pp. 2626–2638, May 2015.

[173] P. Coucheney, E. Hyon, and J. M. Kelif, "Mobile association problem in heterogenous wireless networks with mobility," in *IEEE 24th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC)*, Sept 2013, pp. 3129–3133.

[174] V. Yazici, U. C. Kozat, and M. O. Sunay, "A new control plane for 5G network architecture with a case study on unified handoff, mobility, and routing management," *IEEE Communications Magazine*, vol. 52, no. 11, pp. 76–85, Nov 2014.

[175] M. Giordani, M. Mezzavilla, C. N. Barati, S. Rangan, and M. Zorzi, "Comparative analysis of initial access techniques in 5G mmWave cellular networks," in *Information Science and Systems (CISS), 2016 Annual Conference on*. IEEE, 2016, pp. 268–273.

[176] K. Chandra, R. V. Prasad, I. G. Niemegeers, and A. R. Biswas, "Adaptive beamwidth selection for contention based access periods in millimeter wave WLANs," in *IEEE 11th Consumer Communications and Networking Conference (CCNC)*. IEEE, 2014, pp. 458–464.

[177] M. Samimi, T. S. Rappaport, and G. MacCartney, "Probabilistic omnidirectional path loss models for millimeter-wave outdoor communications," *IEEE Wireless Commu. Letters*, vol. 4, no. 4, pp. 357–360, 2015.

[178] G. MacCartney, T. S. Rappaport, M. Samimi, and S. Sun, "Wideband millimeter-wave propagation measurements and channel models for future wireless communication system design," *IEEE Trans. Comm.*, vol. 63, no. 9, pp. 3029–3056, 2015.

[179] P. A. Eliasi and S. Rangan, "Stochastic dynamic channel models for millimeter cellular systems," in *Proc. IEEE Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*. IEEE, 2015, pp. 209–212.

[180] A. Dhananjay, "Iris: Mitigating phase noise in millimeter wave OFDM systems," Ph.D. dissertation, New York University (NYU), 2015.

[181] S. Schwarz, C. Mehlführer, and M. Rupp, "Calculation of the spatial preprocessing and link adaption feedback for 3GPP UMTS/LTE," in *6th conference on Wireless advanced (WiAD)*. IEEE, 2010.

[182] Next Generation Mobile Networks Alliance, "Optimised backhaul requirements," Tech. Rep., 2008. [Online]. Available: https://www.ngmn.org/uploads/media/NGMN_Optimised_Backhaul_Requirements.pdf

[183] M. Giordani, M. Mezzavilla, and M. Zorzi, "Initial access in 5G mmWave cellular networks," *IEEE Communications Magazine*, vol. 54, no. 11, pp. 40–47, 2016.

[184] M. Giordani and M. Zorzi, "Improved user tracking in 5G millimeter wave mobile networks via refinement operations," in *16th Annual Mediterranean Ad Hoc Networking Workshop (Med-Hoc-Net)*, June 2017.

[185] T. Nitsche, C. Cordeiro, A. Flores, E. Knightly, E. Perahia, and J. Widmer, "IEEE 802.11ad: directional 60 GHz communication for multi-Gigabit-per-second Wi-Fi [Invited Paper]," *IEEE Communications Magazine*, vol. 52, no. 12, pp. 132–141, December 2014.

[186] R. Santosa, B.-S. Lee, C. K. Yeo, and T. M. Lim, "Distributed Neighbor Discovery in Ad Hoc Networks Using Directional Antennas," in *The Sixth IEEE International Conference on Computer and Information Technology*, Sept 2006, pp. 97–97.

[187] D. Liu, L. Wang, Y. Chen, M. Elkashlan, K.-K. Wong, R. Schober, and L. Hanzo, "User association in 5G networks: A survey and an outlook," *IEEE Communications Surveys & Tutorials*, vol. 18, no. 2, pp. 1018–1044, 2016.

[188] C. Jeong, J. Park, and H. Yu, "Random access in millimeter-wave beamforming cellular networks: issues and approaches," *IEEE Communications Magazine*, vol. 53, no. 1, pp. 180–185, January 2015.

[189] V. Desai, L. Krzymien, P. Sartori, W. Xiao, A. Soong, and A. Alkhateeb, "Initial beamforming for mmWave communications," in *48th Asilomar Conference on Signals, Systems and Computers*, 2014, pp. 1926–1930.

[190] L. Wei, Q. Li, and G. Wu, "Exhaustive, Iterative and Hybrid Initial Access Techniques in mmWave Communications," in *IEEE Wireless Communications and Networking Conference (WCNC)*. IEEE, 2017.

[191] J. Choi, "Beam selection in mm-Wave multiuser MIMO systems using compressive sensing," *IEEE Transactions on Communications*, vol. 63, no. 8, pp. 2936–2947, August 2015.

[192] A. Capone, I. Filippini, and V. Sciancalepore, "Context information for fast cell discovery in mm-wave 5G networks," in *21th European Wireless Conference; Proceedings of European Wireless*, 2015.

[193] A. Capone, I. Filippini, V. Sciancalepore, and D. Tremolada, "Obstacle avoidance cell discovery using mm-waves directive antennas in 5G networks," in *IEEE 26th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC)*. IEEE, 2015, pp. 2349–2353.

[194] Q. C. Li, H. Niu, G. Wu, and R. Q. Hu, "Anchor-booster based heterogeneous networks with mmWave capable booster cells," in *IEEE Globecom Workshops*. IEEE, 2013, pp. 93–98.

[195] W. B. Abbas and M. Zorzi, "Context information based initial cell search for millimeter wave 5G cellular networks," in *European Conference on Networks and Communications (EuCNC)*. IEEE, 2016, pp. 111–116.

[196] A. Alkhateeb, Y. H. Nam, M. S. Rahman, J. Zhang, and R. W. Heath, "Initial Beam Association in Millimeter Wave Cellular Systems: Analysis and Design Insights," *IEEE Transactions on Wireless Communications*, vol. 16, no. 5, pp. 2807–2821, May 2017.

[197] Y. Li, J. Luo, M. Castaneda, R. Stirling-Gallacher, W. Xu, and G. Caire, "On the Beamformed Broadcast Signaling for Millimeter Wave Cell Discovery: Performance Analysis and Design Insight," *arXiv preprint arXiv:1709.08483*, 2017.

[198] L. Wei, Q. C. Li, and G. Wu, "Initial Access Techniques for 5G NR: Omni/Beam SYNC and RACH designs," in *International Conference on Computing, Networking and Communications (ICNC)*, March 2018, pp. 249–253.

[199] A. S. Cacciapuoti, "Mobility-Aware User Association for 5G mmWave Networks," *IEEE Access*, vol. 5, pp. 21 497–21 507, 2017.

[200] S. Jayaprakasam, X. Ma, J. W. Choi, and S. Kim, "Robust Beam-Tracking for mmWave Mobile Communications," *IEEE Communications Letters*, vol. 21, no. 12, pp. 2654–2657, Dec 2017.

[201] N. Gonzalez-Prelcic, A. Ali, V. Va, and R. W. Heath, "Millimeter-Wave Communication with Out-of-Band Information," *IEEE Communications Magazine*, vol. 55, no. 12, pp. 140–146, Dec. 2017.

[202] J. Liu, K. Au, A. Maaref, J. Luo, H. Baligh, H. Tong, A. Chassaigne, and J. Lorca, "Initial Access, Mobility, and User-Centric Multi-Beam Operation in 5G New Radio," *IEEE Communications Magazine*, vol. 56, no. 3, pp. 35–41, Mar 2018.

[203] E. Onggosanusi, M. S. Rahman, L. Guo, Y. Kwak, H. Noh, Y. Kim, S. Faxer, M. Harrison, M. Frenne, S. Grant, R. Chen, R. Tamrakar, and a. Q. Gao, "Modular and High-Resolution Channel State Information and Beam Management for 5G New Radio," *IEEE Communications Magazine*, vol. 56, no. 3, pp. 48–55, Mar 2018.

[204] 3GPP, "NR - Radio Resource Control (RRC) protocol specification," TR 38.331, V15.0.0, 2017.

[205] ——, "NR - Physical layer procedures for control," TS 38.213, V15.0.0, 2018.

[206] ——, "NR - Physical layer measurements," TR 38.215, V15.0.0, 2017.

[207] ——, "Remaining details on sync signals," Samsung - Tdoc R1-1721434, 2017.

[208] ——, "Discussion on remaining issues of SS block and SS burst set," Motorola Mobility, Lenovo - Tdoc R1-1714212, 2017.

[209] ——, "SS burst periodicity for initial cell selection in NR," Nokia, Alcatel-Lucent Shanghai Bell - Tdoc R4-1705123, 2017.

[210] ——, "Measurement configuration and reporting considering additional RS," Huawei - Tdoc R2-1703387, 2017.

[211] ——, "Measurement configuration for CSI-RS," Ericsson - Tdoc R2-1704103, 2017.

[212] ——, "NR CSI-RS configuration for RRM measurement," Samsung - Tdoc R2-1709593, 2017.

[213] ——, "Summary of offline discussions on CSI-RS," Huawei, HiSilicom - Tdoc R1-1718947, 2017.

[214] ——, "Measurement reporting for NR SS and CSI-RS ," Huawei - Tdoc R2-1708703, 2017.

[215] ——, "Details of cell quality derivation," Ericsson - Tdoc R2-1704101, 2017.

[216] ——, "Consideration on CSI RS for beam management," ZTE Corporation - Tdoc R2-1708123, 2017.

[217] ——, "Remaining issues on SRS," InterDigital, Inc. - Tdoc R1-1716472, 2017.

[218] ——, "Discussion on remaining issues on SRS design," CATT - Tdoc R1-1712386, 2017.

[219] Ericsson, "5G New Radio: designing for the future," *Ericsson Technology Review*, 2017.

[220] 3GPP, "Framework for beamformed access," Samsung - Tdoc R1-164013, 2016.

[221] ——, "NR - Physical layer procedures for data," TS 38.214, V15.0.0, 2018.

[222] ——, "Discussion on CSI-RS Design," Qualcomm - Tdoc R1-1718546, 2017.

[223] ——, "Study on new radio access technology: Radio access architecture and interfaces," TR 38.801, V14.0.0, 2017.

[224] J. Oueis and E. C. Strinati, "Uplink traffic in future mobile networks: Pulling the alarm," in *International Conference on Cognitive Radio Oriented Wireless Networks*. Springer, 2016, pp. 583–593.

[225] 3GPP, "NR PRACH preamble resource allocation," Ericsson - Tdoc R1-1611905, 2016.

[226] M. Rebato, L. Resteghini, C. Mazzucco, and M. Zorzi, "Study of Realistic Antenna Patterns in 5G mmWave Cellular Scenarios," in *IEEE International Conference on Communications (ICC)*, May 2018. [Online]. Available: https://arxiv.org/abs/1802.01316

[227] Z. Li, S. Han, and A. F. Molisch, "Hybrid beamforming design for millimeter-wave multi-user massive MIMO downlink," in *IEEE International Conference on Communications (ICC)*, May 2016, pp. 1–6.

[228] A. F. Molisch, V. V. Ratnam, S. Han, Z. Li, S. L. H. Nguyen, L. Li, and K. Haneda, "Hybrid Beamforming for Massive MIMO: A Survey," *IEEE Communications Magazine*, vol. 55, no. 9, pp. 134–141, September 2017.

[229] T. Yoo and A. Goldsmith, "Optimality of zero-forcing beamforming with multiuser diversity," in *IEEE International Conference on Communications (ICC)*, vol. 1, May 2005, pp. 542–546.

[230] S. Kutty and D. Sen, "Beamforming for Millimeter Wave Communications: An Inclusive Survey," *IEEE Communications Surveys & Tutorials*, vol. 18, no. 2, pp. 949–973, Second quarter 2016.

[231] A. Valdes-Garcia, S. T. Nicolson, J.-W. Lai, A. Natarajan, P.-Y. Chen, S. K. Reynolds, J.-H. C. Zhan, D. G. Kam, D. Liu, and B. Floyd, "A fully integrated 16-element phased-array transmitter in SiGe BiCMOS for 60-GHz communications," *IEEE journal of solid-state circuits*, vol. 45, no. 12, pp. 2757–2773, 2010.

[232] E. Cohen, M. Ruberto, M. Cohen, O. Degani, S. Ravid, and D. Ritter, "A CMOS bidirectional 32-element phased-array transceiver at 60 GHz with LTCC antenna," *IEEE Transactions on Microwave Theory and Techniques*, vol. 61, no. 3, pp. 1359–1375, 2013.

[233] X. Gao, L. Dai, and A. M. Sayeed, "Low RF-Complexity Technologies to Enable Millimeter-Wave MIMO with Large Antenna Array for 5G Wireless Communications," *IEEE Communications Magazine*, vol. 56, no. 4, pp. 211–217, APRIL 2018.

[234] 3GPP, "Discussion on NR 4-Step Random Access Procedure," Ericsson - Tdoc R1-1718052, 2017.

[235] ——, "Relation between radio link failure and beam failure," Ericsson - Tdoc R1-1705917, 2017.

[236] ——, "Beam failure detection and beam recovery actions," Ericsson - Tdoc R1-1705893, 2017.

[237] ——, "LS on NR PRACH BW Aspects," Tdoc R1-1716814, 2017.

[238] W. B. Abbas, F. Gomez-Cuba, and M. Zorzi, "Millimeter Wave Receiver Efficiency: A Comprehensive Comparison of Beamforming Schemes With Low Resolution ADCs," *IEEE Transactions on Wireless Communications*, vol. 16, no. 12, pp. 8131–8146, Dec 2017.

[239] Y. Zeng, J. Lyu, and R. Zhang, "Cellular-connected UAV: Potential, challenges, and promising technologies," *IEEE Wireless Communications*, vol. 26, no. 1, pp. 120–127, Feb. 2018.

[240] M. Moradi, K. Sundaresan, E. Chai, S. Rangarajan, and Z. Mao, "SkyCore: Moving core to the edge for untethered and reliable UAV-based LTE networks," in *Proc. of ACM MobiCom*, New Delhi, India, Oct. 2018.

[241] A. Merwaday and I. Guvenc, "UAV assisted heterogeneous networks for public safety communications," in *Proc. of IEEE WCNCW*, New Orleans, LA, USA, Mar. 2015.

[242] S. Singh, M. N. Kulkarni, A. Ghosh, and J. G. Andrews, "Tractable model for rate in self-backhauled millimeter wave cellular networks," *IEEE Journal on Selected Areas in Communications*, vol. 33, no. 10, pp. 2196–2211, May 2015.

[243] Z. Xiao, P. Xia, and X. Xia, "Enabling UAV cellular with millimeter-wave communication: Potentials and approaches," *IEEE Communications Magazine*, vol. 54, no. 5, pp. 66–73, May 2016.

[244] T. Cuvelier and R. Heath, "mmWave MU-MIMO for aerial networks," in *Proc. of IEEE ISWCS*, Lisbon, Portugal, Aug. 2018.

[245] W. Khawaja, O. Ozdemir, and I. Guvenc, "Temporal and spatial characteristics of mmWave propagation channels for UAVs," in *Proc. of IEEE GSMM*, Boulder, CO, USA, May 2018.

[246] J. Zhao, G. Gao, L. Kuang, Q. Wu, and W. Jia, "Channel tracking with flight control system for UAV mmWave MIMO communications," *IEEE Communications Letters*, vol. 22, no. 6, pp. 1224–1227, June 2018.

[247] G. Bielsa, M. Mezzavilla, J. Widmer, and S. Rangan, "Performance assessment of off-the-shelf mmWave radios for drone communications," in *Proc. of IEEE WoWMoM*, Washington, DC, USA, June 2019.

[248] DJI Matrice 600 Pro. (2019) https://www.dji.com/matrice600-pro.

[249] Facebook, Terragraph. (2019) https://terragraph.com/product.

[250] mmBAC Demo. (2019) https://youtu.be/Swnf5JyfqY0.

[251] D. López-Pérez, M. Ding, H. Claussen, and A. H. Jafari, "Towards 1 Gbps/UE in Cellular Systems: Understanding Ultra-Dense Small Cell Deployments," *IEEE Communications Surveys & Tutorials*, vol. 17, no. 4, pp. 2078–2101, Fourth quarter 2015.

[252] R. Mijumbi, J. Serrat, J.-L. Gorricho, N. Bouten, F. De Turck, and R. Boutaba, "Network function virtualization: State-of-the-art and research challenges," *IEEE Communications Surveys & Tutorials*, vol. 18, no. 1, pp. 236–262, Sept 2016.

[253] H. S. Dhillon and G. Caire, "Wireless backhaul networks: Capacity bound, scalability analysis and design guidelines," *IEEE Transactions on Wireless Communications*, vol. 14, no. 11, pp. 6043–6056, 2015.

[254] 3GPP, "NR; Study on integrated access and backhaul," TR 38.874, V15.0.0, 2018.

[255] C. Saha, M. Afshang, and H. S. Dhillon, "Bandwidth Partitioning and Downlink Analysis in Millimeter Wave Integrated Access and Backhaul for 5G," *IEEE Transactions on Wireless Communications*, vol. 17, no. 12, pp. 8195–8210, Dec 2018.

[256] A. Ometov, D. Moltchanov, M. Komarov, S. V. Volvenko, and Y. Koucheryavy, "Packet Level Performance Assessment of mmWave Backhauling Technology for 3GPP NR Systems," *IEEE Access*, vol. 7, pp. 9860–9871, 2019.

[257] V. Gambiroza, B. Sadeghi, and E. W. Knightly, "End-to-end performance and fairness in multihop wireless backhaul networks," in *Proceedings of the 10th Annual International Conference on Mobile Computing and Networking*, ser. MobiCom '04.   ACM, 2004, pp. 287–301.

[258] M. Alicherry, R. Bhatia, and L. E. Li, "Joint channel assignment and routing for throughput optimization in multi-radio wireless mesh networks," in *Proceedings of the 11th Annual International Conference on Mobile Computing and Networking*, ser. MobiCom '05.   ACM, 2005.

[259] 3GPP, "Evolved Universal Terrestrial Radio Access (E-UTRA) and Evolved Universal Terrestrial Radio Access Network (E-UTRAN); Overall description," TS 36.300 (Rel. 15), 2018.

[260] X. Ge, H. Cheng, M. Guizani, and T. Han, "5G wireless backhaul networks: challenges and research advances," *IEEE Network*, vol. 28, no. 6, pp. 6–11, Nov 2014.

[261] 3GPP, "Study on Integrated Access and Backhaul for NR," AT&T, Qualcomm, Samsung - Tdoc RP-171880, 2017.

[262] 3GPP, "Service requirements for next generation new services and markets," TS 22.261 (Rel. 15), 2018.

[263] D. Yuan, H.-Y. Lin, J. Widmer, and M. Hollick, "Optimal Joint Routing and Scheduling in Millimeter-Wave Cellular Networks," in *IEEE Conference on Computer Communications (INFOCOM)*.   IEEE, 2018.

[264] Y. Niu, C. Gao, Y. Li, L. Su, D. Jin, Y. Zhu, and D. O. Wu, "Energy-efficient scheduling for mmwave backhauling of small cells in heterogeneous cellular networks," *IEEE Transactions on Vehicular Technology*, vol. 66, no. 3, pp. 2674–2687, March 2017.

[265] C. Saha, M. Afshang, and H. S. Dhillon, "Integrated mmWave Access and Backhaul in 5G: Bandwidth Partitioning and Downlink Analysis," in *IEEE International Conference on Communications (ICC)*, May 2018.

[266] A. Mesodiakaki, A. Kassler, E. Zola, M. Ferndahl, and T. Cai, "Energy efficient line-of-sight millimeter wave small cell backhaul: 60, 70, 80 or 140 GHz?" in *IEEE 17th International Symposium on A World of Wireless, Mobile and Multimedia Networks (WoWMoM)*, June 2016.

[267] 3GPP, "Way Forward - IAB Architecture for L2/3 relaying," Qualcomm Inc, KDDI, AT&T, Nokia, Nokia Shangai Bell, Huawi, Ericsson, Intel, LG Electronics, CMCC, Samsung - Tdoc R3-181502, 2018.

[268] J. Postel, "Transmission Control Protocol," IETF, RFC 793, Sep. 1981. [Online]. Available: https://rfc-editor.org/rfc/rfc793.txt

[269] K. Liu and J. Y. Lee, "On Improving TCP Performance over Mobile Data Networks," *IEEE Trans. Mobile Comput.*, vol. 15, no. 10, pp. 2522–2536, Oct. 2016.

[270] R. Li, M. Shariat, and M. Nekovee, "Transport protocols behaviour study in evolving mobile networks," in *IEEE International Symposium on Wireless Communication Systems (ISWCS)*, Sep 2016, pp. 456–460.

[271] J. Iyengar and M. Thomson, "QUIC: A UDP-based multiplexed and secure transport," IETF, Working Draft: draft-ietf-quic-transport-08, Dec. 2017. [Online]. Available: https://tools.ietf.org/id/draft-ietf-quic-transport-08.txt

[272] R. Stewart, "Stream Control Transport Protocol," IETF, RFC 6582, Sep. 2007. [Online]. Available: https://rfc-editor.org/rfc/rfc6582.txt

[273] M. Xiao, S. Mumtaz, Y. Huang, L. Dai, Y. Li, M. Matthaiou, G. K. Karagiannidis, E. Björnson, K. Yang, C. L. I, and A. Ghosh, "Millimeter Wave Communications for Future Mobile Networks," *IEEE J. on Sel. Areas Commun.*, vol. 35, no. 9, pp. 1909–1935, Sept 2017.

[274] C. Callegari, S. Giordano, M. Pagano, and T. Pepe, "A survey of congestion control mechanisms in Linux TCP," in *Distributed Computer and Communication Networks*. Springer, Mar 2014, pp. 28–42.

[275] A. Gurtov, T. Henderson, S. Floyd, and Y. Nishida, "The NewReno modification to TCP's Fast Recovery algorithm," IETF, RFC 6582, Apr. 2012. [Online]. Available: https://rfc-editor.org/rfc/rfc6582.txt

[276] S. Floyd, "HighSpeed TCP for Large Congestion Windows," RFC 3649, Dec. 2003. [Online]. Available: https://rfc-editor.org/rfc/rfc3649.txt

[277] S. Ha, I. Rhee, and L. Xu, "CUBIC: A new TCP-friendly high-speed TCP variant," *ACM Operating Systems Review*, vol. 42, no. 5, pp. 64–74, Jul. 2008.

[278] N. Cardwell, Y. Cheng, C. S. Gunn, S. H. Yeganeh, and V. Jacobson, "BBR: Congestion-based congestion control," *ACM Queue*, vol. 14, no. 5, pp. 20–53, Sep. 2016.

[279] S. Floyd, J. Mahdavi, M. Mathis, and D. A. Romanow, "TCP Selective Acknowledgment options," IETF, RFC 2018, Oct. 1996. [Online]. Available: https://rfc-editor.org/rfc/rfc2018.txt

[280] Y. Gong, D. Rossi, C. Testa, S. Valenti, and M. D. Täht, "Fighting the bufferbloat: on the coexistence of AQM and low priority congestion control," *Computer Networks*, vol. 65, pp. 255–267, June 2014.

[281] P. Mach and Z. Becvar, "Mobile edge computing: A survey on architecture and computation offloading," *IEEE Commun. Surveys Tuts.*, vol. 19, no. 3, pp. 1628–1656, third quarter 2017.

[282] D. Murray, T. Koziniec, K. Lee, and M. Dixon, "Large MTUs and internet performance," in *IEEE 13th International Conference on High Performance Switching and Routing*, June 2012, pp. 82–87.

[283] M. Mezzavilla, D. Chiarotto, D. Corujo, M. Wetterwald, and M. Zorzi, "Evaluation of Jumboframes feasibility in LTE access networks," in *IEEE International Conference on Communications (ICC)*, June 2013, pp. 5964–5968.

[284] D. A. Hayes, D. Ros, and Ö. Alay, "On the importance of TCP splitting proxies for future 5G mmWave communications," in *IEEE Conference on Local Computer Networks (LCN)*. IEEE, 2019.

[285] H. Balakrishnan, V. N. Padmanabhan, S. Seshan, and R. H. Katz, "A comparison of mechanisms for improving TCP performance over wireless links," *IEEE/ACM Transactions on Networking*, vol. 5, no. 6, pp. 756–769, Dec 1997.

[286] J. Border, M. Kojo, J. Griner, G. Montenegro, and Z. Shelby, "Performance Enhancing Proxies Intended to Mitigate Link-Related Degradations," RFC 3135, June 2001.

[287] J. H. Saltzer, D. P. Reed, and D. D. Clark, "End-to-end arguments in system design," *ACM Trans. Comput. Syst.*, vol. 2, no. 4, pp. 277–288, Nov 1984.

[288] F. Ren and C. Lin, "Modeling and improving TCP performance over cellular link with variable bandwidth," *IEEE Transactions on Mobile Computing*, vol. 10, no. 8, pp. 1057–1070, Aug 2011.

[289] K. Brown and S. Singh, "M-TCP: TCP for mobile cellular networks," *ACM SIGCOMM Computer Communication Review*, vol. 27, no. 5, pp. 19–43, Oct 1997.

[290] H. Balakrishnan, S. Seshan, E. Amir, and R. H. Katz, "Improving TCP/IP performance over wireless networks," in *Proceedings of the 1st annual international conference on Mobile computing and networking*. ACM, 1995, pp. 2–11.

[291] A. Bakre and B. Badrinath, "I-TCP: Indirect TCP for mobile hosts," in *Proceedings of the 15th International Conference on Distributed Computing Systems*. IEEE, 1995, pp. 136–143.

[292] M. Kim, S. Ko, H. Kim, S. Kim, and S. Kim, "Exploiting Caching for Millimeter-Wave TCP Networks: Gain Analysis and Practical Design," *IEEE Access*, vol. 6, pp. 69 769–69 781, 2018.

[293] 3GPP, "NR - Radio Link Control (RLC) protocol specification," TS 38.322, V15.0.0, 2017.

[294] D. Skordoulis, Q. Ni, H. H. Chen, A. P. Stephens, C. Liu, and A. Jamalipour, "IEEE 802.11n MAC frame aggregation mechanisms for next-generation high-throughput WLANs," *IEEE Wireless Communications*, vol. 15, no. 1, pp. 40–47, Feb 2008.

[295] D. Borman, R. T. Braden, V. Jacobson, and R. Scheffenegger, "TCP Extensions for High Performance," RFC 7323, Tech. Rep., Sep. 2014. [Online]. Available: https://rfc-editor.org/rfc/rfc7323.txt

[296] B. Veal, K. Li, and D. Lowenthal, "New Methods for Passive Estimation of TCP Round-trip Times," in *Proceedings of the 6th International Conference on Passive and Active Network Measurement*, ser. PAM'05, 2005, pp. 121–134.

[297] H. Jiang, Y. Wang, K. Lee, and I. Rhee, "Tackling bufferbloat in 3G/4G networks," *Proceedings of the 2012 ACM conference on Internet measurement*, pp. 329–342, 2012.

[298] F. Chiariotti, S. Kucera, A. Zanella, and H. Claussen, "Analysis and Design of a Latency Control Protocol for Multi-Path Data Delivery With Pre-Defined QoS Guarantees," *IEEE/ACM Transactions on Networking*, vol. 27, no. 3, pp. 1165–1178, June 2019.

[299] C. Raiciu, S. Barre, C. Pluntke, A. Greenhalgh, D. Wischik, and M. Handley, "Improving Datacenter Performance and Robustness with Multipath TCP," in *Proceedings of the ACM SIGCOMM 2011 Conference*, ser. SIGCOMM '11, 2011, pp. 266–277.

[300] V. Petrov, D. Solomitckii, A. Samuylov, M. A. Lema, M. Gapeyenko, D. Moltchanov, S. Andreev, V. Naumov, K. Samouylov, M. Dohler, and Y. Koucheryavy, "Dynamic multi-connectivity performance in ultra-dense urban mmwave deployments," *IEEE Journal on Selected Areas in Communications*, vol. 35, no. 9, pp. 2038–2055, Sep. 2017.

[301] A. Ford, C. Raiciu, M. Handley, and O. Bonaventure, "TCP Extensions for Multipath Operation with Multiple Addresses," RFC 6824, 2013.

[302] C. Raiciu, M. Handley, and D. Wischik, "Coupled congestion control for multipath transport protocols," RFC 6356, 2011.

[303] M. Scharf and A. Ford, "MPTCP Application Interface Considerations," RFC 6897, 2012.

[304] R. Khalili, N. Gast, M. Popovic, and J. Y. L. Boudec, "MPTCP is Not Pareto-Optimal: Performance Issues and a Possible Solution," *IEEE/ACM Transactions on Networking*, vol. 21, no. 5, pp. 1651–1665, Oct 2013.

[305] M. V. Pedersen, J. Heide, and F. H. Fitzek, "Kodo: An open and research oriented network coding library," in *International Conference on Research in Networking*. Springer, 2011, pp. 145–152.

[306] Measy, "60 GHz video transmitter," 2017. [Online]. Available: http://www.measy.com.cn/product/showproduct143_en.htm

[307] H. Singh, J. Oh, C. Kweon, X. Qin, H. R. Shao, and C. Ngo, "A 60 GHz wireless network for enabling uncompressed video communication," *IEEE Communications Magazine*, vol. 46, no. 12, pp. 71–78, December 2008.

[308] O. Abari, D. Bharadia, A. Duffield, and D. Katabi, "Enabling high-quality untethered virtual reality," in *14th USENIX Symposium on Networked Systems Design and Implementation (NSDI 17)*. USENIX Association, 2017, pp. 531–544.

[309] D. Vukobratović, C. Khirallah, V. Stanković, and J. S. Thompson, "Random Network Coding for Multimedia Delivery Services in LTE/LTE-Advanced," *IEEE Transactions on Multimedia*, vol. 16, no. 1, pp. 277–282, January 2014.

[310] J. Jin, B. Li, and T. Kong, "Is Random Network Coding Helpful in WiMAX?" in *IEEE INFOCOM 2008 - The 27th Conference on Computer Communications*, April 2008.

[311] A. Tassi, C. Khirallah, D. Vukobratović, F. Chiti, J. S. Thompson, and R. Fantacci, "Resource Allocation Strategies for Network-Coded Video Broadcasting Services Over LTE-Advanced," *IEEE Transactions on Vehicular Technology*, vol. 64, no. 5, pp. 2186–2192, May 2015.

[312] W. Song and W. Zhuang, "Performance analysis of probabilistic multipath transmission of video streaming traffic over multi-radio wireless devices," *IEEE Transactions on Wireless Communications*, vol. 11, no. 4, pp. 1554–1564, April 2012.

[313] H. Halbauer, P. Rugelandand, R. Tanoand, M. Tercero, A. Vijay, Y. Li, M. Filippouand, H. Miao, J. Widmerand, C. Fiandrino, I. Siaudand, A.-M. Ulmer-Moll, M. Shariat, J. Lorcaand, and Y. Zou, "Evaluations of the concepts for the 5G architecture and integration," mmMAGIC Deliverable D3.2, June 2017.

[314] H. Schulzrinne, A. Rao, R. Lanphier, M. Westerlund, and M. Stiemerling, "Real-Time Streaming Protocol Version 2.0," Internet Requests for Comments, RFC 7826, December 2016.

[315] R. Ahlswede, N. Cai, S. Y. R. Li, and R. W. Yeung, "Network information flow," *IEEE Transactions on Information Theory*, vol. 46, no. 4, pp. 1204–1216, July 2000.

[316] T. Ho, M. Medard, R. Koetter, D. R. Karger, M. Effros, J. Shi, and B. Leong, "A random linear network coding approach to multicast," *IEEE Transactions on Information Theory*, vol. 52, no. 10, pp. 4413–4430, October 2006.

[317] E. Magli, M. Wang, P. Frossard, and A. Markopoulou, "Network coding meets multimedia: A review," *IEEE Transactions on Multimedia*, vol. 15, no. 5, pp. 1195–1212, August 2013.

[318] J. Heide, M. V. Pedersen, F. H. Fitzek, and M. Médard, "On code parameters and coding vector representation for practical RLNC," in *IEEE International Conference on Communications (ICC)*, 2011.

[319] T. Wiegand, G. J. Sullivan, G. Bjontegaard, and A. Luthra, "Overview of the H.264/AVC video coding standard," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 13, no. 7, pp. 560–576, July 2003.

[320] H. Schwarz, D. Marpe, and T. Wiegand, "Overview of the Scalable Video Coding Extension of the H.264/AVC Standard," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 17, no. 9, pp. 1103–1120, September 2007.

[321] D. Vukobratovic and V. Stankovic, "Unequal error protection random linear coding strategies for erasure channels," *IEEE Transactions on Communications*, vol. 60, no. 5, pp. 1243–1252, May 2012.

[322] R. Julien, H. Schwarz, and M. Wien, "Joint Scalable Video Model 9.19.15 (JSVM 9.19.15)," 2011. [Online]. Available: https://github.com/floriandejonckheere/jsvm

[323] A. Detti, G. Bianchi, C. Pisa, F. S. Proto, P. Loreti, W. Kellerer, S. Thakolsri, and J. Widmer, "SVEF: an open-source experimental evaluation framework for H.264 scalable video streaming," in *IEEE Symposium on Computers and Communications*, July 2009, pp. 36–41.

[324] "FFmpeg," 2017. [Online]. Available: https://ffmpeg.org

[325] M. Iwamura, "NGMN view on 5G architecture," in *IEEE 81st Vehicular Technology Conference (VTC Spring)*, May 2015.

[326] M. Zorzi, A. Zanella, A. Testolin, M. D. F. D. Grazia, and M. Zorzi, "Cognition-based networks: A new perspective on network optimization using learning and distributed intelligence," *IEEE Access*, vol. 3, pp. 1512–1530, 2015.

[327] R. Li, Z. Zhao, X. Zhou, G. Ding, Y. Chen, Z. Wang, and H. Zhang, "Intelligent 5G: When Cellular Networks Meet Artificial Intelligence," *IEEE Wireless Communications*, vol. 24, no. 5, pp. 175–183, October 2017.

[328] S. Chinchali, P. Hu, T. Chu, M. Sharma, M. Bansal, R. Misra, M. Pavone, and K. Sachin, "Cellular network traffic scheduling with deep reinforcement learning," in *National Conference on Artificial Intelligence (AAAI)*, 2018.

[329] Y. He, F. R. Yu, N. Zhao, H. Yin, H. Yao, and R. C. Qiu, "Big data analytics in mobile cellular networks," *IEEE Access*, vol. 4, pp. 1985–1996, March 2016.

[330] O-RAN Alliance White Paper, "O-RAN: Towards an Open and Smart RAN," 2018. [Online]. Available: https://www.o-ran.org/resources

[331] F. Chiariotti, M. Condoluci, T. Mahmoodi, and A. Zanella, "Symbiocity: Smart cities for smarter networks," *Transactions on Emerging Telecommunications Technologies (ETT)*, June 2017, e3206 ett.3206. [Online]. Available: http://dx.doi.org/10.1002/ett.3206

[332] P. Hunt, D. Robertson, R. Bretherton, and M. C. Royle, "The SCOOT on-line traffic signal optimisation technique," *Traffic Engineering & Control*, vol. 23, no. 4, Apr. 1982.

[333] V. Pejovic and M. Musolesi, "Anticipatory mobile computing: A survey of the state of the art and research challenges," *ACM Comput. Surv.*, vol. 47, no. 3, pp. 47:1–47:29, Apr. 2015.

[334] C. Jiang, H. Zhang, Y. Ren, Z. Han, K. C. Chen, and L. Hanzo, "Machine learning paradigms for next-generation wireless networks," *IEEE Wireless Communications*, vol. 24, no. 2, pp. 98–105, April 2017.

[335] A. Imran, A. Zoha, and A. Abu-Dayya, "Challenges in 5G: how to empower SON with big data for enabling 5G," *IEEE Network*, vol. 28, no. 6, pp. 27–33, Nov 2014.

[336] K. Winstein and H. Balakrishnan, "TCP Ex Machina: Computer-generated Congestion Control," in *Proceedings of the ACM SIGCOMM 2013 Conference*.   Hong Kong, China: ACM, 2013, pp. 123–134.

[337] M. Gadaleta, F. Chiariotti, M. Rossi, and A. Zanella, "D-DASH: A Deep Q-Learning Framework for DASH Video Streaming," *IEEE Transactions on Cognitive Communications and Networking*, vol. 3, no. 4, pp. 703–718, Dec 2017.

[338] R. Becker, R. Cáceres, K. Hanson, S. Isaacman, J. M. Loh, M. Martonosi, J. Rowland, S. Urbanek, A. Varshavsky, and C. Volinsky, "Human mobility characterization from cellular network data," *Communications of the ACM*, vol. 56, no. 1, pp. 74–82, Jan 2013.

[339] R. A. Becker, R. Caceres, K. Hanson, J. M. Loh, S. Urbanek, A. Varshavsky, and C. Volinsky, "A tale of one city: Using cellular network data for urban planning," *IEEE Pervasive Computing*, vol. 10, no. 4, pp. 18–26, April 2011.

[340] W. Dong, N. Duffield, Z. Ge, S. Lee, and J. Pang, "Modeling cellular user mobility using a leap graph," in *International Conference on Passive and Active Network Measurement*. Springer, 2013, pp. 53–62.

[341] T. X. Tran, A. Hajisami, P. Pandey, and D. Pompili, "Collaborative Mobile Edge Computing in 5G Networks: New Paradigms, Scenarios, and Challenges," *IEEE Communications Magazine*, vol. 55, no. 4, pp. 54–61, April 2017.

[342] E. Bastug, M. Bennis, and M. Debbah, "Living on the edge: The role of proactive caching in 5G wireless networks," *IEEE Communications Magazine*, vol. 52, no. 8, pp. 82–89, Aug 2014.

[343] T. S. J. Darwish and K. A. Bakar, "Fog based intelligent transportation big data analytics in the internet of vehicles environment: Motivations, architecture, challenges, and critical issues," *IEEE Access*, vol. 6, pp. 15 679–15 701, 2018.

[344] M. Habib ur Rehman, P. P. Jayaraman, S. u. R. Malik, A. u. R. Khan, and M. Medhat Gaber, "RedEdge: A Novel Architecture for Big Data Processing in Mobile Edge Computing Environments," *Journal of Sensor and Actuator Networks*, vol. 6, no. 3, 2017. [Online]. Available: http://www.mdpi.com/2224-2708/6/3/17

[345] M. Peng, S. Yan, K. Zhang, and C. Wang, "Fog-computing-based radio access networks: issues and challenges," *IEEE Network*, vol. 30, no. 4, pp. 46–53, July 2016.

[346] T. Chen, M. Matinmikko, X. Chen, X. Zhou, and P. Ahokangas, "Software defined mobile networks: concept, survey, and research directions," *IEEE Communications Magazine*, vol. 53, no. 11, pp. 126–133, November 2015.

[347] L. E. Li, Z. M. Mao, and J. Rexford, "Toward software-defined cellular networks," in *European Workshop on Software Defined Networking*, Oct 2012, pp. 7–12.

[348] A. Gudipati, D. Perry, L. E. Li, and S. Katti, "SoftRAN: Software Defined Radio Access Network," in *Proceedings of the Second ACM SIGCOMM Workshop on Hot Topics in Software Defined Networking*, ser. HotSDN '13. Hong Kong, China: ACM, 2013, pp. 25–30.

[349] P. K. Agyapong, M. Iwamura, D. Staehle, W. Kiess, and A. Benjebbour, "Design considerations for a 5G network architecture," *IEEE Communications Magazine*, vol. 52, no. 11, pp. 65–75, Nov. 2014.

[350] X. An, F. Pianese, I. Widjaja, and U. G. Acer, "dMME: Virtualizing LTE mobility management," in *IEEE 36th Conference on Local Computer Networks (LCN)*, 2011, pp. 528–536.

[351] A. S. Rajan, S. Gobriel, C. Maciocco, K. B. Ramia, S. Kapury, A. Singhy, J. Ermanz, V. Gopalakrishnanz, and R. Jana, "Understanding the bottlenecks in virtualizing cellular core network functions," in *IEEE 21st International Workshop on Local and Metropolitan Area Networks*. IEEE, 2015, pp. 1–6.

[352] C. L. I, J. Huang, R. Duan, C. Cui, J. X. Jiang and L. Li, "Recent Progress on C-RAN Centralization and Cloudification," *IEEE Access*, vol. 2, pp. 1030–1039, Aug. 2014.

[353] P. Andres-Maldonado, P. Ameigeiras, J. Prados-Garzon, J. J. Ramos-Munoz, and J. M. Lopez-Soler, "Virtualized MME design for IoT support in 5G systems," *Sensors*, vol. 16, no. 8, p. 1338, 2016.

[354] N. Bui and J. Widmer, "Data-Driven Evaluation of Anticipatory Networking in LTE Networks," *IEEE Transactions on Mobile Computing*, vol. 17, no. 10, pp. 2252–2265, Oct 2018.

[355] Y. Xu, W. Xu, F. Yin, J. Lin, and S. Cui, "High-Accuracy Wireless Traffic Prediction: A GP-Based Machine Learning Approach," in *IEEE Global Communications Conference (GLOBECOM)*, Dec 2017, pp. 1–6.

[356] R. Sivakumar, E. Ashok Kumar, and G. Sivaradje, "Prediction of traffic load in wireless network using time series model," in *International Conference on Process Automation, Control and Computing*, July 2011, pp. 1–6.

[357] C. Qiu, Y. Zhang, Z. Feng, P. Zhang, and S. Cui, "Spatio-Temporal Wireless Traffic Prediction With Recurrent Neural Network," *IEEE Wireless Communications Letters*, vol. 7, no. 4, pp. 554–557, Aug 2018.

[358] S. Pandi, F. H. P. Fitzek, C. Lehmann, D. Nophut, D. Kiss, V. Kovacs, A. Nagy, G. Csorvasi, M. Toth, T. Rajacsis, H. Charaf, and R. Liebhart, "Joint Design of Communication and Control for Connected Cars in 5G Communication Systems," in *IEEE Globecom Workshops (GC Wkshps)*, Dec 2016, pp. 1–7.

[359] S. Jain, A. Kumar, S. Mandal, J. Ong, L. Poutievski, A. Singh, S. Venkata, J. Wanderer, J. Zhou, M. Zhu, J. Zolla, U. Hölzle, S. Stuart, and A. Vahdat, "B4: Experience with a Globally-deployed Software Defined WAN," in *Proceedings of the ACM SIGCOMM 2013 Conference on SIGCOMM*, ser. SIGCOMM '13. New York, NY, USA: ACM, 2013, pp. 3–14. [Online]. Available: http://doi.acm.org/10.1145/2486001.2486019

[360] L. Cui, F. R. Yu, and Q. Yan, "When big data meets software-defined networking: SDN for big data and big data for SDN," *IEEE Network*, vol. 30, no. 1, pp. 58–65, January 2016.

[361] E. Dahlman and S. Parkvall, "NR - The New 5G Radio-Access Technology," in *IEEE 87th Vehicular Technology Conference (VTC Spring)*, June 2018, pp. 1–6.

[362] Telecom Italia, Fondazione Bruno Kessler, "Open big data initiative," 2014. [Online]. Available: https://dandelion.eu/datamine/open-big-data/

[363] Z. Ali, N. Baldo, J. Mangues-Bafalluy, and L. Giupponi, "Machine learning based handover management for improved QoE in LTE," in *IEEE/IFIP Network Operations and Management Symposium (NOMS)*, April 2016, pp. 794–798.

[364] K. Poularakis, G. Iosifidis, G. Smaragdakis, and L. Tassiulas, "One step at a time: Optimizing SDN upgrades in ISP networks," in *IEEE Conference on Computer Communications (INFOCOM)*, May 2017, pp. 1–9.

[365] B. Heller, R. Sherwood, and N. McKeown, "The controller placement problem," in *Proceedings of the First Workshop on Hot Topics in Software Defined Networks*, ser. HotSDN '12. Helsinki, Finland: ACM, 2012, pp. 7–12.

[366] S. Deb, P. Monogioudis, J. Miernik, and J. P. Seymour, "Algorithms for Enhanced Inter-Cell Interference Coordination (eICIC) in LTE HetNets," *IEEE/ACM Transactions on Networking*, vol. 22, no. 1, pp. 137–150, Feb 2014.

[367] T. Zhang, A. Bianco, and P. Giaccone, "The role of inter-controller traffic in SDN controllers placement," in *IEEE Conference on Network Function Virtualization and Software Defined Networks (NFV-SDN)*, Nov 2016, pp. 87–92.

[368] S. E. Schaeffer, "Graph clustering," *Computer Science Review*, vol. 1, no. 1, pp. 27 – 64, 2007.

[369] M. C. Nascimento and A. C. de Carvalho, "Spectral methods for graph clustering – a survey," *European Journal of Operational Research*, vol. 211, no. 2, pp. 221 – 231, 2011.

[370] A. Blum, J. Hopcroft, and R. Kannan, "Foundations of data science," *Vorabversion eines Lehrbuchs*, 2016.

[371] F. D. Malliaros and M. Vazirgiannis, "Clustering and community detection in directed networks: A survey," *Physics Reports*, vol. 533, no. 4, pp. 95 – 142, 2013.

[372] U. Von Luxburg, "A tutorial on spectral clustering," *Statistics and computing*, vol. 17, no. 4, pp. 395–416, 2007.

[373] P. Bradley, K. Bennett, and A. Demiriz, "Constrained k-means clustering," *Microsoft Research, Redmond*, pp. 1–8, 2000.

[374] K. Thaalbi, M. T. Missaoui, and N. Tabbane, "Performance analysis of clustering algorithm in a C-RAN architecture," in *13th International Wireless Communications and Mobile Computing Conference (IWCMC)*, June 2017, pp. 1717–1722.

[375] D. Mishra, P. C. Amogh, A. Ramamurthy, A. A. Franklin, and B. R. Tamma, "Load-aware dynamic RRH assignment in Cloud Radio Access Networks," in *IEEE Wireless Communications and Networking Conference*, April 2016, pp. 1–6.

[376] O. G. Aliu, A. Imran, M. A. Imran, and B. Evans, "A survey of self organisation in future cellular networks," *IEEE Communications Surveys & Tutorials*, vol. 15, no. 1, pp. 336–361, 2013.

[377] D. J. MacKay, "Bayesian interpolation," *Neural computation*, vol. 4, no. 3, pp. 415–447, May 1992.

[378] Q. Shi, M. Abdel-Aty, and J. Lee, "A Bayesian ridge regression analysis of congestion's impact on urban expressway safety," *Accident Analysis & Prevention*, vol. 88, pp. 124–137, Mar. 2016.

[379] L. Breiman, "Random forests," *Machine learning*, vol. 45, no. 1, pp. 5–32, October 2001.

[380] R. W. Douglass, D. A. Meyer, M. Ram, D. Rideout, and D. Song, "High resolution population estimates from telecommunications data," *EPJ Data Science*, vol. 4, no. 1, p. 4, Dec. 2015.

[381] C. E. Rasmussen, "Gaussian processes in machine learning," in *Advanced lectures on machine learning*. Springer, 2004, pp. 63–71.

[382] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, October 2011.

[383] G. Wang, "Downlink Shared Channel Evaluation of LTE System," Master of Science Thesis in Communication Engineering, Chalmers University of Technology, Gothenburg, Sweden, 2013.

[384] 3GPP, "Evolved Universal Terrestrial Radio Access (E-UTRA); Physical layer procedures," TS 36.213 - V15.1.0, 2018.

[385] 3GPP, "Study on channel model for frequencies from 0.5 to 100 GHz," TR 38.901, V15.0.0, 2018.

[386] ——, "Study on latency reduction techniques for LTE," TR 36.881, V14.0.0, 2016.

[387] S. Sur, I. Pefkianakis, X. Zhang, and K.-H. Kim, "WiFi-Assisted 60 GHz Wireless Networks," in *Proc. of ACM MobiCom*, Snowbird, Utah, USA, 2017, pp. 28–41.

# List of Publications

## Journals

[388] M. Polese, M. Giordani, M. Mezzavilla, S. Rangan, and M. Zorzi, "Improved Handover Through Dual Connectivity in 5G mmWave Mobile Networks," *IEEE Journal on Selected Areas in Communications*, vol. 35, no. 9, pp. 2069–2084, September 2017.

[389] M. Polese, R. Jana, and M. Zorzi, "TCP and MP-TCP in 5G mmWave Networks," *IEEE Internet Computing*, vol. 21, no. 5, pp. 12–19, September 2017.

[390] M. Mezzavilla, M. Zhang, M. Polese, R. Ford, S. Dutta, S. Rangan, and M. Zorzi, "End-to-End Simulation of 5G mmWave Networks," *IEEE Communications Surveys & Tutorials*, vol. 20, no. 3, pp. 2237–2263, Third quarter 2018.

[391] M. Mezzavilla, M. Polese, A. Zanella, A. Dhananjay, S. Rangan, C. Kessler, T. S. Rappaport, and M. Zorzi, "Public Safety Communications above 6 GHz: Challenges and Opportunities," *IEEE Access*, vol. 6, pp. 316–329, 2018.

[392] M. Dalla Cia, F. Mason, D. Peron, F. Chiariotti, M. Polese, T. Mahmoodi, M. Zorzi, and A. Zanella, "Using Smart City Data in 5G Self-Organizing Networks," *IEEE Internet of Things Journal*, vol. 5, no. 2, pp. 645–654, April 2018.

[393] M. Zhang, M. Polese, M. Mezzavilla, J. Zhu, S. Rangan, S. Panwar, and a. M. Zorzi, "Will TCP Work in mmWave 5G Cellular Networks?" *IEEE Communications Magazine*, vol. 57, no. 1, pp. 65–71, January 2019.

[394] M. Giordani, M. Polese, A. Roy, D. Castor, and M. Zorzi, "Standalone and Non-Standalone Beam Management for 3GPP NR at mmWaves," *IEEE Communications Magazine*, vol. 57, no. 4, pp. 123–129, April 2019.

[395] ——, "A Tutorial on Beam Management for 3GPP NR at mmWave Frequencies," *IEEE Communications Surveys & Tutorials*, vol. 21, no. 1, pp. 173–196, First quarter 2019.

[396] M. Polese, F. Chiariotti, E. Bonetto, F. Rigotto, A. Zanella, and M. Zorzi, "A survey on recent advances in transport layer protocols," *IEEE Communications Surveys & Tutorials*, pp. 1–1, 2019.

[397] F. Meneghello, M. Calore, D. Zucchetto, M. Polese, and A. Zanella, "IoT: Internet of Threats? A survey of practical security vulnerabilities in real IoT devices," *IEEE Internet of Things Journal*, pp. 1–1, 2019.

[398] M. Polese, R. Jana, V. Kounev, K. Zhang, S. Deb, and M. Zorzi, "Machine Learning at the Edge: A Data-Driven Architecture with Applications to 5G Cellular Networks," *submitted to IEEE Transactions on Mobile Computing*, 2019. [Online]. Available: https://arxiv.org/pdf/1808.07647.pdf

[399] M. Polese, M. Giordani, T. Zugno, A. Roy, S. Goyal, D. Castor, and M. Zorzi, "Integrated Access and Backhaul in 5G mmWave Networks: Potentials and Challenges," *submitted to IEEE Communications Magazine*, 2019. [Online]. Available: https://arxiv.org/pdf/1906.01099.pdf

[400] M. Giordani, M. Polese, M. Mezzavilla, S. Rangan, and M. Zorzi, "Towards 6G Networks: Use Cases and Technologies," *submitted to IEEE Communications Magazine*, 2019. [Online]. Available: https://arxiv.org/pdf/1903.12216.pdf

# Conference Proceedings

[401] M. Polese, M. Centenaro, A. Zanella, and M. Zorzi, "M2M massive access in LTE: RACH performance evaluation in a Smart City scenario," in *IEEE International Conference on Communications (ICC)*, May 2016, pp. 1–6.

[402] M. Polese, M. Mezzavilla, and M. Zorzi, "Performance Comparison of Dual Connectivity and Hard Handover for LTE-5G Tight Integration," in *Proceedings of the 9th EAI International Conference on Simulation Tools and Techniques*, ser. SIMUTOOLS'16, Prague, Czech Republic, 2016, pp. 118–123.

[403] F. Chiariotti, D. D. Testa, M. Polese, A. Zanella, G. M. D. Nunzio, and M. Zorzi, "Learning methods for long-term channel gain prediction in wireless networks," in *International Conference on Computing, Networking and Communications (ICNC)*, Jan 2017, pp. 162–166.

[404] M. Polese, R. Jana, and M. Zorzi, "TCP in 5G mmWave Networks: Link Level Retransmissions and MP-TCP," in *IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS)*, May 2017.

[405] E. Lovisotto, E. Vianello, D. Cazzaro, M. Polese, F. Chiariotti, D. Zucchetto, A. Zanella, and M. Zorzi, "Cell Traffic Prediction Using Joint Spatio-Temporal Information," in *6th International Conference on Circuits and Systems Technologies (MOCAST)*, May 2017.

[406] M. Zhang, M. Polese, M. Mezzavilla, S. Rangan, and M. Zorzi, "ns-3 Implementation of the 3GPP MIMO Channel Model for Frequency Spectrum Above 6 GHz," in *Proceedings of the Workshop on ns-3*. Porto, Portugal: ACM, 2017, pp. 71–78. [Online]. Available: http://doi.acm.org/10.1145/3067665.3067678

[407] T. Azzino, M. Drago, M. Polese, A. Zanella, and M. Zorzi, "X-TCP: a cross layer approach for TCP uplink flows in mmwave networks," in *16th Annual Mediterranean Ad Hoc Networking Workshop (Med-Hoc-Net)*, June 2017.

[408] M. Dalla Cia, F. Mason, D. Peron, F. Chiariotti, M. Polese, T. Mahmoodi, M. Zorzi, and A. Zanella, "Mobility-aware Handover Strategies in Smart Cities," in *International Symposium on Wireless Communication Systems (ISWCS)*, August 2017.

[409] M. Polese, M. Mezzavilla, S. Rangan, and M. Zorzi, "Mobility Management for TCP in mmWave Networks," in *Proceedings of the 1st ACM Workshop on Millimeter-Wave Networks and Sensing Systems 2017*, ser. mmNets '17. Snowbird, Utah, USA: ACM, 2017, pp. 11–16.

[410] M. Gentil, A. Galeazzi, F. Chiariotti, M. Polese, A. Zanella, and M. Zorzi, "A deep neural network approach for customized prediction of mobile devices discharging time," in *IEEE Global Communications Conference (GLOBECOM)*, Dec 2017, pp. 1–6.

[411] M. Polese, M. Mezzavilla, M. Zhang, J. Zhu, S. Rangan, S. Panwar, and M. Zorzi, "milliProxy: A TCP proxy architecture for 5G mmWave cellular systems," in *51st Asilomar Conference on Signals, Systems, and Computers*, Oct 2017, pp. 951–957.

[412] M. Polese, M. Mezzavilla, S. Rangan, C. Kessler, and M. Zorzi, "mmwave for future public safety communications," in *Proceedings of the First CoNEXT Workshop on ICT Tools for Emergency Networks and DisastEr Relief*, ser. I-TENDER '17. Incheon, Republic of Korea: ACM, 2017, pp. 44–49. [Online]. Available: http://doi.acm.org/10.1145/3152896.3152905

[413] M. Drago, T. Azzino, M. Polese, C. Stefanovic, and M. Zorzi, "Reliable Video Streaming over mmWave with Multi Connectivity and Network Coding," in *International Conference on Computing, Networking and Communications (ICNC)*, March 2018, pp. 508–512.

[414] T. Zugno, M. Polese, and M. Zorzi, "Integration of Carrier Aggregation and Dual Connectivity for the ns-3 mmWave Module," in *Proceedings of the 10th Workshop on ns-3*, ser. WNS3 '18.  Surathkal, India: ACM, 2018, pp. 45–52. [Online]. Available: http://doi.acm.org/10.1145/3199902.3199909

[415] M. Polese and M. Zorzi, "Impact of Channel Models on the End-to-End Performance of Mmwave Cellular Networks," in *IEEE 19th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, June 2018, pp. 1–5.

[416] M. Giordani, M. Polese, A. Roy, D. Castor, and M. Zorzi, "Initial access frameworks for 3GPP NR at mmWave frequencies," in *17th Annual Mediterranean Ad Hoc Networking Workshop (Med-Hoc-Net)*, June 2018, pp. 1–8.

[417] M. Polese, M. Giordani, A. Roy, S. Goyal, D. Castor, and M. Zorzi, "End-to-End Simulation of Integrated Access and Backhaul at mmWaves," in *IEEE 23rd International Workshop on Computer Aided Modeling and Design of Communication Links and Networks (CAMAD)*, Sep. 2018, pp. 1–7.

[418] M. Polese, M. Giordani, A. Roy, D. Castor, and M. Zorzi, "Distributed Path Selection Strategies for Integrated Access and Backhaul at mmWaves," in *IEEE Global Communications Conference (GLOBECOM)*, Dec 2018.

[419] M. Rebato, M. Polese, and M. Zorzi, "Multi-Sector and Multi-Panel Performance in 5G mmWave Cellular Networks," in *IEEE Global Communications Conference (GLOBECOM)*, Dec 2018, pp. 1–6.

[420] M. Polese, R. Jana, V. Kounev, K. Zhang, S. Deb, and M. Zorzi, "Exploiting spatial correlation for improved user prediction in 5G cellular networks," in *Proceedings of the Information Theory and Applications Workshop*, ser. ITA '19, San Diego, 2019.

[421] W. Xia, M. Polese, M. Mezzavilla, G. Loianno, S. Rangan, and M. Zorzi, "Millimeter Wave Remote UAV Control and Communications for Public Safety Scenarios," in *Proceedings of the 1st International Workshop on Internet of Autonomous Unmanned Vehicles*, ser. IAUV '19, Boston, MA, 2019.

[422] M. Polese, T. Zugno, and M. Zorzi, "Implementation of Reference Public Safety Scenarios in ns-3," in *Proceedings of the 2019 Workshop on ns-3*, ser. WNS3 2019.  Florence, Italy: ACM, 2019, pp. 73–80. [Online]. Available: http://doi.acm.org/10.1145/3321349.3321356

[423] A. De Biasio, F. Chiariotti, M. Polese, A. Zanella, and M. Zorzi, "A QUIC Implementation for ns-3," in *Proceedings of the Workshop on ns-3*, ser. WNS3 2019.  Florence, Italy: ACM, 2019, pp. 1–8. [Online]. Available: http://doi.acm.org/10.1145/3321349.3321351

[424] T. Zugno, M. Polese, M. Lecci, and M. Zorzi, "Simulation of Next-generation Cellular Networks with ns-3: Open Challenges and New Directions," in *Proceedings of the 2019 Workshop on Next-Generation Wireless with ns-3*, ser. WNGW 2019.  Florence, Italy: ACM, 2019, pp. 38–41. [Online]. Available: http://doi.acm.org/10.1145/3337941.3337951

[425] M. Polese, F. Restuccia, A. Gosain, J. Jornet, S. Bhardwaj, V. Ariyarathna, S. Mandal, K. Zheng, A. Dhananjay, M. Mezzavilla, J. Buckwalter, M. Rodwell, X. Wang, M. Zorzi, A. Madanayake, and T. Melodia, "MillimeTera: Toward A Large-Scale Open-Source mmWave and Terahertz Experimental Testbed," in *Proceedings of the 3rd ACM Workshop on Millimeter-Wave Networks and Sensing Systems*, ser. mmNets '19.  Los Cabos, Mexico: ACM, 2019.

[426] L. Bertizzolo, M. Polese, L. Bonati, A. Gosain, M. Zorzi, and T. Melodia, "mmBAC: Location-aided mmWave Backhaul Management for UAV-based Aerial Cells," in *Proceedings of the 3rd ACM Workshop on Millimeter-Wave Networks and Sensing Systems*, ser. mmNets '19.  Los Cabos, Mexico: ACM, 2019.

[427] M. Drago, M. Polese, S. Kucera, D. Kozlov, V. Kirillov, and M. Zorzi, "QoS Provisioning in 60 GHz Communications by Physical and Transport Layer Coordination," in *IEEE 16th International Conference on Mobile Ad Hoc and Sensor Systems (MASS)*, Nov 2019.

# Book Chapters

[428] M. Polese, M. Giordani, and M. Zorzi, "3GPP NR: the standard for 5G cellular networks," in *5G Italy White eBook: from Research to Market*, 2018.

# Acknowledgments

The three years I spent carrying out the research of my Ph.D. have been some of the most enjoyable so far. For this, I would like to thank all of the people that have been part of this journey.

First, this would not have been possible without my advisor, Michele Zorzi, who convinced me to apply to the program, and allowed me to grow as a researcher, while connecting to a worldwide network with the best minds working on wireless communications. I would also like to thank Andrea Zanella, who trusted me as tutor for his class for three years in a row, giving me the opportunity of working with very talented studens on exciting projects. A big thanks also goes to Michele Rossi and Lorenzo Vangelista, for all the stimulating discussions. I really enjoyed being part of the SIGNET Lab, and, in particular, I would like to thank Marco Giordani, Federico Chiariotti, and Mattia Rebato, with whom I shared many projects, papers and discussions. Thanks to all my other coauthors from Padova, who gave me precious advice, or trusted me with their projects and theses (thank you!): Tommaso Zugno, Matteo Drago, Mattia Lecci, Paolo Testolina, Daniel Zucchetto, Marco Centenaro, Davide Del Testa, Giorgio Maria Di Nunzio, Tommy Azzino, Federico Mason, Davide Peron, Massimo Dalla Cia, Francesca Meneghello, Matteo Calore, Enrico Vianello, Enrico Lovisotto, Davide Cazzaro, Mattia Gentil, Alessandro Galezzi, Elia Bonetto, Filippo Rigotto, and Alvise De Blasio.

The people I have worked with around the world have made this an unique experience. I would like to thank my coauthors Sundeep Rangan, Marco Mezzavilla, Menglei Zhang, Sourjya Dutta, Russell Ford, Aditya Dhananjay, William Xia, Giuseppe Loianno, Shivendra Panwar, and Theodore Rappaport from NYU; Jing Zhu from Intel; Rittwik Jana, Velin Kounev, Supratim Deb, and Ke Zhang from AT&T; Tommaso Melodia, Leonardo Bonati, Lorenzo Bertizzolo, Francesco Restuccia, Salvatore D'Oro, Abimayu Gosain, and Josep Jornet from Northeastern; Arnab Roy, Sanjay Goyal, Douglas Castor from InterDigital; Stepan Kucera and Dmitry Kozlov from Nokia Bell Labs; Toktam Mahmoodi from King's College; and Cedomir Stefanovic from Aalborg University. I also would like to thank Mihaela Beluri, Tom Henderson, Nada Golmie, Natale Patriciello, Chris Slezak, and Emrecan Demiros for all the interesting exchanges.

On the personal side, it was a pleasure to share this experience with friends that I met years ago in the Department's classrooms and with whom I hope to keep in touch in the future, even if we will be spread all around the world. Silvia, Gian, Davide, Paolo, Laura, Alberto and Andrea (even if remote), thanks for being part of this! Thanks to all the new and old friends from Padova, Oderzo, New York and Boston: Davide, Erin, Chiara (also - thanks for hosting me when I needed it), Umberto, Federico, Martina, Mattia, Paolo, Federico, Francesco, Leonardo, Matteo, Marco, Tommaso, Riccardo, Sebastiano, Daniel, Filippo, Giovanna, Marco, Sara, Pasquale, Leonardo, Lorenzo, Sila, Daniel, Deniz, Volkan, Bernard, and Velin.

Last but not least, special thanks go to my parents, Diego and Stella, who always supported me and let me freely pursue this path (grazie!), and to Giulia, one of the most precious findings of the last part of this Ph.D., across Padova, Boston, and Brisbane.