

MATTEO PAGIN

MODELS, ALGORITHMS AND PROTOCOLS FOR INNOVATIVE DEPLOYMENT SOLUTIONS IN 5G AND 6G CELLULAR NETWORKS

A THESIS SUBMITTED IN PARTIAL FULFILLMENT OF THE REQUIREMENTS FOR THE
PH.D. DEGREE IN INFORMATION ENGINEERING

OF THE
DEPARTMENT OF INFORMATION ENGINEERING
UNIVERSITY OF PADOVA
SUPERVISED BY
PROF. MICHELE ZORZI

XXXVII CYCLE
ACADEMIC YEAR 2023/2024

Abstract

Cellular networks are constantly evolving in order to support the ever-increasing number of mobile users, and the corresponding growth in wireless data traffic, coupled with the emergence of new applications. Specifically, the last generation of mobile networks, i.e., 5G, brought high peak performance and extreme flexibility, making it possible to support a diverse set of applications with heterogeneous yet stringent requirements. One of 5G's main novelties is represented by the support for millimeter wave (mmWave) frequencies, which unlocks an unprecedented amount of previously unused radio resources. In turn, the latter enables extremely high data rates and low latencies. Moreover, it is envisioned that the upcoming generation, i.e., 6G, will unleash additional bandwidth, by further expanding the supported spectrum bands to include terahertz (THz) frequencies as well. However, despite their theoretical potential, mmWave and THz frequencies exhibit harsh propagation conditions which make it challenging to provide ubiquitous high speed wireless connectivity. To fill this gap, this thesis studies innovative deployment solutions to overcome the unfavorable propagation characteristics of mmWave and THz communications, paving the way for their widespread use in the context of 6G cellular networks.

In particular, this thesis (i) presents novel simulation tools, which model innovative coverage enhancement technologies such as Intelligent Reflective Surfaces (IRSs), Amplify-and-Forward (AF) relays, and Non-Terrestrial Networkss (NTNs); (ii) presents novel simulation models which improve the computational complexity of Multiple Input, Multiple Output (MIMO) simulations; (iii) introduces schemes for optimizing Integrated Access and Back-haul (IAB) networks; (iv) analyzes the potential of mixed mmWave and THz links for wireless backhauling; and (v) analyzes the impact of non-ideal control channels in IRS-aided deployments, and introduces algorithms for mitigating the corresponding performance degradation.

This thesis adopts a system-level approach, thus characterizing the network behavior in an end-to-end fashion, and capturing the interplay between the physical signal propagation and the different layers of the communications protocol stack. Results demonstrate the effectiveness of the proposed

solutions, which pave the way towards ubiquitous high-performance mobile networks.

Sommario

Le reti cellulari sono in costante evoluzione, sia per sostenere il numero crescente di utenti mobili e la corrispondente crescita del traffico di dati, sia per far fronte all'emergenza di nuovi scenari d'uso. In particolare, l'ultima versione delle reti mobili, ovvero 5G, offre alte prestazioni di picco ed estrema flessibilità, permettendo di sostenere un diverso insieme di applicazioni con caratteristiche eterogenee. Una delle principali novità della tecnologia 5G è rappresentata dal supporto alle frequenze mmWave, che forniscono l'accesso ad una quantità inedita di risorse radio precedentemente non utilizzate. Queste ultime consentono trasmissioni dati a velocità estremamente elevate e latenze particolarmente basse. Inoltre, si prevede che la prossima generazione di reti cellulari, ovvero 6G, supporterà ulteriori bande includendo anche le frequenze THz.

Tuttavia, nonostante il loro potenziale teorico, le frequenze mmWave e THz presentano condizioni di propagazione estremamente sfavorevoli, che rendono difficile fornire una connessione wireless veloce ed onnipresente. A questo fine, questa tesi studia infrastrutture innovative per ovviare all'intrinseca copertura limitata delle comunicazioni mmWave e THz, aprendo la strada ad un loro utilizzo diffuso nel contesto delle reti cellulari 6G.

In particolare, questa tesi (i) presenta nuovi strumenti di simulazione, che modellano tecnologie di miglioramento della copertura innovative come IRSs, ripetitori AF e NTNs; (ii) presenta nuovi modelli di simulazione che migliorano la complessità computazionale delle simulazioni MIMO; (iii) introduce schemi di ottimizzazione per reti IAB; (iv) analizza il potenziale di collegamenti misti mmWave e THz per wireless backhauling; (v) analizza l'impatto di canali di controllo non ideali in reti supportate da IRSs e introduce algoritmi per mitigare la conseguente perdita di prestazioni.

Questa tesi adotta un approccio di studio a livello di sistema, quindi caratterizzando il comportamento della rete nel suo complesso e catturando l'interazione tra la propagazione del segnale fisico ed i diversi livelli dello stack protocollare. I risultati presentati dimostrano l'efficacia delle soluzioni proposte, che aprono quindi la strada verso reti cellulari ubique con alte prestazioni.

Contents

1	<i>Introduction</i>	1
2	<i>Simulation tools for future cellular networks</i>	3
2.1	<i>Modeling AF and IRS relays-aided wireless channels</i>	7
2.1.1	<i>The TR 38.901 Channel Model for 5G NR</i>	8
2.1.2	<i>A Signal Model for the IRS</i>	10
2.1.3	<i>A Signal Model for the AF Relay</i>	11
2.1.4	<i>A Full-Stack Simulator for IRS/AF Relays</i>	11
2.1.4.1	Implementation of the IRS/AF Signal Model	12
2.1.4.2	Integration of the IRS/AF Signal Model in the Simulator	16
2.1.5	<i>Performance Evaluation</i>	17
2.1.5.1	Simulation Setup	18
2.1.6	<i>Numerical Results</i>	19
2.2	<i>Modeling non-terrestrial wireless channels</i>	25
2.2.1	<i>Scenarios and Path Loss Condition</i>	26
2.2.2	<i>Path Loss</i>	26
2.2.3	<i>Atmospheric Absorption</i>	27
2.2.4	<i>Scintillation</i>	27
2.2.4.1	Ionospheric Scintillation	28
2.2.4.2	Tropospheric Scintillation	28
2.2.5	<i>Fast Fading</i>	28
2.2.6	<i>Antenna Model</i>	29
2.2.7	<i>Coordinate System</i>	29
2.2.8	<i>Implementation in ns-3</i>	30
2.2.8.1	Small-scale fading	31
2.2.8.2	Coordinate systems	31
2.2.8.3	Channel condition	32

2.2.8.4	Path loss and shadowing	32
2.2.8.5	Geocentric mobility models	33
2.2.8.6	Antenna models	33
2.2.9	<i>Examples and Comparisons</i>	34
2.2.9.1	Link-Level Results	34
2.2.9.2	Frequency Test	36
2.2.9.3	Mobility Test	37
2.2.9.4	End-to-End Performance	38
2.3	<i>Improving the scalability of wireless channel simulation in ns-3</i>	40
2.3.1	<i>Efficient MIMO modeling with the Eigen library</i>	41
2.3.2	<i>A performance-oriented MIMO statistical channel model</i>	43
2.3.3	<i>Path loss, Shadowing, and LoS Condition</i>	44
2.3.3.1	Antenna and Beamforming Gain	45
2.3.3.2	Fast Fading	47
2.3.4	<i>Benchmarks, examples and use cases</i>	50
2.3.4.1	Examples and Benchmarks	50
2.3.4.2	Use Cases	53
2.4	<i>Conclusions and future work</i>	53
3	<i>Towards wireless-backhauled next-generation cellular networks</i>	55
3.1	<i>Semi-centralized framework for resource management in 5G NR Integrated Access and Backhaul</i>	58
3.1.1	<i>Contributions</i>	60
3.1.2	<i>IAB networks</i>	61
3.1.2.1	Network topology	61
3.1.2.2	Multiple access schemes and scheduling	63
3.1.2.3	System model	63
3.1.3	<i>Semi-centralized resource allocation scheme for IAB networks</i>	65
3.1.3.1	MWM for ST graphs	65
3.1.4	<i>Semi-centralized resource partitioning scheme</i>	70
3.1.5	<i>Implementation of semi-centralized allocation schemes in mmWave IAB networks</i>	71
3.1.5.1	Simulation scenario and parameters	79
3.1.6	<i>Performance evaluation</i>	81
3.1.6.1	Throughput	81

Contents

3.1.6.2	Latency	84
3.1.6.3	Network congestion	86
3.1.6.4	Performance with TCP traffic	86
3.1.6.5	Further considerations	88
3.2	<i>Risk-averse learning for latency minimization in mmWave integrated access and backhaul networks</i>	91
3.2.1	<i>System Model</i>	92
3.2.2	<i>Problem Formulation</i>	94
3.2.2.1	Preliminaries on Conditional Value at Risk (CVaR)	95
3.2.2.2	Optimization Problem	95
3.2.3	<i>Our proposed solution: Safehaul</i>	96
3.2.3.1	Multi-Armed Bandit Formulation	97
3.2.3.2	Latency and CVaR Estimation	98
3.2.3.3	Consensus	98
3.2.3.4	Implementation of Safehaul	99
3.2.3.5	Regret Analysis	100
3.2.4	<i>Simulation setup</i>	102
3.2.4.1	Extensions to Sionna’s physical layer module .	103
3.2.5	<i>Performance Evaluation</i>	106
3.2.5.1	Scenario 1: Average Network Performance .	109
3.2.5.2	Scenario 2: Impact of the Network Size	110
3.2.5.3	Scenario 3: Impact of the number of BS-donors	111
3.2.5.4	Scenario 4: Impact of the risk parameter α . .	112
3.2.5.5	Scenario 5: Performance in different topologies	113
3.2.5.6	Scenario 6: Network resilience	116
3.2.6	<i>Related work</i>	118
3.3	<i>High-capacity integrated access and backhaul networks using sub-terahertz links</i>	124
3.3.1	<i>System Model</i>	125
3.3.1.1	Channel Models	126
3.3.2	<i>Sum-rate optimization via THz Link Selection</i>	128
3.3.2.1	Backhaul Scheduler	128
3.3.3	<i>Performance Evaluation</i>	130
3.3.3.1	Simulation Setup	130
3.3.3.2	Simulation Scenario	131
3.3.3.3	Results	132

3.4 Conclusions and future work	137
4 Downlink Clustering-Based Scheduling of IRS-Assisted Communications With Reconfiguration Constraints	139
4.1 Prior Work	139
4.2 Contributions	141
4.2.1 Organization and Notation	143
4.3 System Model	143
4.3.1 IRS Model	145
4.4 Sum Capacity Optimization Problem	146
4.5 Heuristic Sum Capacity Maximization	148
4.5.1 Optimal Individual IRS Configurations	148
4.5.2 Clustering-based TDMA Scheduling	150
4.6 Distance-Based Clustering Algorithms	151
4.7 Capacity-Based Clustering Algorithms	153
4.7.1 Capacity-Weighted Clustering (CWC)	153
4.7.2 One-Shot Capacity-Based Clustering (OSCBC)	156
4.7.3 Inverse Capacity-Weighted Clustering (ICWC)	157
4.7.4 Computational Complexity	157
4.8 Numerical Results	159
4.8.1 Simulation Parameters	159
4.8.2 Performance Metrics	162
4.8.3 Scheduling Performance	163
4.9 Conclusions	171
5 Conclusions	173
<i>Acronyms</i>	175
<i>Publications</i>	181
<i>Bibliography</i>	185

1 *Introduction*

Mobile networks play a key role in our society and are poised to become ever more important in the coming years. In fact, the International Telecommunications Union (ITU) foresees that in 2030 and beyond wireless broadband will be ubiquitous, and will be required to provide connectivity not only to humans, but also to a plethora of intelligent devices such as wearables, road vehicles, Unmanned Aerial Systems (UASs) and robots [1]. Moreover, novel use cases such as holographic communications, Extended Reality (XR) and tactile applications will further exacerbate the throughput and latency requirements which were posed by enhanced Mobile Broadband (eMBB) and Ultra-Reliable Low-Latency Communication (URLLC) [2].

To meet these goals, future cellular systems will further evolve 5th generation (5G) networks, which have introduced a flexible, virtualized architecture, the support for mmWave communications and the use of massive MIMO (m-MIMO) technologies [3]. Notably, the research community is considering a more central role for mmWaves, a further expansion of the spectrum towards the THz band, and an Artificial Intelligence (AI)-native network design, with the goal of achieving autonomous data-centric orchestration and management of the network [4], possibly down to the air interface [5].

The THz and mmWave bands offer large chunks of untapped bandwidth which operators can leverage to meet the Tb/s peak rates that are envisioned by the ITU [1]. However, this portion of the spectrum is plagued by unfavorable propagation characteristics, comprising a marked free-space propagation loss and susceptibility to blockages [6, 7], which make it challenging to harvest its potential. Although the harsh propagation environment can be partially mitigated by using directional links and densifying network deployments [8], the support for mmWave and THz bands entails a major redesign not only of the physical layer, but of the whole cellular protocol stack [9]. For instance, the intrinsic directionality of the communication requires ad hoc control procedures [10], while the frequent transitions between Line-of-

1 Introduction

Sight (LoS) and Non-Line-of-Sight (NLoS) conditions call for an ad hoc transport layer design, such as novel Transmission Control Protocol (TCP) algorithms [11]. In addition, as the network progressively becomes increasingly complex and heterogeneous, the push for spectrum expansion will be coupled with an AI-native design which, thanks to the ongoing virtualization, will not be limited to the radio link level, but will encompass the orchestration of large scale deployments as well [12].

Nevertheless, how to design, test and eventually deploy management and orchestration algorithms is an open research challenge [13]. First, the training data must accurately capture the interplay of the whole protocol stack with the wireless channel. Furthermore, optimization frameworks such as Deep Reinforcement Learning (DRL) also call for preliminary testing in isolated yet realistic environments, with the goal of minimizing the performance degradation to actual network deployments [14, 15].

2 *Simulation tools for future cellular networks*

The preliminary, yet reliable performance evaluation of disruptive research ideas is a vital step in identifying the most promising key technology enablers for future cellular networks, before further research and development, and eventual deployment in real-world networks. To this end, while providing reliable performance estimates, experiments with real testbeds are often impractical due to limitations in the scalability and flexibility of the involved platforms, as well as the high cost of hardware components. At the same time, analytical models usually introduce strong assumptions for the sake of tractability. End-to-end simulations fill these gaps, by coupling (possibly simplified) analytical models with the protocol stack accuracy exhibited by testbeds and/or digital twins.

Additionally, the ever increasing complexity and heterogeneous nature of wireless networks is poised to be coupled with an AI-native design which, thanks to the ongoing virtualization, will not be limited to the radio link level, but will encompass the orchestration of large scale deployments as well [12]. Nevertheless, how to design, test and eventually deploy management and orchestration algorithms is an open research challenge [13]. First, the training data must accurately capture the interplay of the whole protocol stack with the wireless channel. Secondly, optimization frameworks such as DRL also call for preliminary testing in isolated yet realistic environments, with the goal of minimizing the performance degradation to actual network deployments [14, 15]. Both these requirements are met by end-to-end simulations, which will thus play a fundamental role in designing and evaluating the performance of the next generation of cellular networks.

Recently, both academic and industry researchers have been strongly favoring open-source simulators [16], i.e., simulators which are made freely available for both use, modification and redistribution. In general, these char-

2 Simulation tools for future cellular networks

acteristics of these softwares foster decentralized development and open collaboration. In this class of simulators, ns-3 is the facto standard in the wireless research space, thanks to the already available modules for 5G NR [17, 18], IEEE 802.11ad/ay/ax [19–21] and its implementation of the 3GPP TR 38.901 statistical channel model [22]. Nevertheless, ns-3 currently lacks physical propagation models for most disruptive 6G deployment solutions. Most notably, ns-3 lacks channel models for AF and IRSs relays (also referred to as “smart relays”), and/or NTN. Both these technologies are expected to play a key role in achieving the ubiquitous connectivity target set for 6G networks. Moreover, the ns-3 TR 38.901 [23] channel modeling framework exhibits limitation in its scalability, thus rendering infeasible the simulation of large-scale deployments and/or terminals featuring massive-MIMO arrays.

Despite their widespread use, the short- and medium-term suitability of end-to-end network simulators as performance evaluation tools will largely depend on their scalability for realistically-sized deployments, and on the accuracy of their channel model [24]. In fact, system-level simulators generally abstract the actual link-level transmission via an error model, which maps the Signal-to-Interference-plus-Noise Ratio (SINR) of the wireless link to a packet error probability [25]. Eventually, the latter is used to determine whether the packet has been successfully decoded by the receiver. As a consequence, the accuracy of network simulators heavily depends on the reliability of the SINR estimation, especially when considering the mmWave and THz bands. Indeed, these portions of the spectrum entail a major redesign not only of the physical layer, but also of the whole cellular protocol stack [9], which makes it paramount to accurately model the peculiar propagation characteristics which they exhibit. For instance, the intrinsic directionality of the communication requires ad hoc control procedures [10], while the frequent transitions between LoS and NLoS conditions call for an ad hoc transport layer design, such as novel TCP algorithms [11].

To fill these gaps, in the first part of this chapter we introduce the design and implementation of channel models for innovative deployment solutions in ns-3. In particular, Section 2.1 describes a new channel model for ns-3 for IRS/AF-aided communications which is based on the current 3GPP channel model for 5G networks standardized in [23]. Then, we present how the former can be used in conjunction with the ns3-mmwave module [17], which models the Physical (PHY) and Medium Access Control (MAC) layers of the

5G NR protocol stack to achieve an end-to-end simulation framework for smart relays. The latter incorporates the interplay with the 5G NR protocol stack and relative control tasks, as well as the impact of the upper (including transport and application) layers. Then, we leverage this novel framework to conduct an extensive simulation campaign to study the performance of IRS/AF nodes for relaying connectivity requests from end users, compared to a baseline solution in which relays are not deployed. We demonstrate that IRSs and AF relays are valid solutions, especially in small networks, even though high-EIRP AF relays are required to support more aggressive traffic applications. Based on our simulations, we provide guidelines towards the optimal dimensioning of IRS and AF configurations, in terms of number of antenna elements and amplification power.

Section 2.2 presents a new open-source module for ns-3 that implements the NTN channel model based on the 3GPP specifications described in TR 38.811 [26]. While the described implementation is mainly related to the channel and the physical layer, a deep understanding of the propagation model is the first step towards proper protocol design [27], which makes our module a valuable and accurate tool in the study of NTNs. Specifically, the module introduces: (i) new simulation scenarios for NTN; (ii) a new Path Loss (PL) model for the air/space channel in a wide range of frequencies (from 0.5 GHz to 100 GHz); (iii) the characterization of atmospheric absorptions; (iv) a new fast fading model for the space environment; (v) an antenna model for both terrestrial and non-terrestrial nodes; and (vi) a new coordinate system to account for the Geocentric Cartesian coordinate system of satellites.

The second part of this chapter focuses on novel solutions to improve the scalability of the ns-3 5G NR simulation framework. Specifically, Section 2.3 presents optimizations to the ns-3 implementation of the TR 38.901 channel model of [28], both at the codebase and at the design level, which aim to provide wireless researchers with the tools for simulating future dense wireless scenarios in a computationally efficient manner. Specifically, we significantly improve the runtime of simulations involving the 3GPP TR 38.901 channel model [29] by porting the intensive linear algebra operations to the open-source library Eigen [30]. To this end, we also design and implement a set of common linear algebra APIs, which increase the modularity of the spectrum module with respect to the underlying data structures and algorithms. Ad-

2 Simulation tools for future cellular networks

ditionally, we propose a simplified channel model, based on [29], which aims to provide an additional order of magnitude of runtime reduction, at the cost of a slight accuracy penalty. Profiling results show that the support for Eigen, coupled with further TR 38.901 optimizations, leads to a decrease of up to 5 times in the simulation time of typical MIMO scenarios. Furthermore, the proposed performance-oriented channel model further improved the runtime of simulations, which now take as low as 6 % with respect to the full TR 38.901 channel model, with a negligible loss in accuracy.

2.1 Modeling AF and IRS relays-aided wireless channels

Disruptive technologies, such as IRSs and AF relays, have been proposed as promising alternatives to overcome the coverage issues of mmWave networks with energy efficiency in mind [31]. The former are meta-surfaces that can be programmed to favorably alter an Electromagnetic (EM) field towards an intended destination. Most notably, IRSs are nodes which passively beamform the impinging signal, without amplification, thus being able to guarantee minimum capacity requirements in dead spots with lower power consumption compared to IAB [32]. AF relays, instead, are envisioned to capture an incident electromagnetic wave coming from a base station, to actively amplify the received signal, and to re-radiate it towards a target area to be served. They are candidates for achieving higher capacity with respect to IRS nodes, at the expense of higher cost and amplification noise [33].

Despite their research hype, whether these technologies will be able to fulfill 5G (and beyond) service requirements and, if so, how to properly dimension IRS/AF systems, are still crucial issues that remain unsolved. While field experiments with real hardware are infeasible due to scalability and flexibility concerns, as well as the high cost of testbed components, computer-based simulations represent a viable approach for testing and calibrating IRS/AF deployments. Prior works, e.g., [34, 35], have addressed this task, though focusing on link-level analyses, which typically adopt conservative assumptions on the system architecture, and should be taken as a lower bound for more representative end-to-end performance studies. To fill this gap, in this section we provide a more comprehensive system-level performance evaluation of IRS/AF deployments using a new simulation framework that operates end-to-end, thus incorporating the interplay with the 5G NR protocol stack and relative control tasks, as well as the impact of the upper (including transport and application) layers.

To this end, we first provide a mathematical model for the IRS and AF relay channels (Secs. 2.1.2 and 2.1.3, respectively), based on the standard 3GPP channel model for 5G networks (Section 2.1.1).

Notation. We use boldface upper- and lower-case letters to refer to matrices and vectors, respectively, while lower-case letters denote scalars. We use \mathbf{I}_N to denote the identity matrix of order N , $[\Phi]_{j,k}$ to indicate the (j,k) -th entry of matrix Φ , $\text{diag}(\phi_1, \dots, \phi_N)$ to indicate an $N \times N$ diagonal matrix with entries

2 Simulation tools for future cellular networks

$\{\phi_j | j = 1, \dots, N\}$. We use the superscripts T, H and * for transposition, Hermitian transposition, and conjugation, respectively.

2.1.1 The TR 38.901 Channel Model for 5G NR

Throughout this work, we consider the 3GPP TR 38.901 Spatial Channel Model (SCM) standardized in [23], leaving the study of more accurate channel models (e.g., based on ray tracing measurements) as part of our future research. This choice is motivated by the fact that TR 38.901 supports a wide range of frequencies, from 0.5 to 100 GHz, and can be integrated with realistic beamforming models. Furthermore, it is suggested and adopted by the 3rd Generation Partnership Project (3GPP) itself for the performance evaluation of 5G networks via system-level simulations.

In particular, the TR 38.901 model outlines the procedures for generating a channel matrix \mathbf{H} whose entries $\mathbf{H}_{p,q}(t, \tau)$ correspond to the impulse response of the channel between the p -th radiating element of the antenna array of the signal source (S), and the q -th radiating element of the antenna array of its destination (D), at time t and with delay τ . To model multipath fading, each of these terms is computed as the superposition of N different clusters, each of which consists of M rays that arrive (depart) to (from) the antenna arrays with specific angles and powers. Based on [23], and using the simplifications proposed in [22], the generic entry $\mathbf{H}_{p,q}(t, \tau)$ of the channel matrix can then be computed as:

$$\begin{aligned} \mathbf{H}_{p,q}(t, \tau) &= \sum_{n=1}^N \sqrt{\frac{P_n}{M}} \sum_{m=1}^M \bar{\mathbf{F}}_{rx}(\theta_{n,m}^A, \phi_{n,m}^A) \\ &\quad \times \begin{bmatrix} e^{j\Phi_{n,m}^{\theta,\theta}} & \sqrt{K_{n,m}^{-1}} e^{j\Phi_{n,m}^{\theta,\phi}} \\ \sqrt{K_{n,m}^{-1}} e^{j\Phi_{n,m}^{\phi,\theta}} & e^{j\Phi_{n,m}^{\phi,\phi}} \end{bmatrix} \bar{\mathbf{F}}_{tx}(\theta_{n,m}^D, \phi_{n,m}^D) \quad (2.1) \\ &\quad \times e^{j\bar{\mathbf{k}}_{rx,n,m}^T \bar{\mathbf{d}}_{rx,p}} e^{j\bar{\mathbf{k}}_{tx,n,m}^T \bar{\mathbf{d}}_{tx,q}} e^{j2\pi v_n t} \delta(\tau - \tau_n). \end{aligned}$$

For a complete description of the specific terms appearing in Eq. (2.1) we refer the interested reader to [22].

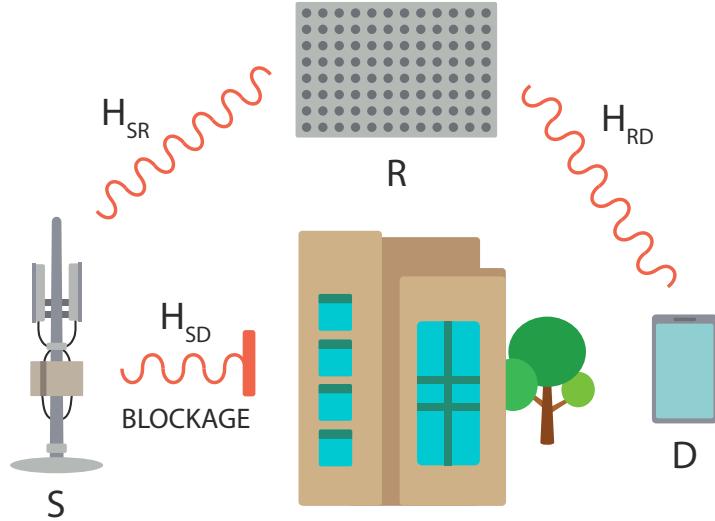


Figure 2.1: A typical urban scenario where a relay (R) can be used to bridge the signal from a source (S) to a destination (D), that would otherwise communicate in NLoS, i.e., the direct link between S and D is blocked due to obstacles such as buildings and/or vegetation.

Then, a frequency-flat path gain term is added to each channel coefficient as a function of the carrier frequency f_c and the distance d between the endpoints, i.e.,

$$\text{PL}(d, f_c) = A \log_{10}(d) + B + C \log_{10}(f_c) + X \text{ [dB]}, \quad (2.2)$$

where model parameters A, B and C depend on the propagation conditions and the type of environment, and X is an optional term to represent shadowing [22].

We consider the transmission of a single data stream x_S , i.e., a sequence of signals, from a source S to a destination D via a relay R, as depicted in Figure 2.1. Then, the channel matrix is combined with the beamforming vectors used at S and D, in order to obtain the SINR experienced at D. In particular, let x_S be the signal transmitted from S to D, and w_S, w_D and w_I be the beamforming vectors used at S, D and the I-th interferer, respectively. Moreover, we define the following matrices: \mathbf{H}_{SD} is the channel matrix between the source and the destination, \mathbf{H}_{ID} is the channel matrix between the I-th interferer and the destination, \mathbf{H}_{IR} is the channel matrix from the I-th interferer to the relay, \mathbf{H}_{SR} is the channel matrix between the source and the

2 Simulation tools for future cellular networks

relay, and \mathbf{H}_{RD} is the channel matrix between the relay and the destination. In a relay-free environment, the signal received at the User Equipment (UE) is computed as:

$$y_D = \mathbf{w}_D^T \mathbf{H}_{SD} \mathbf{w}_S x_S + \sum_{I=1}^N \mathbf{w}_D^T \mathbf{H}_{ID} \mathbf{w}_I x_I + \mathbf{w}_D^T \mathbf{n}_D. \quad (2.3)$$

where \mathbf{n}_D represents the circularly symmetric complex Gaussian noise vector with correlation matrix $\sigma_N^2 \mathbf{I}$, and $\mathbf{w}_D^T \mathbf{H}_{ID} \mathbf{w}_I x_I$ is the signal received from the I-th interferer. Accordingly, the SINR at D reads:

$$\Lambda = \frac{\|\mathbf{w}_D^T \mathbf{H}_{SD} \mathbf{w}_S\|^2 \sigma_S^2}{\sum_{I=1}^N \|\mathbf{w}_D^T \mathbf{H}_{ID} \mathbf{w}_I\|^2 \sigma_I^2 + \sigma_N^2}, \quad (2.4)$$

where σ_S^2 and σ_I^2 are the powers of the intended and the I-th interfering signals, respectively.

2.1.2 A Signal Model for the IRS

An IRS is a planar surface made of N_R low-cost passive reflecting elements that can be programmed to alter an EM field, for example to achieve three-dimensional beamforming towards an intended destination. The working principle is similar to that of a conventional relay, the main difference being that while the latter amplifies the received signal before retransmitting it, an IRS reflects and beamforms the signal without introducing any amplification, thus saving power compared with other relaying solutions [32].

In particular, each element of the IRS acts as an antenna that captures and reflects the incoming signals, introducing a phase shift on the baseband-equivalent signal. We denote with $\phi_n = e^{j\theta_n}$, $n = 1, \dots, N_R$, the reflection coefficient of the n -th IRS element, where $\theta_n \in [-\pi, \pi]$ is the induced, controllable phase shift. Adopting a complex baseband notation, the signal $z \in \mathbb{C}^{N_R \times 1}$ reflected by an IRS (denoted as R), impinged with a signal x_S originating from a source S, reads

$$z = \Phi \mathbf{H}_{SR} \mathbf{w}_S x_S, \quad (2.5)$$

where Φ is a diagonal matrix defined as $\Phi \doteq \text{diag}(\phi_1, \dots, \phi_{N_R})$, and typically referred to as *IRS configuration*. Therefore, the signal received at the intended

destination D (under a far-field assumption with respect to the IRS) can be expressed as

$$y_D = \mathbf{w}_D^T \mathbf{H}_{RD} \Phi \mathbf{H}_{SR} \mathbf{w}_S x_S + \mathbf{w}_D^T \mathbf{H}_{SD} \mathbf{w}_S x_S + \mathbf{w}_D^T \mathbf{n}_D. \quad (2.6)$$

2.1.3 A Signal Model for the AF Relay

AF relays have been studied in the context of cooperative communications as a means to regenerate a relayed signal through amplification, with the goal of improving the system capacity. Unlike IRSs, AF relays feature a non-negligible power consumption, and introduce noise amplification.

In this work we consider as AF relay a device equipped with M_T transmit and M_R receive antennas. Therefore, the signal received at D is:

$$\begin{aligned} y_D = & \mathbf{w}_D^T \mathbf{H}_{RD} \Phi \mathbf{H}_{SR} \mathbf{w}_S x_S + \mathbf{w}_D^T \mathbf{H}_{SD} \mathbf{w}_S x_S \\ & + \mathbf{w}_D^T \mathbf{H}_{RD} \Phi \mathbf{n}_R + \mathbf{w}_D^T \mathbf{n}_D, \end{aligned} \quad (2.7)$$

where in this case matrix Φ also accounts for the amplification gain, and its structure depends on the specific relay design. Moreover, \mathbf{n}_R represents the circularly symmetric complex Gaussian noise vector with covariance matrix $\sigma_{N_R}^2 \mathbf{I}_{M_R}$. Then, the power of the noise term relayed by the AF relay to receiver D and measured after the combiner at the UE, is

$$\begin{aligned} \hat{\sigma}_{N_R}^2 &= \left(\mathbf{w}_D^T \mathbf{H}_{RD} \Phi \right) \left(\mathbf{w}_D^T \mathbf{H}_{RD} \Phi \right)^H \sigma_{N_R}^2 \\ &= \mathbf{w}_D^T \mathbf{H}_{RD} \Phi \Phi^H \mathbf{H}_{RD}^H \mathbf{w}_D^* \sigma_{N_R}^2. \end{aligned} \quad (2.8)$$

2.1.4 A Full-Stack Simulator for IRS/AF Relays

Despite the availability of accurate sub-6 GHz and mmWave channel models, analytical evaluations of the 5G NR protocol stack introduce several assumptions in the system architecture, and are generally not desirable [36]. Additionally, 5G/6G cellular networks are rapidly shifting towards open and controllable network configurations, which further introduce unprecedented data-driven programmability [37]. In these regards, computer simulators are emerging as a valuable tool to let researchers better understand the performance of wireless networks, and dimension them accordingly [38].

2 Simulation tools for future cellular networks

Several simulators for 5G cellular and vehicular networks are available in the literature [17, 18, 39–44]. However, they provide a detailed characterization of either the lower (i.e., at the link level) or the upper (i.e., at the system level) layers of the 5G NR protocol stack. Notably, the latter sacrifice PHY layer accuracy to reduce the computational complexity, but incorporate accurate models of the remainder of the protocol stack, thus enabling scalable end-to-end simulations. Despite the many software-based evaluation platforms available, to the best of our knowledge there are no end-to-end simulators for IRSs and AF relays. In [45], the authors presented an open-source module for IAB, even though it was not extended to support passive relays like IRSs. Moreover, the authors in [46] presented an ns-3 IRS module, but their work focused on vehicular networks, and did not consider the case of AF relays. In this work we fill this gap by proposing an ns-3-based simulator for IRSs and AF relays. Arguably, the main effect of the presence of these entities is the alteration of the wireless channel between the communication endpoints. Accordingly, our simulator extends the ns-3 mmwave module [17] (among the most popular 5G-oriented NR-compliant frameworks to simulate 5G networks) by implementing a new signal model for IRS and AF relays, following the characterization in Secs. 2.1.2 and 2.1.3, respectively, which is then used to compute the SINR experienced by signals transmitted over a relayed wireless link.

2.1.4.1 Implementation of the IRS/AF Signal Model

In line with [22], we assume that the transmission of the signal x_S occurs over a frequency-selective wireless channel as 5G NR supports network operations with a bandwidth up to 400 MHz, when using FR2 [47]. Therefore, the evaluation of the SINR requires, among other things, the computation of the Power Spectrum Density (PSD) of the useful component of the signal at D, i.e., \mathcal{P}_{rx} , starting from that of the input signal \mathcal{P}_{tx} . Additionally, we consider that both the transmitter and the receiver feature m-MIMO arrays equipped with multiple antenna elements, and use the beamforming vectors w_S and

\mathbf{w}_D , respectively. Under these assumptions, the input-output relationship in (2.3) becomes [32, 48]:

$$y_D = \mathbf{w}_D^T \mathbf{H}_{RD} \Phi \mathbf{H}_{SR} \mathbf{w}_S x_S + \mathbf{w}_D^T \mathbf{H}_{SD} \mathbf{w}_S x_S + \tilde{n} + \sum_{I=1}^N \mathbf{w}_D^T \mathbf{H}_{RD} \Phi \mathbf{H}_{IR} \mathbf{w}_I x_I + \sum_{I=1}^N \mathbf{w}_D^T \mathbf{H}_{ID} \mathbf{w}_I x_I, \quad (2.9)$$

where in turn \tilde{n} is defined as:

$$\tilde{n} = \begin{cases} \mathbf{w}_D^T \mathbf{n}_D & \text{if IRS,} \\ \mathbf{w}_D^T \mathbf{n}_D + \mathbf{w}_D^T \mathbf{H}_{RD} \Phi \mathbf{n}_R & \text{if AF,} \end{cases}$$

where matrix Φ is the relay matrix, i.e., a matrix which fully encodes the effect of the relay, i.e., either IRS or AF, as described in Secs. 2.1.2 and 2.1.3 for the single user case, respectively, over the wireless channel. Notably, S and D are either in NLoS (in this case they communicate via the relay, and we consider the direct link towards D to be unavailable), or in LoS (in this case they do not use the relay). Accordingly, assuming that the source of interest is in NLOS with respect to its intended destination, (2.9) becomes:

$$y_D = \mathbf{w}_D^T \mathbf{H}_{RD} \Phi \mathbf{H}_{SR} \mathbf{w}_S x_S + \sum_{\hat{I} \in I_{LOS}} \mathbf{w}_D^T \mathbf{H}_{ID} \mathbf{w}_{\hat{I}} x_{\hat{I}} + \sum_{\bar{I} \in I_{NLOS}} \mathbf{w}_D^T \mathbf{H}_{RD} \Phi \mathbf{H}_{IR} \mathbf{w}_{\bar{I}} x_{\bar{I}} + \tilde{n}, \quad (2.10)$$

where I_{LOS} and I_{NLOS} are the two disjoint sets of interferers which experience either an LoS or an NLoS channel towards D, respectively. Then, the PSD of the useful component of the signal at the receiver can be written as:

$$\mathcal{P}_{rx}(t, f) = \mathcal{P}_{tx}(t, f) \|\mathbf{w}_D^T \mathbf{H}_{RD} \Phi \mathbf{H}_{SR} \mathbf{w}_S\|^2. \quad (2.11)$$

Based on the above definitions, our simulator computes the PSD by checking whether the communication from S to D involves a relay. If so, the PSD is computed according to the following steps.

1. *Channel matrices generation.* After having identified S and D as the two endpoints of the communication, the channel matrices \mathbf{H}_{SR} and \mathbf{H}_{RD} are computed based on (2.1) [23].

2. Configuration of the relay and beamforming vectors. We assume that the choice of the beamforming vectors for both S and D (\mathbf{w}_D and \mathbf{w}_S), as well as the relay configuration (Φ), consist in the choice of a *codeword* from a pre-defined *codebook*. The latter is computed offline by first defining a set of beam directions $\{\omega_{n,m}\}$ which scan a given angular sector via steps of Half Power Beamwidth (HPBW). In particular, let $\mathbf{a}_{n,m}$ be the steering vector corresponding to direction $\omega_{n,m}$. This is computed as:

$$\mathbf{a}_{n,m} = \left[1, \dots, e^{j\frac{2\pi}{\lambda}d(i_H \sin \alpha_n \sin \beta_m + i_V \cos \beta_m)}, \dots, e^{j\frac{2\pi}{\lambda}d((N_H-1) \sin \alpha_n \sin \beta_m + (N_V-1) \cos \beta_m)} \right]^T, \quad (2.12)$$

where $0 \leq i_H \leq N_H$ ($0 \leq i_V \leq N_V$) is the horizontal (vertical) index of an antenna element, N_H and N_V are the number of antenna elements in the horizontal and vertical direction, respectively, and α_n and β_m represent the azimuth and the elevation angles of $\omega_{n,m}$, respectively. Then, we define the codebook for the UEs, Next Generation Node Bases (gNBs) and AF relays as the set $\{(\sqrt{N_H N_V})^{-1} \mathbf{a}_{n,m}\}$, while the IRS codebook is defined as $\{\mathbf{a}_{n,m}\}$.

Moreover, we assume that the devices do not have full channel knowledge, i.e., they do not know the realizations of \mathbf{H}_{SR} and \mathbf{H}_{RD} . Then, in line with the 5G NR beam management procedure [49], the choice of the codeword in the codebook is performed via exhaustive search, i.e., by repeatedly sending pilot signals, and measuring the SINR experienced with various configurations of the codebook. Eventually, we choose the combination of \mathbf{w}_D , \mathbf{w}_S , and Φ yielding the highest SINR.

Notably, this procedure is not repeated at each transmission opportunity. Instead, \mathbf{w}_D , \mathbf{w}_S , and Φ are stored and re-used for the whole channel coherence time, to mimic the actual 5G NR beam management procedure, and also reduce the complexity of the simulations. Furthermore, the evaluation of the SINR is performed by neglecting the small-scale fading terms, to further reduce the overhead. The small-scale fading will be eventually incorporated in Step 4 of the model.

3. *Long-term computation.* Along the lines of [22], the PSD of the transmitted signal x_S at D can be expressed as:

$$\begin{aligned}\mathcal{P}_{rx}(t, f) &= \\ &= \mathcal{P}_{tx}(t, f) \|\mathbf{w}_D^T \mathbf{H}_{RD} \boldsymbol{\Phi} \mathbf{H}_{SR} \mathbf{w}_S\|^2 \\ &= \mathcal{P}_{tx}(t, f) \|\mathbf{w}_D^T \mathbf{H}_{SRD} \mathbf{w}_S\|^2 \\ &= \mathcal{P}_{tx}(t, f) \left\| \sum_{d=1}^{N_D} \sum_{s=1}^{N_S} w_d^D h_{d,s}^{SRD}(t, f) w_s^S \right\|^2.\end{aligned}\quad (2.13)$$

In Eq. (2.13), \mathbf{H}_{SRD} is the equivalent channel matrix between S and D, whose generic entry $h_{d,s}^{SRD}(t, f)$ is:

$$\begin{aligned}h_{d,s}^{SRD}(t, f) &= [\mathbf{H}_{RD}(t, f) \boldsymbol{\Phi} \mathbf{H}_{SR}(t, f)]_{d,s} \\ &= \sum_{n=1}^{N_{RD}} \sum_{m=1}^{N_{SR}} \sum_{k=1}^{N_R} \sum_{l=1}^{N_R} h_{d,k,n}^{RD} \phi_{k,l} h_{l,s,m}^{SR} \\ &\quad \times e^{j2\pi v_n t} e^{j2\pi \tau_n f} \\ &\quad \times e^{j2\pi v_m t} e^{j2\pi \tau_m f},\end{aligned}\quad (2.14)$$

where N_{RD} and N_{SR} are the number of multipath clusters in \mathbf{H}_{RD} and \mathbf{H}_{SR} , respectively. Moreover, w_s^S and w_d^D denote entries s and d of vectors \mathbf{w}_S and \mathbf{w}_D , respectively. Then, Step 3 consists in the evaluation of the long-term fading:

$$L_{n,m} \doteq \sum_{d=1}^{N_D} \sum_{s=1}^{N_S} \sum_{k=1}^{N_R} \sum_{l=1}^{N_R} w_d^D h_{d,k,n}^{RD} \phi_{k,l} h_{l,s,m}^{SR} w_s^S. \quad (2.15)$$

4. *Small-scale fading and path loss.* The small-scale fading terms are combined with the terms $L_{n,m}$ to compute the overall fading component of the PSD of interest:

$$\tilde{\mathcal{P}}_{rx}(t, f) = \mathcal{P}_{tx}(t, f) \left\| \sum_{n=1}^{N_{RD}} \sum_{m=1}^{N_{SR}} L_{n,m} E_{n,m} \right\|^2, \quad (2.16)$$

where

$$E_{n,m} \doteq e^{j2\pi v_n t} e^{j2\pi \tau_n f} e^{j2\pi v_m t} e^{j2\pi \tau_m f}. \quad (2.17)$$

Additionally, the path loss is computed as in (2.2). Since the useful signal received at D experiences two channels (from S to R, and from R to D) as a cascade, as described in (2.9), two path loss terms are added (in dB), to obtain the final PSD of x_S at D as:

$$\begin{aligned}\mathcal{P}_{rx}(t, f)[\text{dB}] &= \text{PL}(d_{SR}, f_c)[\text{dB}] \\ &\quad + \text{PL}(d_{RD}, f_c)[\text{dB}] + \tilde{\mathcal{P}}_{rx}(t, f)[\text{dB}].\end{aligned}\tag{2.18}$$

5. *Interference and SINR.* As the last step, we evaluate the PSDs $\{\mathcal{P}_i(t, f)\}_{i=1,\dots,N_I}$ of the N_I interfering signals at D. To do so, we follow Steps 1–4 as for the useful component of the signal. However, the beamforming configurations are not optimized as described in Step 2. That is to say, each interferer uses the beamforming vector yielding the highest SINR towards its intended destination, while R and D employ the same configurations used in the previous steps. Finally, the SINR is evaluated as:

$$\Lambda(t, f) = \frac{\mathcal{P}_{rx}(t, f)}{\sum_{i=1}^{N_I} \mathcal{P}_i(t, f) + \mathcal{P}_n(t, f)},$$

where $\mathcal{P}_n(t, f)$ is the PSD of the thermal noise at D.

2.1.4.2 Integration of the IRS/AF Signal Model in the Simulator

In Section 2.1.4.1 we described how our simulator computes the channel (in terms of PSD) in case of IRS/AF relays, which is then used to calculate the end-to-end SINR at the destination D. Notice that the SINR can refer to either the SINR relative to the whole bandwidth, for narrowband signals over frequency-flat channels, or the SINR experienced over a single subcarrier, for wideband signals transmitted over frequency-selective channels. In the second case, the SINRs corresponding to the various frequency chunks are then mapped into a single SINR value, according to additional maps obtained from link-level simulations [25]. Based on that, our simulator defines a Link-to-System Mapping (L2SM), i.e., a table which associates a given SINR to a MAC-layer Transport Block (TB) error rate [50], in turn used to decide whether the TB has been correctly received or not.

The upper layers of the 5G NR protocol stack are modeled based on the ns3-mmwave module [17], which implements a custom PHY layer supporting the NR frame structures and numerologies, and a MAC layer with ad hoc

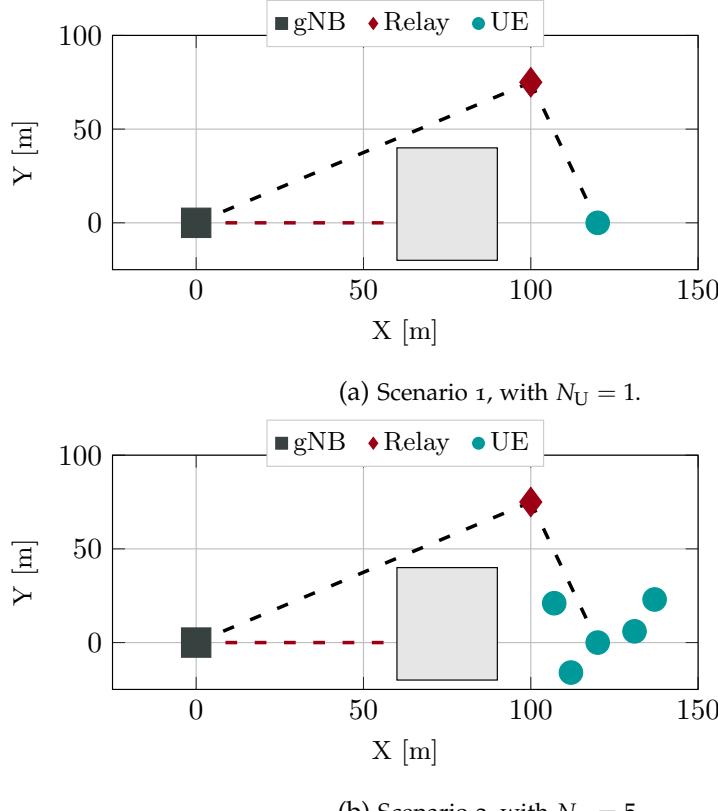


Figure 2.2: Simulation scenarios, where we deploy one gNB, N_U UEs and, possibly, a relay. A building (the gray rectangle) blocks the direct link (dashed red line) from the gNB to the UEs. In turn, the relay guarantees a LoS link (dashed black line) to all the devices.

beamforming and scheduling policies. The Radio Link Control (RLC) and Packet-Data Convergence Protocol (PDCP) layers implement network functions such as packet segmentation, retransmissions and/or reassembly.

2.1.5 Performance Evaluation

In this section we describe our simulation setup and parameters (Section 2.1.5.1), and evaluate the performance of IRSs and AF relays, considering full-stack network metrics as a function of different antenna array configurations (Section 2.1.6).

Table 2.1: Simulation parameters.

Parameter	Value
Carrier frequency	28 GHz
Total bandwidth	100 MHz
Number of UEs (N_U)	{1, 5}
gNB antenna array	8H×8V
gNB max RF power	33 dBm
UE antenna array	2H×1V
IRS antenna array	{10H×20V, 20H×40V, 40H×80V, 60H×120V}
AF antenna array	{4H × 4V, 8H × 8V, 16H × 16V}
AF amplification	40 dB
Antenna radiation pattern	[23, Table 7.3-1]
UDP source rate	50 Mbps

2.1.5.1 *Simulation Setup*

In our simulations we consider two simple yet realistic urban canyon scenarios, where we deploy a single gNB, N_U UEs, with $N_U = 1$ (5) in Scenario 1 (2), as illustrated in Figure 2.2, and a single relay, which can be either an IRS or an AF relay. The wireless channel is modeled as an Urban Macro (UMa) link [23]. The LoS/NLoS condition depends on the geometry of the scenario. In particular, we assume that the direct wireless link between the UEs and the gNB is blocked by a building, as illustrated in Figure 2.2, which introduces an additional penetration loss modeled based on [23, Section 7.4.3.1]. The end nodes can still communicate in LoS via the relay. Furthermore, we assume that at each Transmission Time Interval (TTI) the relays can arbitrarily switch configuration to serve a given user and that their phase shifters have infinite resolution, i.e., we do not account for quantization loss.

Our simulation parameters are reported in Table 2.1. Specifically, the UEs download User Datagram Protocol (UDP) data, modeled as a constant bit-rate stream of 50 Mbps, from a remote server. We assume that, at each transmission opportunity towards the generic k -th UE, both AF and IRS relays can use their optimal configuration, i.e., the codeword yielding the highest end-to-end SINR towards UE k . The system operates at 28 GHz, with a total bandwidth of 100 MHz, to be shared among all the devices in Time Division Multiple Access (TDMA). The gNB is equipped with an antenna array of 64

elements, and uses a power of 33 dBm. For the IRS, we consider a number of reflecting elements from 200 to 7200. For the AF relay, we consider antenna arrays from 16 to 256 elements.

2.1.6 Numerical Results

We now compare the end-to-end performance of IRS- and AF-relay assisted networks in terms of:

- *SINR*. It is a measure of the quality of the channel. It depends on PHY-layer characteristics, including the relative distance between the transmitter, the receiver and the relay (if applicable), the operating frequency, the propagation conditions, and the channel bandwidth.
- *End-to-end throughput*. It is measured as the total number of received bytes per user divided by the total simulation time.
- *End-to-end latency*. It is measured from the time each packet is generated at the application layer to when it is successfully received. Accordingly, it accounts for both transmission and queuing times.
- *Packet Error Rate (PER)*. It is measured as the ratio between the number of packets delivered with errors and the total number of transmitted packets.

The IRS/AF performance will be evaluated against a baseline scenario (referred to as “gNB-only”) in which there is no intermediate relay.

SINR Our analysis starts with the SINR statistics depicted in Figure 2.3, relative to Scenario 1 with $N_U = 1$. First, in Figure 2.3a we observe that the presence of the relay improves the SINR (on average up to +55 dB) compared to the “gNB only” baseline, in which the UE communicates in NLoS. Notably, as depicted in Figure 2.3b, both IRS and AF relays provide an end-to-end SINR gain which scales proportionally with respect to the number of radiating elements at the relay. For the IRS, this effect is given by the beamforming gain, as well as by the fact that the power collected by the IRS is proportional to its surface area, which in turn is proportional to the number of radiating elements [51].

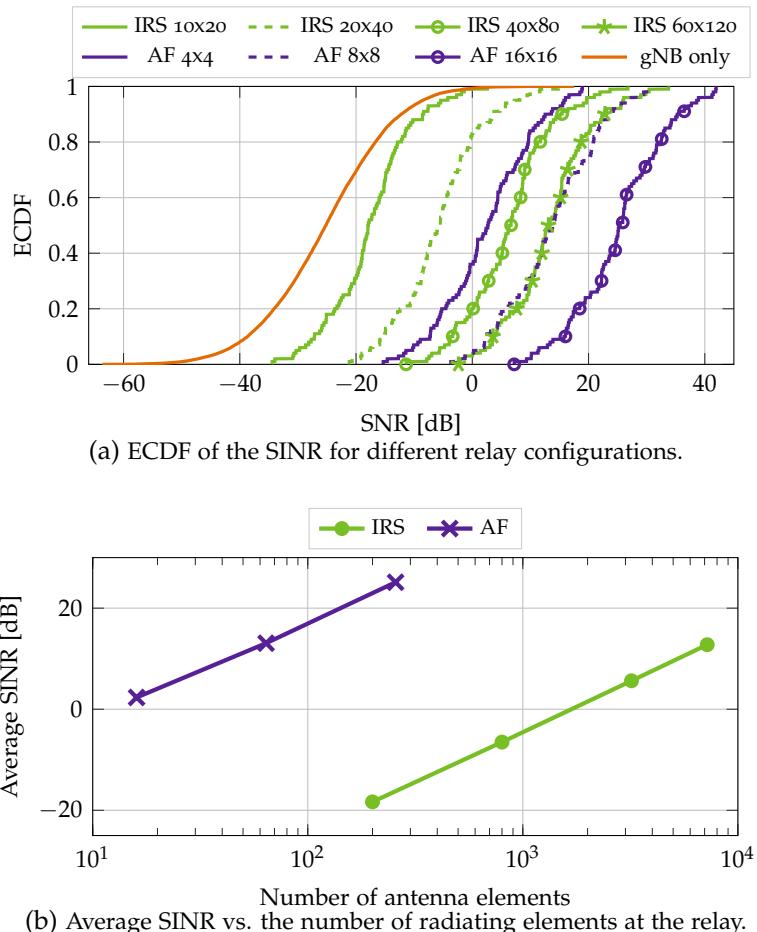


Figure 2.3: SINR statistics for Scenario 1.

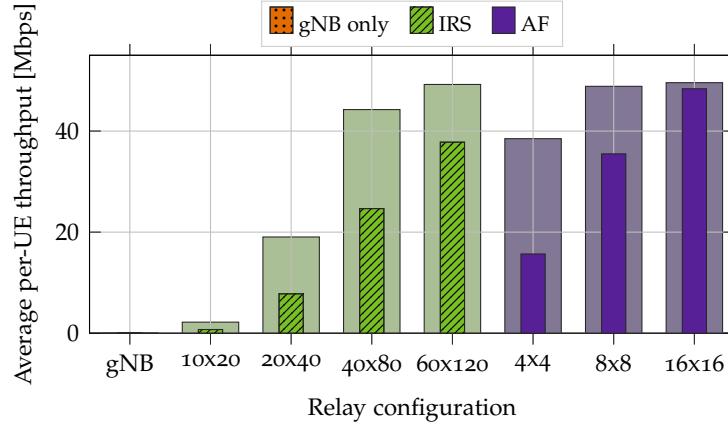


Figure 2.4: End-to-end per-UE throughput at the application layer in Scenario 1 (wide bars) and Scenario 2 (narrow bars) for different relay configurations.

The AF-assisted configurations always outperform the IRS-assisted ones in terms of SINR (on average up to +40 dB, with the same number of antennas): this is expected since the AF relay amplifies the signal, thus achieving a higher end-to-end gain. Notice that the SINR is below 0 dB when the IRS is made of fewer than 800 elements, which justifies the use of very large IRS panels. Indeed, an IRS panel of 60×120 elements provides an average SINR of 13 dB, which is enough to support reliable transmissions as long as communication requirements are not too extreme, as we will demonstrate in the following paragraphs.

End-to-end throughput In Figure 2.4 we plot the end-to-end throughput experienced at the application layer, thus considering the impact of the whole 5G NR protocol stack. When $N_U = 1$ (Scenario 1) the average throughput is an indication of the ergodic capacity. We see that the throughput for the “gNB only” baseline is zero, given the very low SINR (below the sensitivity threshold of most commercial receivers) experienced at the physical layer. Interestingly, even though the AF relay with 16×16 antennas guarantees, on average, 15 dB higher SINR than an IRS with 60×120 elements (from Figure 2.3a), we see that the end-to-end throughput of the two configurations is comparable. This demonstrates that, in a simple scenario with only one UE, an average SINR of 15 dB is enough to satisfy all traffic requests. In this case, the IRS is more desirable than an AF relay given its simplicity. Also,

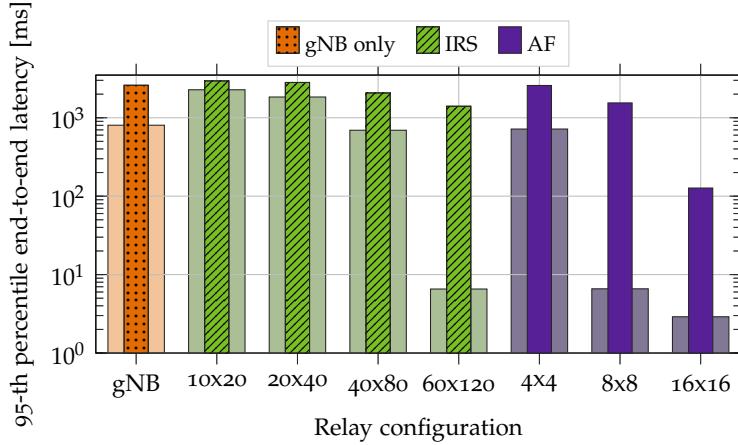


Figure 2.5: 95-th percentile of the end-to-end latency at the application layer in Scenario 1 (wide bars) and Scenario 2 (narrow bars) for different relay configurations.

it is not convenient to further increase the IRS size, given that the throughput is already maximized and equal to the UDP source rate (50 Mbps in our simulations).

When $N_U = 5$ (Scenario 2) the average per-UE throughput decreases significantly with respect to Scenario 1 due to the fact that, in a multi-user scenario, radio resources must be shared among UEs, which may lead to channel congestion. This result validates the accuracy and realism of our ns-3 framework. Nevertheless, this effect is less pronounced for very large antenna panels. For example, for an AF relay of 4×4 antennas, the per-UE throughput drops by almost 60%, while considering an array of 16×16 elements the per-UE throughput decreases by only 2%. Even in Scenario 2, AF-assisted networks can still sustain the application source rate, as long as at least 16×16 antennas are used. On the other hand, IRSs are constrained by the limited SINR available at the PHY layer, and are never able to achieve the full source rate offered by the application. The maximum achievable throughput is around 40 Mbps for 60×120 elements, i.e., -20% compared to the case of $N_U = 1$.

End-to-end latency Finally, in Figure 2.5 we plot the 95-th percentile of the end-to-end latency experienced at the application layer. We can see that the performance is generally poor even in the simple scenario in which only one UE is deployed (Scenario 1), where the latency is higher than 100 ms for most

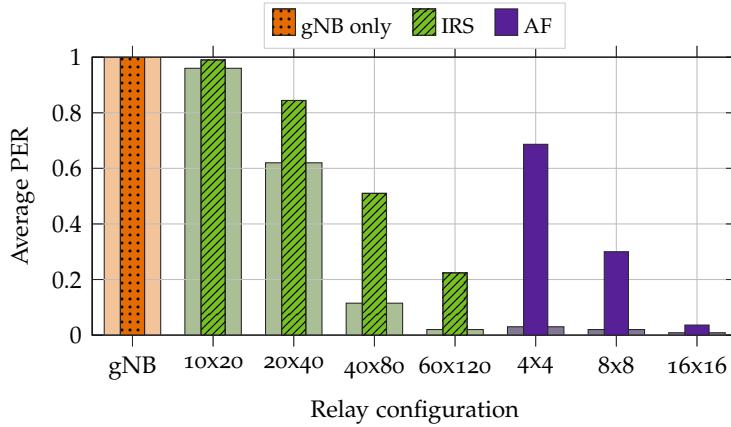


Figure 2.6: Average PER at the application layer in Scenario 1 (wide bars) and Scenario 2 (narrow bars) for different relay configurations.

relay configurations, suggesting that in these cases the system is unstable. In fact, the use of relays featuring small antenna panels results in very high levels of queuing and buffering, which leads to latency degradation. This issue can be solved by configuring larger IRS and AF relays, despite the increased system complexity. For example, an IRS of 60×120 elements and an AF relay with $\geq 8 \times 8$ elements can guarantee an end-to-end latency lower than 10 ms, which is in line with most 5G application requirements. Notice that the latency for the “gNB only” configuration is not particularly representative, as it is relative to only the correctly received packets. In fact, without the relay, transmissions are in NLoS and result in several packet losses (see the PER in Figure 2.6), which makes the system less congested; the (few) packets that make it to the application layer are then transmitted with lower delay. Nevertheless, the latency is still more than two orders of magnitude higher than considering the best IRS and AF configurations, an indication that relays are desirable in these types of networks.

When $N_U = 5$ (Scenario 2), the latency is generally higher compared to when $N_U = 1$. This is expected since UEs are competing for the available resources. In addition, using UDP as transport protocol, thus with a full buffer source traffic model, each end-to-end flow does not self-regulate to the actual network conditions, thus congestion arises. Better performance could be achieved considering non-UDP traffic: for example, the congestion control mechanism available in TCP would regulate the source traffic, and prevent network congestion and buffer overflow.

2 Simulation tools for future cellular networks

Notice that, even considering the most aggressive IRS architecture with 60×120 elements, the latency is on average above 1000 ms, vs. 6.5 ms in Scenario 1. This is due to the fact that, in this scenario, more than 20% of the packets are lost and retransmitted (see Figure 2.6), which increases the packet delay. For an AF relay with 16×16 antennas, instead, the latency is more than 10 times lower and equal to around 130 ms on average, with a PER as low as 3%, which can still support some key target communication requirements. We can conclude that IRS-assisted networks, though consuming less power, are not appropriate in this scenario, unless very large IRS panels are used.

2.2 Modeling non-terrestrial wireless channels

Satellites have been used from the 1990s to provide basic services such as phone and Internet access. Notably, Geostationary Equatorial Orbit (GEO) satellites orbit at 35 786 km, and offer global coverage at limited costs, despite the huge propagation delays. Only in the early 2010s have the costs for satellite launch and maintenance become low enough to allow for huge constellations of satellites to be launched in the Low Earth Orbit (LEO) [52], providing wide coverage on the Earth, while promoting low latency. Besides satellites, both High Altitude Platforms (HAPs) and Unmanned Aerial Vehicles (UAVs) stand out as valid cost-effective alternatives for NTNs. UAVs, flying at low altitudes (typically no more than 1 km), can guarantee on-demand support for ground networks, for example providing immediate assistance when cellular towers are overwhelmed or unavailable. HAPs operate in the stratosphere (from 20 to 50 km), and can be used to shape large coverage beams in unpopulated areas, or provide services like backhauling and Mobile Edge Cloud (MEC), e.g., to gather and process data generated on the ground [53, 54].

Despite these premises, however, communication using space or airborne vehicles introduces new challenges compared to a terrestrial base station, including (i) severe PL due to the longer propagation distance, (ii) additional attenuation from the atmosphere, such as scintillation, rain and clouds, (iii) Doppler shift due to the orbital mobility of satellites, and (iv) additional delays, mainly for propagation. While experiments with real testbeds are impractical due to limitations in the scalability and flexibility of platforms, as well as the high cost of hardware components, the option to test network configurations via simulations in a sandbox environment facilitates the research process. Furthermore, an open-source simulator encourages research in the field, and offers industries and research institutions a better way to categorize and evaluate technologies. However, the de-facto standard end-to-end simulator ns-3 currently implements the TR 38.901 channel model only, thus lacking its NTN extension outlined in TR 38.811 [26]. To fill this gap, we implement the 3GPP NTN model in ns-3.

2 Simulation tools for future cellular networks

2.2.1 Scenarios and Path Loss Condition

Similarly to the cellular channel model in [29], the NTN model gives the option to run simulations in four scenarios, to represent different propagation environments. Specifically:

- Dense Urban: Extremely dense environment, with tall buildings acting as potential blockers.
- Urban: City environment, with buildings.
- Suburban: Small city, with up to two-storey buildings.
- Rural: Open field environment, with little or no buildings.

For satellites, only outdoor communication is possible, since attenuation from buildings would be enough to make the signal unusable. For HAPs or UAVs, instead, indoor communication is feasible, even though not yet implemented in our module.

The 3GPP defines both LoS and NLoS propagation, where the probability depends on the scenario and the elevation angle. The latter is defined as the angle between the horizon plane of the ground terminal and the vector pointing to the NTN platform.

2.2.2 Path Loss

The basic path loss (in dB) can be written as

$$PL_b = FSPL(d, f_c) + SF + CL(\alpha, f_c). \quad (2.19)$$

The first term is the free space path loss, which can be calculated as

$$FSPL(d, f_c) = 32.45 + 20\log_{10}(f_c) + 20\log_{10}(d), \quad (2.20)$$

where f_c is the carrier frequency in GHz and d is the distance between the transmitter and the receiver in meters. Notice that, while the channel model is valid for frequencies from 0.5 GHz to 100 GHz, two frequency bands are targeted in NTN, i.e., the S-band for frequencies below 6 GHz and the Ka-band for frequencies of 20 (30) GHz for downlink (uplink) transmissions, thereby in the millimeter-wave spectrum [55].

In Eq. (2.19), SF represents the Shadow Fading (SF), and is modeled as a log-normal random variable of zero mean and variance σ_{SF}^2 , i.e., $SF \sim N(0, \sigma_{SF}^2)$. In order to calculate this variance, the model requires four parameters: the type of scenario, the frequency, the path loss condition (LoS or NLoS), and the elevation angle. These parameters are used to find the correct entry in a table, given in [26]. A similar process is needed to calculate the clutter loss $CL(\alpha, f_c)$.

2.2.3 Atmospheric Absorption

The attenuation introduced by the presence of atmospheric gasses was a marginal factor in the terrestrial channel. This is no longer the case for the NTN channel, where atmospheric absorption plays a crucial role in the overall link budget. A complete and accurate characterization of the atmospheric attenuation is given in the ITU model [56], and depends on a set of parameters which is usually difficult to retrieve in simulations, such as absolute humidity, dry air pressure, water-vapour density and water-vapour partial pressure. Therefore, the 3GPP offers a simplified method considering only ground users placed at the sea level, with an elevation angle fixed to 90 degrees, and considering mean annual global values for the rest of the parameters. For elevation angles different than 90 degrees, the calculation is straightforward. Given the zenith attenuation $A_{zenith}(f_c)$, the additional path loss due to atmospheric gasses is

$$PL_{A,dB}(\alpha, f_c) = \frac{A_{zenith}(f_c)}{\sin(\alpha)}, \quad (2.21)$$

where α is the actual elevation angle. Atmospheric absorption should be considered only for frequencies above 10 GHz, or for any frequency in case of $\alpha < 10$ degrees.

2.2.4 Scintillation

Scintillation corresponds to the rapid fluctuation in amplitude and phase of the received signal, caused by the variation of the refractive index of the channel. Specifically, scintillation depends on location, time of the day, season, and solar and geomagnetic activity. Stronger levels of scintillation are

2 Simulation tools for future cellular networks

observed only at high latitudes (more than 60 degrees), in auroral and polar regions.

While a complete absorption model would unnecessarily complicate system-level simulations, the 3GPP recommends a simplified model for scintillation [26], which is structured into two components: ionospheric scintillation and tropospheric scintillation.

2.2.4.1 Ionospheric Scintillation

Ionospheric scintillation is modeled based on the Gigahertz Scintillation Model [57]. While for the purpose of system-level simulations ionospheric scintillation is generally negligible at mid latitudes (between 20 and 60 degrees) or at above-6 GHz frequencies, in all other latitudes and conditions it is modeled as

$$PL_{S,dB} = \left(\frac{f_c}{4} \right)^{-1.5} \frac{P_{fluc}(4 \text{ GHz})}{\sqrt{2}}, \quad (2.22)$$

where f_c is the carrier frequency, and $P_{fluc}(4 \text{ GHz})$ is a scaling factor representing the ionospheric attenuation level at 99% of the time observed in Hong Kong between March 1977 and March 1978 at a frequency of 4 GHz [26, Figure 6.6.6.1.4-1].

2.2.4.2 Tropospheric Scintillation

Unlike ionospheric scintillation, the effect of tropospheric scintillation increases with the frequency, and becomes significant above 10 GHz. Furthermore, it increases at low elevation due to the longer path of the signal. In these conditions, tropospheric scintillation is due to sudden changes in the refractive index due to the variation of temperature, water vapor content, and barometric pressure. For system-level simulations, the additional attenuation due to tropospheric scintillation is modeled as the attenuation level at 99% of the time observed in Tolouse at 20 GHz, reported in [26, Figure 6.6.6.2.1-1].

2.2.5 Fast Fading

As far as the fading is concerned, the 3GPP introduces both a flat-fading and a frequency-selective model. The flat-fading model, however, can be applied only if specific conditions are met, including (i) minimum elevation angle of

20 degrees, (ii) quasi-LOS propagation, (iii) communication in the S-band, (iv) channel bandwidth of at most 5 MHz, and (v) rural, suburban or urban scenario. Hence, we consider the more general, though complex, frequency-selective fading model for NTN. Then, the fading is based on the TR 38.901 model for cellular networks [29] (already implemented in ns-3 in the mmwave module [58]), but with different parameters as described in [26, Section 6.7.2].

2.2.6 Antenna Model

Different antenna models are defined, depending on the device (satellite, HAP or UAV, and ground terminal). For satellites, the 3GPP suggests to use a circular aperture antenna model. Circular aperture antennas are reflector antennas that offer circular polarization. The normalized antenna gain pattern is given by

$$G(\theta) = \begin{cases} 1 & \theta = 0; \\ 4 \left| \frac{J_1(k \cdot \ell \cdot \sin \theta)}{k \cdot \ell \cdot \sin \theta} \right|^2 & 0 < |\theta| \leq 90^\circ, \end{cases} \quad (2.23)$$

where $J_1(\cdot)$ is the Bessel function of the first kind and first order, ℓ is the radius of the antenna's circular aperture and, given a carrier frequency f_c , the value of k is equal to $k = 2\pi f_c / c$, where c is the speed of light in vacuum.

When considering flying vehicles that are not satellites, such as HAPs or UAVs, the lower distance makes it possible to use Uniform Planar Array (UPA) antennas, that is the current standard for UEs and eNB/gNB nodes in cellular networks according to the TR 38.901 model [29].

Finally, for terrestrial terminals, the 3GPP suggests to use either UPA antennas or Very Small Aperture Terminal (VSAT) antennas. The latter model, in particular, is a common choice in satellite communication, and consists of a circular reflector antenna of small size (less than 1 m of diameter), to be typically placed on roofs pointing at the sky. The VSAT radiation pattern is the same as that of circular aperture antennas for satellites, and is given in Eq. 2.23.

2.2.7 Coordinate System

In general, the 3GPP defines a simple Cartesian coordinate system where the position of each node is uniquely described by a set of three values, (x, y, z) ,

2 Simulation tools for future cellular networks

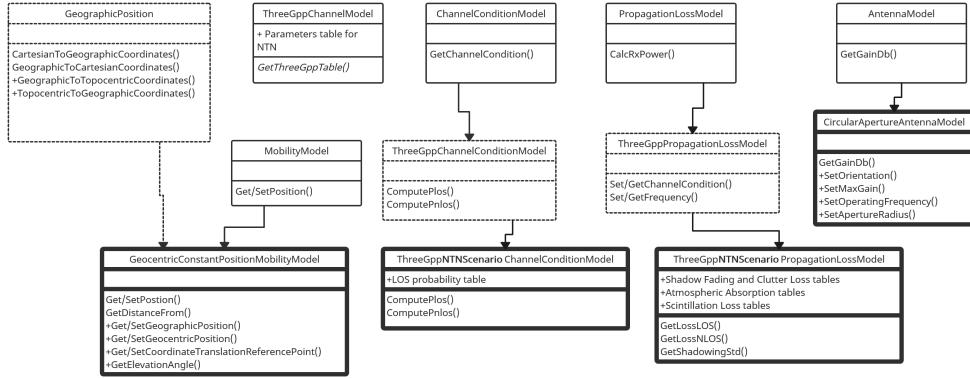


Figure 2.7: Simplified UML diagram, which depicts the most significant changes which we introduced to ns-3. Bold classes represent the newly implemented ones, while dotted classes are pre-existing ones that have been modified.

where x and y define the ground plane, and z represents the height of the node. While this model is accurate enough to describe scenarios where nodes are deployed at close distance (e.g., a few hundreds of meters), it is not for scenarios where end nodes are placed hundred (or thousands) of kilometers apart such as in the NTN environment. In this case, the Earth's curvature, as well as the elevation angle, play an important role. Therefore, the 3GPP suggests to use a Geocentric Cartesian coordinate system (or Earth-Centered Earth-Fixed (ECEF) system), where the position of a node is still described by three values (x, y, z) , but now the origin of the axes lays in the center of the Earth, which is modeled as a sphere of radius $R = 6371$ km. The x-y plane defines the equatorial plane, with the x-axis pointing at 0-degree longitude, the y-axis pointing at 90-degree longitude, and the z-axis pointing at the geographical North Pole. Then, terrestrial nodes on the surface of the Earth are deployed so that $\sqrt{x^2 + y^2 + z^2} = R$, while aerial/space nodes flying/orbiting around the Earth are deployed so that $\sqrt{x^2 + y^2 + z^2} > R$.

2.2.8 Implementation in ns-3

The proposed ns-3 implementation of the 3GPP TR 38.811 [26] NTN channel model is based upon the TR 38.901 model presented in [58]. Despite the fact that the newly introduced methods and classes are designed with the goal of introducing minimal changes to the existing ns-3 APIs, some modifications

to the existing code structure are still required. A schematic of these changes, which we make publicly available¹, can be found in Figure 2.7. The remainder of this section describes more in detail our implementation of the NTN channel model of [26] in ns-3.

2.2.8.1 Small-scale fading

The most significant modifications to the pre-existing ns-3 classes concern the `ThreeGppChannelModel` class, which computes the small-scale propagation phenomena in the form of a complex channel matrix. Indeed, conversely from the procedure implemented in [58], in the NTN channel model of [26] most channel parameters depend on all the propagation scenario, LoS condition, carrier frequency and elevation angle variables. To account for this, we store the small-scale fading parameters in a *nested map*. The choice of this data-structure is motivated by the good trade-off between code efficiency and readability which it provides, considering that the possible combinations of the input parameters is particularly high, i.e., 144.

In particular, we change the signature of the `GetThreeGppTable()` method to include the `MobilityModel` instances of both transmitting and receiving nodes. In such a way, we account for the dependence of the small scale parameters with respect to the elevation angle, which the model of [26] exhibits. Finally, we include the required angular scaling factors for the propagation scenarios that have a lower number of clusters than the ones described in TR 38.901 [29].

2.2.8.2 Coordinate systems

Instead of the coordinate system described in [29], the NTN channel model of [26] considers the *Geocentric Cartesian* (or ECEF) coordinate system. We introduce this reference system in ns-3 via the `GeographicPositions` class, which provides methods to translate points represented using the coordinate system of [29] to/from those of [26].

Moreover, with usability in mind, we also implement a geographic coordinate system which allows ns-3 user to specify positions using the system of [26] in a more convenient manner. This auxiliary reference system represents positions as points exhibiting a relative altitude from their projection

¹<https://gitlab.com/mattiasandri/ns-3-ntn/-/tree/ntn-dev>

2 Simulation tools for future cellular networks

on the surface of the Earth. That is to say, any location on, or possibly above, Earth is referenced by a longitude ϕ , a latitude λ and an altitude h . To this end, we implement in the `GeographicPositions` class the methods `GeographicToTopocentricCoordinates` and `TopocentricToGeographicCoordinates`, which can be used to translate positions between geocentric and non-geocentric geographic coordinate systems.

For the conversion between any of the newly introduced models, and the cartesian reference system of [29], we introduce a *reference point* of translation between the two classes of coordinate systems, following the procedure outlined in [59, Ch. 4].

2.2.8.3 Channel condition

To model the channel condition for the NTN propagation scenarios, we create the classes:

- `ThreeGppNTNDenseUrbanChannelConditionMode`;
- `ThreeGppNTNUrbanChannelConditionMode`;
- `ThreeGppNTNSuburbanChannelConditionMode`; and
- `ThreeGppNTNRuralChannelConditionMode`.

Each of these derives from the base class `ThreeGppChannelConditionModel`, which in turn implements the `ChannelConditionModel` interface.

These channel condition classes interact with the remainder of the spectrum module as follows. Whenever the `GetChannelCondition` method is called, the newly introduced NTN `ChannelConditionModel` classes compute the channel state and cache it, along with its generation time. Then, the following calls to `GetChannelCondition` retrieve the previously stored value, if it has not expired. Otherwise, they compute a new LoS condition.

2.2.8.4 Path loss and shadowing

We implement the path loss and shadowing models of [26] in four different classes, as depicted in Figure 2.7. The latter extend the `ThreeGppPropagationLossModel` class, which in turn implements the `PropagationLossModel` interface.

The classes which implement this interface shall override the `DoCalcRxPower`, returning the received power based on the positions of the communicating endpoints, and when considering frequency-flat phenomena only. In the case of the NTN propagation scenarios of [26], these phenomena comprise the typical free space path loss, on top of tropospheric and ionospheric scintillation, shadow fading, clutter loss, and atmospheric absorption.

2.2.8.5 Geocentric mobility models

Along with the geographic coordinate systems, we implement a new mobility model, i.e., `GeocentricConstantPositionMobilityModel`, which allows ns-3 users to position nodes using real world coordinates. Specifically, the latter class stores positions via the variable `m_position`, which specifies their geographic coordinates.

When using these mobility models, the position of a node can be retrieved (set) using the methods `GetGeographicPostion` (`SetGeographicPostion`) and `GetGeocentricPosition` (`SetGeocentricPosition`). In turn, these methods rely on the functionality provided by the class `GeographicPosition` for translating between different coordinate systems.

Notably, the conversion from geocentric cartesian or geographic coordinates, to the coordinate systems used by ns-3, uses by default the so-called reference point “Null Island” (0,0,0). Nevertheless, ns-3 users are given the possibility of tuning this value by using the `GeocentricConstantPosition-MobilityModel` attribute `SetCoordinateTranslationReferencePoint`.

2.2.8.6 Antenna models

The circular aperture reflector antenna model currently implemented in ns-3, i.e., `ParabolicAntennaModel`, is based on a parabolic approximation of the main lobe radiation pattern, as described in [60] and [61]. This simplification reduces the computational complexity of the field pattern calculation, by avoiding the Bessel functions evaluations that the circular aperture antenna would require, and using trigonometric approximations instead.

As part of our contributions, we leverage the efficient implementation of the Bessel functions which has been introduced with C++17 to implement an exact circular aperture reflector antenna model. Specifically, we introduce this functionality extending the `AntennaModel` via the `CircularAperture-`

2 Simulation tools for future cellular networks

AntennaModel class. The latter allows ns-3 users to steer the pointing direction of the antenna via the SetOrientation and SetInclination methods. Similarly, the operating frequency and the aperture radius can be tuned by using the methods SetOperatingFrequency and SetApertureRadius.

2.2.9 Examples and Comparisons

In this section we validate the accuracy of our ns-3 module for the NTN channel, and compare simulation results with the calibration reported in TR 38.821 [62]. Furthermore, we provide numerical results to measure link-level and end-to-end performance (including throughput and packet drop ratio). We focus on satellites, even though the model is valid for different NTN scenarios.

2.2.9.1 Link-Level Results

While ns-3 enables system-level simulations, an evaluation of the link-level performance is still useful to validate the technical accuracy of our module. Hence, in this section we run link-level simulations to compare the calibration results from the 3GPP [62] with results from our module.

The 3GPP identifies 30 calibration study cases [62, Tab. 6.1.1.1-9], which include a combination of different satellite orbits, frequency bands, and antenna configurations for the ground terminal. Link-level calibration results are reported in [62, Tab. 6.1.1.2], including results for the Free Space Path Loss (FSPL), atmospheric loss (AL) and scintillation loss (SL), and the Carrier-to-Noise Ratio (CNR). Specifically, the CNR is calculated as described in [62, Section 6.1.3.1].

For this comparison we selected four study cases, considering both Up-link (UL) and Downlink (DL) transmissions, that illustrate four representative NTN scenarios. Specifically:

- Study Case 1 (SC1): GEO satellite, 45 degrees of elevation, VSAT antenna for the ground terminal, Ka-band.
- Study Case 6 (SC6): LEO satellite at 600 km, 90 degrees of elevation, VSAT antenna for the ground terminal, Ka-band.
- Study Case 9 (SC9): LEO satellite at 600 km, 90 degrees of elevation, UPA antenna for the ground terminal, S-band.

Table 2.2: Link-level comparison between the 3GPP calibration results (“3GPP”) and those obtained in simulations (“Obtained”). Header acronyms: Free Space Path Loss (FSPL), Atmospheric Loss (AL), Scintillation Loss (SL), Carrier-to-Noise Ratio (CNR). All values are in dB.

SC	Tx	Source	FSPL	AL	SL	CNR
1	DL	3GPP	210.6	1.2	1.1	11.6
		Obtained	210.6	1.4	1.1	11.3
1	UL	3GPP	214.1	1.1	1.1	0.5
		Obtained	214.2	1.4	1.1	0.1
6	DL	3GPP	179.1	0.5	0.3	8.5
		Obtained	179.9	0.5	0.3	8.6
6	UL	3GPP	182.6	0.5	0.3	18.4
		Obtained	182.6	0.5	0.3	18.4
9	DL	3GPP	159.1	0.1	2.2	6.6
		Obtained	159.1	0.0	2.2	6.7
9	UL	3GPP	159.1	0.1	2.2	2.8
		Obtained	159.1	0.0	2.2	2.4
14	DL	3GPP	164.5	0.1	2.2	7.2
		Obtained	164.5	0.0	2.2	7.3
14	UL	3GPP	164.5	0.1	2.2	-2.6
		Obtained	164.5	0.0	2.2	-3

- Study Case 14 (SC14): LEO satellite at 1200 km, 90 degrees of elevation, UPA antenna for the ground terminal, S-band.

The complete list of parameters used in the calibration can be found in [62, Section 6.1]. In Tab. 2.2 we report the calibration results (“3GPP”) and those from our simulations (“Obtained”). We can see that, despite some minor variations, numerical result are compatible under all metrics, thereby validating the accuracy of our module. As expected, the FSPL increases as the distance between the ground terminal and the satellite , as well as the carrier frequency, increase, due to the more severe effect of atmospheric losses. In particular, the impact of the carrier frequency is quite significant: for LEO satellites, for example, the FSPL grows from around 160 dB in the S-band (SC6) to around 180 dB in the Ka-band (SC9). In any case, we can see that, even considering long-range GEO satellites in the Ka-band, the CNR is large enough to support adequate levels of communication, especially in downlink. We shed light on two main trends. First, uplink communication is generally worse than downlink, except for SC6. This is reasonable, and due to the

Table 2.3: Simulation parameters.

Parameter	Value
Frequency	20 GHz ÷ 100 GHz
Satellite orbit	GEO
Satellite altitude	35 786 km
Elevation angle	90 deg
Tx. mode	Downlink
Transmit power	37.5 dBm
Satellite antenna	Circular aperture (Gain: 58.5 dB)
Terminal antenna	VSAT (Gain: 39.7 dB)
Scenario	Suburban

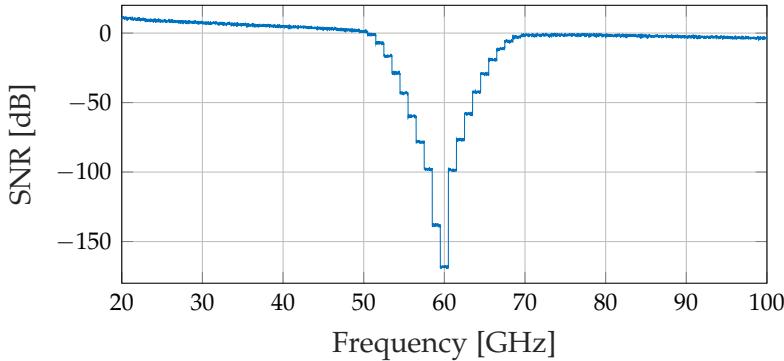


Figure 2.8: The SNR for different carrier frequencies. We consider a GEO satellite, with the parameters in Tab. 2.3.

fact that ground terminals are more constrained in terms of power availability, capacity, and size (e.g., for antenna deployment). Second, according to the 3GPP model, LEO satellites have more severe hardware constraints than GEO satellites (e.g., LEO's effective isotropic radiated power (EIRP) in the Ka-band is as low as 36 dBW, vs. 66 dBW for GEO): as a result, the CNR for SC6 is around 50% lower than for SC1, which makes LEO communication more difficult.

2.2.9.2 Frequency Test

While the 3GPP identifies the S-band (at 2 GHz) and the Ka-band (at 20 GHz for DL and 30 GHz for UL) as frequencies of interests, the NTN channel model is valid for a wide range of frequencies, from 0.5 GHz to 100 GHz. Therefore, in Figure 2.8 we plot the Signal-to-Noise Ratio (SNR), which is an

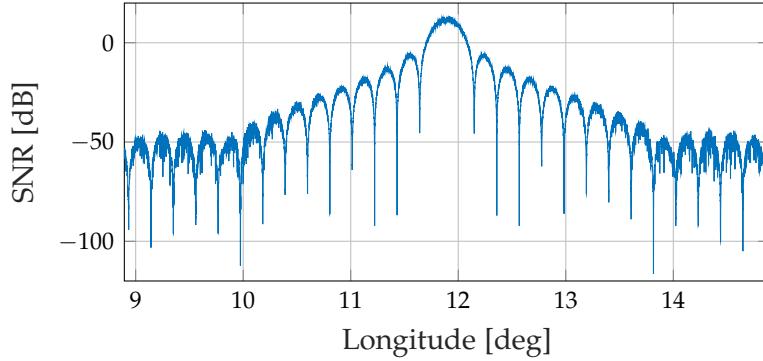


Figure 2.9: The SNR for different angular positions of a GEO satellite, with the parameters in Tab. 2.3.

indication of the quality of the channel, as a function of the carrier frequency, with a resolution of 8 MHz. The other parameters are summarized in Tab. 2.3.

We observe that the SNR decreases linearly (in the log scale) as the frequency increases, with a deep degradation at 60 GHz. This is because of the impact of atmospheric absorptions described in Section 2.2.3, more specifically the additional signal attenuation experienced at 60 GHz due to oxygen absorption (as large as 15 dB/km).

2.2.9.3 Mobility Test

Our new mobility model for NTN (see Section 2.2.8.5) allows to change the position of the nodes during the simulation, thus to evaluate the effect of different parameters such as the elevation angle, the antenna radiation pattern, and the altitude.

First, we run simulations where we iteratively change the coordinates of a satellite, which traces an arc of 6 degrees in the GEO orbit (from 8.8 to 14.8 degrees of longitude). The receiving node on the ground is deployed so that it is perpendicular to the satellite in the mid point of its trajectory. The rest of the parameters are set as in Tab. 2.3. Notably, the orientation of the antenna is not changed during the simulation, so that satellite and ground terminal are perfectly aligned only when the former is exactly perpendicular to the latter. In Figure 2.9 we plot the corresponding SNR, which resembles the circular antenna radiation pattern of the satellite as per Eq. (2.23), which therefore defines the power profile of the received signal.

Table 2.4: End-to-end performance. Acronyms: Orbit Type (OT), Transmit power (TxP), Throughput (TP), Drop Rate (DR), Frequency Band (FB), Terminal antenna (TA).

SC	1	6	9	14
OT	GEO	LEO (600 km)		LEO (1200 km)
TxP	37.52 dBm	21.52 dBm	48.77 dBm	54.77 dBm
FB	Ka-band	Ka-band	S-band	S-band
TA	VSAT	VSAT	UPA	UPA
TP	3.811 Mbit/s	3.286 Mbit/s	4.101 Mbit/s	5.161 Mbit/s
DR	0.61	0.67	0.45	0.36

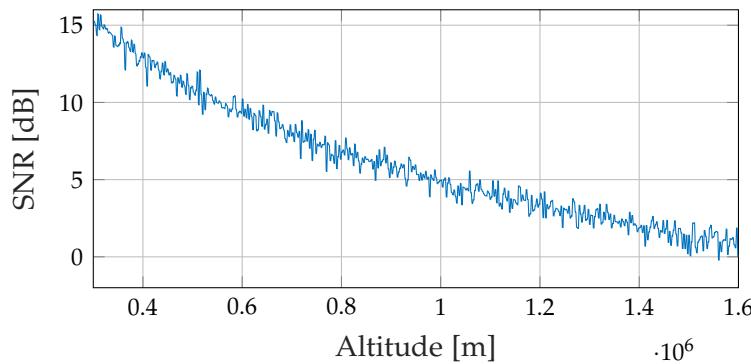


Figure 2.10: The SNR for different values of the altitude of a LEO satellite, with the parameters in Tab. 2.3.

Second, in Figure 2.10 we plot the SNR for different values of the altitude of the satellite, from 300 to 1600 km to consider different LEO satellite architectures. We see that the SNR decreases as the altitude of the satellite increases, even though it is consistently above 0 dB in all configurations.

2.2.9.4 End-to-End Performance

Unlike other simulators, ns-3 incorporates an accurate model of the whole ISO/OSI protocol stack, thus enabling scalable end-to-end simulations. Notably, end-to-end results can be collected to validate and measure the performance of communication networks, so ns-3 stands out as a valid tool to dimension NTN systems too. To do so, we consider a downlink application transmitting data (in the form of UDP packets) at a constant rate of 10 Mbit/s. We test the four study cases presented in Section 2.2.9.1, so to consider both GEO and LEO satellites, and different altitudes, antenna con-

2.2 Modeling non-terrestrial wireless channels

figurations, and both communication in the Ka- and S-band. Simulations results are reported in Tab. 2.4 in terms of end-to-end throughput (TP) and packet drop rate (DR), which is defined as the ratio between the number of received packets and the total number of packets sent at the application layer. We observe that both GEO-to-ground and LEO-to-ground communication is feasible, provided that the satellite operates with large-scale antennas offering fine-grained beams on the ground, and at high transmit power. Notice that the DR is quite significant, especially in the Ka-band, which requires the design of appropriate Automatic Repeat reQuest (ARQ) schemes to deal with retransmissions.

2.3 *Improving the scalability of wireless channel simulation in ns-3*

Channel models range from simple models that just consider a propagation loss component combined with Nakagami-m or Rayleigh fading but fail to capture the spatial dimension of the channel and the interactions with beam-forming [63], to deterministic models that are very accurate in specific scenarios but are much more complex and require a precise characterization of the environment [64]. To address the complexity-accuracy trade-off, the 3GPP has adopted a stochastic channel model for simulations of 5G and beyond networks [29]. Stochastic channel models are generic, thanks to their stochastic nature, and can model interactions with multiple-antenna arrays. The latter was included in ns-3 thanks to the efforts of Tommaso Zugno in the 2019 Google Summer of Code [28], and later extended to address vehicular scenarios in [65] and industrial scenarios in [66]. As a consequence, the current spatial channel model implemented in ns-3 is very accurate for simulations in line with 3GPP specifications for a wide range of frequencies. However, it represents the main bottleneck in terms of computational complexity when considering large-scale simulations with many multi-antenna nodes, especially when equipped with large antenna arrays. This is because of the intrinsic complexity in the generation of the channel model according to 3GPP specifications, and the need to deal with inefficient tensor structures. In fact, the channel matrix in the ns-3 implementation of the 3GPP spatial channel model is currently implemented as a 3D structure made of nested vectors, whose dimensions depend on the number of the transmit antennas, receive antennas, and clusters.

The design of computationally efficient yet accurate channel models has been a topic of interest also in the Wireless Local Area Network (WLAN) space. The authors of [67, 68] present a frequency-selective channel for WLANs, and use Exponential Effective SNR Mapping (EESM) L2SM to integrate their model with the ns-3 system-level Wi-Fi implementation. Moreover, they develop a framework which leverages cached statistical channel matrix realizations to directly estimate the effective SNR, thus further improving the computational efficiency of the model. Specifically, the latter is modeled as a parameterized log-SGN random variable. They extend their work in [69], by accounting for the channel correlation over time. More-

2.3 Improving the scalability of wireless channel simulation in ns-3

over, [70] compares statistical channel models for the 60 GHz band with the Quasi Deterministic (QD) Ray Tracer (RT) of [71].

In the remainder of this section, we summarize the efforts carried out in the 2022 Google Summer of Code to further optimize the code in ns-3 in two directions: 1) improving the efficiency of the code by allowing the use of the Eigen library, and 2) proposing a new performance-oriented MIMO channel model for reduced complexity in ns-3 large-scale simulations.

2.3.1 Efficient MIMO modeling with the Eigen library

The use of multiple antennas both at the transmitter and at the receiver, a fundamental feature of modern wireless systems, makes a scalar representation of the channel impulse response insufficient. Instead, MIMO channels are usually represented in the form of a complex matrix $\mathbf{H} \in \mathbb{C}^{U \times S}$, whose elements depict the channel impulse response between the U and S radiating elements of the transmitting and receiving antenna arrays, respectively [29]. This peculiarity significantly increases the computational complexity of MIMO channel models, compared to Single Input Single Output (SISO) ones, since the complex gain of the channel must be evaluated for each pair of transmit and receive antennas. Notably, previous analyses identified in statistical channel models the main bottleneck for system-level MIMO wireless simulations. In typical m-MIMO 5G scenarios, where the devices feature a high number of antennas, the channel matrix generation and the computation of the beamforming gain represent up to 90% of the simulation time [24].

In light of these limitations, as the first of our contributions, we optimized the implementation of the 3GPP TR 38.901 model in ns-3 introduced in [28]. First, we observed that, as of ns-3.37, part of the trigonometric operations of the `GetNewChannel` method of the `ThreeGppChannelModel` class are unnecessarily repeated for each pair of transmitting and receiving radiating elements. This represents a significant inefficiency, since the inputs of these functions, i.e., the angular parameters of the propagation clusters, depend on the cluster index only. Moreover, the standard library `sin` and `cos` functions are particularly demanding to evaluate. Therefore, we cached the trigonometric evaluations of these terms prior to the computation of \mathbf{H} 's coefficients, effectively

reducing the complexity of the trigonometric operations from $\mathcal{O}(U \times S \times N)$ to $\mathcal{O}(N)$, where N is the number of propagation clusters.

Then, we focused on improving the algebra manipulations of the channel matrix performed in the `ThreeGppSpectrumPropagationLossModel` by introducing the support for the open-source library `Eigen` in `ns-3`. `Eigen` is a linear algebra C++ template library that offers fast routines for algebra primitives such as matrix multiplication, decomposition and space transformation [30], and is used by many open-source frameworks such as `TensorFlow`.

We set `Eigen` as an optional, external `ns-3` dependency, with the goal of minimizing future code maintenance efforts, and thus mimicking the support for other third-party libraries. To get `Eigen`, `ns-3` users can either rely on packet managers, i.e., install the package `libeigen3-dev` (`eigen`) for Linux (Mac) systems, or manually install the library by following the official instructions². Then, `Eigen` can be enabled via a custom flag defined in the `macros-and-definitions.cmake` file, and its presence in the system is shown to the user by exposing whether it has been found or not via the `ns3--config-table.cmake` file. The latter also defines the preprocessor definition `HAVE_EIGEN3`, which is used in the `ns-3` source files to discern `Eigen`'s availability. Finally, the linking of `Eigen` with the `ns-3` source files is taken care of by the `CMake` configuration file provided by the library itself, as suggested in the related `ns-3` guide.

To prevent the need for `Eigen` to be installed in the host system, we developed a common set of APIs between the `Eigen`- and the Standard Template Library (STL)-based data structures and primitives. Thanks to this choice, the remainder of the `spectrum` code is completely abstracted with respect to the presence of the library. Given that most of the needed operators can not be overloaded for STL C++ vectors (for instance, `operator()`), the common interface for both `Eigen` and STL's based vectors and matrices has been implemented by defining ad hoc structs with custom operators. In particular, we defined:

- The complex vector type `PhasedArrayModel::ComplexVector`. This data-structure is defined as an `std::vector` of `std::complex<double>` whenever `Eigen` is not installed, and as an `Eigen` vector of `std::complex<double>` otherwise. The set of APIs includes operators `[]` and `!=`, which can be

²<https://gitlab.com/libeigen/eigen/-/blob/master/INSTALL>

used to access the vector entries and to compare pairs of vectors, respectively. Additionally, we defined the STL-like methods `size`, `norm` and `resize`, which return the vector size, its \mathcal{L}^2 -norm, and allow the user to resize the underlying container, respectively. These definitions follow the typical STL notation, as it is supported by Eigen as well.

- The complex matrix type `MatrixBasedChannelModel::Complex2DVector`. In this case, the underlying type is a nested `std::vector` of `std::complex<double>` for when Eigen is disabled, and an Eigen matrix whose entries are of type `std::complex<double>` otherwise.

In this case, we aligned the notation to the APIs provided by Eigen. Specifically, the matrix elements can be accessed via the operator `()`, which takes as arguments the row and column indices of the entry, while the method `resize` allows users to resize matrices by specifying the number of rows and columns. In turn, these can be accessed via the `rows` and `columns` methods, respectively.

- The 3D matrix `MatrixBasedChannelModel::Complex3DVector`. This data structure is defined, regardless of Eigen's availability, as an `std::vector` of `MatrixBasedChannelModel::Complex2DVector`. In this case, the only method provided is `MultiplyMatByLeftAndRightVec`, which computes a product of the type $\mathbf{w}_T \mathbf{H} \mathbf{w}_R^T$, where $\mathbf{H} \in \mathbb{C}^{U \times S}$, $\mathbf{w}_T \in \mathbb{C}^{1 \times U}$ and $\mathbf{w}_R \in \mathbb{C}^{1 \times S}$. Notably, this computationally demanding evaluation, which is required for computing the beamforming gain in `ThreeGppSpectrumPropagationLossModel`, leverages Eigen's optimized algorithms whenever the library is installed in the host system.

Finally, we remark that the support for Eigen in the ns-3 codebase can possibly be further extended to improve the efficiency of other linear algebra operations, such as the Singular Value Decomposition (SVD) which is used in the `mmwave` and `nr` modules to compute optimal beamformers, and the matrix-by-matrix multiplications needed for relayed channels [72].

2.3.2 A performance-oriented MIMO statistical channel model

The second approach to reduce computational complexity we propose in this section is a MIMO channel model for simulating large m-MIMO scenarios,

implemented in the class `TwoRaySpectrumPropagationLossModel`. The goal of this auxiliary model is to offer a faster, albeit slightly less accurate, statistical channel model than the 3GPP TR 38.901 framework of [28] by preventing the need for the computation of the complete channel matrix. In line with [29], the frequency range of applicability of this model is 0.5 – 100 GHz, although the framework can be possibly extended to support higher frequencies as well.

The overall channel model design follows the approach of [73], i.e., the end-to-end channel gain is computed by combining several loss and gain terms which account for both large- and small-scale propagation phenomena, and the antenna and beamforming gains. In particular, let T be a device transmitting a signal x with power P_T^x , and R be another device in the simulation (which may or may not be the intended destination of x). The proposed model implements the `PhasedArraySpectrumPropagationLossModel` interface by estimating P_R^x , i.e., the power of x received at R , as follows:

$$\begin{aligned} P_R^x[dBm] &= P_T^x[dBm] - PL_{T,R}[dB] \\ &\quad + S_{T,R}[dB] + G_{T,R}[dB] + F_{T,R}[dB], \end{aligned} \tag{2.24}$$

where the terms $PL_{T,R}$ and $S_{T,R}$ represent the path loss and the shadowing, respectively, while $G_{T,R}$ and $F_{T,R}$ denote the antenna and beamforming gain and the small-scale fading, respectively. The remainder of this section describes in detail how each of these terms is computed.

2.3.3 Path loss, Shadowing, and LoS Condition

The large-scale propagation phenomena are modeled according to the 3GPP TR 38.901 model [29], since its implementation of [28] is not computationally demanding. Nevertheless, the channel model can in principle be coupled with arbitrary classes extending the `ChannelConditionModel` interface.

Specifically, we first determine the 3GPP scenario. Then, for each link we set the LoS condition in a stochastic manner, using the class extending `ThreeGppChannelConditionModel` which corresponds to the chosen scenario.

Then, we compute the path loss using the 3GPP TR 38.901 formula

$$PL_{T,R} = A \log_{10}(d) + B + C \log_{10}(f_C)[dB], \tag{2.25}$$

2.3 Improving the scalability of wireless channel simulation in ns-3

where d is the 3D distance between the transmitter and the receiver, f_C is the carrier frequency, and A, B and C are model parameters which depend on the specific scenario and the LoS condition.

To account for the presence of blockages, an optional log-normal shadowing component $S_{T,R}$ and an outdoor-to-indoor penetration loss term are added to $PL_{T,R}$.

2.3.3.1 Antenna and Beamforming Gain

The combined array and beamforming gain is computed using the approach of [74]. The proposed model supports the presence of multiple antenna elements at the transmitter and at the receiver, and arbitrary analog beamforming vectors and antenna radiation patterns. Therefore, ns-3 users can use this model in conjunction with any class that implements the `AntennaModel` interface. In this implementation, we focus on UPAs, although the methodology is general and can be applied to arbitrary antenna arrays.

Let θ and φ be the relative zenith and azimuth angles between transmitter and receiver, respectively, and let $\mathbf{w}(\theta_0, \varphi_0)$ denote the beamforming vector pointing towards the steering direction (θ_0, φ_0) . We denote with $U = U_h U_v$ the total, horizontal, and vertical number of antenna elements, respectively, and with d_h, d_v their spacing in the horizontal and vertical domains of the array, respectively.

Considering first isotropic antennas, the gain pattern of a UPA, in terms of received power relative to a single radiating element, can be expressed as [75]

$$G_{T,R}^{iso}(\theta, \varphi) = \left| \mathbf{a}_i^T(\theta, \varphi) \mathbf{w}(\theta_0, \varphi_0) \right|^2, \quad (2.26)$$

where $\mathbf{a}_i(\theta, \varphi)$ is the array response vector, whose generic entry m, n with $m \in \{0, \dots, U_v - 1\}, n \in \{0, \dots, U_h - 1\}$ reads

$$a_i(\theta, \varphi)_{m,n} = \exp\left(j \frac{2\pi}{\lambda} m d_v \cos(\theta)\right) \exp\left(j \frac{2\pi}{\lambda} n d_h \sin(\theta) \sin(\varphi)\right).$$

In this work, which supports arbitrary antennas, each antenna element (m, n) actually exhibits a generic radiation pattern $g(\theta, \varphi)_{m,n}$ towards direction (θ, φ) . In particular, we assume that $g(\theta, \varphi)_{m,n}$ is constant for all the elements of the array, i.e., $g(\theta, \varphi)_{m,n} \equiv g(\theta, \varphi)$. Accordingly, we compute

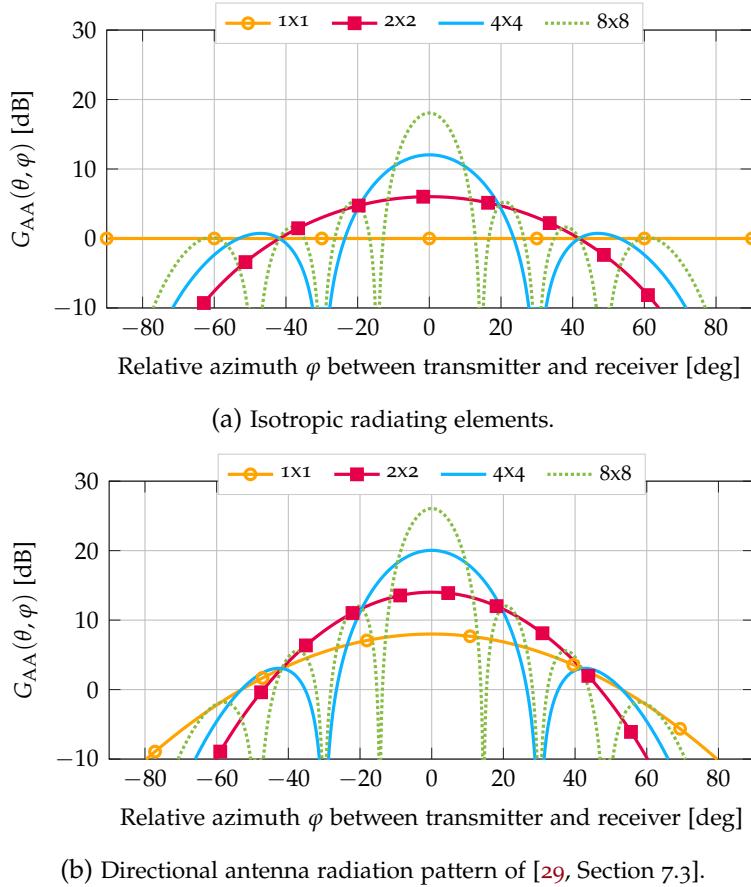


Figure 2.11: Overall array and beamforming gain of a UPA, for isotropic and 3GPP [29, Section 7.3] radiating elements and {1x1, 2x2, 4x4, 8x8} antenna configurations. The steering direction is fixed to $(\theta_0, \varphi_0) = (0^\circ, 0^\circ)$, and $\theta \equiv 0^\circ$.

$G_{T,R}(\theta, \varphi)$ in the `ComputeBeamformingGain` function of the `TwoRaySpectrum-PropagationLossModel` class as

$$G_{T,R}(\theta, \varphi) = G_{T,R}^{iso}(\theta, \varphi) |g(\theta, \varphi)|^2. \quad (2.27)$$

Figures 2.11a and 2.11b report $G_{T,R}(\theta, \varphi)$ for both the isotropic (`IsotropicAntennaModel`) and the 3GPP (`ThreeGppAntennaModel`) radiation patterns, respectively.

It can be noted that our model abstracts the computation of the received signal power as a SISO keyhole channel [76], which is then combined with the spatial antenna gain patterns at the transmitter/receiver to obtain the

2.3 Improving the scalability of wireless channel simulation in ns-3

received power. This approximation is possibly imprecise when considering NLoS links, due to the lack of a dominant multipath component. To account for this limitation, we introduce a multiplicative correction factor η which scales the beamforming gain as $G'_{T,R}(\theta, \varphi) \equiv \eta G_{T,R}(\theta, \varphi)$. In line with [77], we set $\eta = 1/19$.

2.3.3.2 Fast Fading

The widely used Rayleigh and Rician distributions fail, even in their generalized forms, to capture the intrinsic bimodality exhibited by mmWave scenarios [78–80]. Therefore, in our implementation we model fast fading using the more general Fluctuating Two-Ray (FTR) model of [81]. This fading model assumes that the received signal comprises two dominant specular components and a mixture of scattered paths, thus modeling the amplitude of the received signal V_r as

$$V_r = V_1 \sqrt{\xi} \exp(j\phi_1) + V_2 \sqrt{\xi} \exp(j\phi_2) + X + jY, \quad (2.28)$$

where ϕ_1, ϕ_2 are statistically independent random phases, distributed as $\phi_i \sim \mathcal{U}[0, 2\pi]$. X and Y are independent Gaussian random variables, i.e., $X, Y \sim \mathcal{N}(0, \sigma^2)$, which represent the diffuse component of the received signal, which is assumed to be the superposition of multiple weak scattered waves with independent phase. Finally, ξ is a unit-mean Gamma distributed random variable with rate m and Probability Density Function (PDF)

$$f_\xi(u) = \frac{m^m u^{m-1}}{\Gamma(m)} \exp(-mu). \quad (2.29)$$

In our implementation, $F_{T,R} = |V_r|^2$ is sampled via the `GetFtrFastFading` function of the `TwoRaySpectrumPropagationLossModel` class.

The FTR fading model is usually expressed as a function of the Gamma rate m and the auxiliary parameters

$$K \doteq \frac{V_1^2 + V_2^2}{2\sigma^2} \quad (2.30)$$

$$\Delta \doteq \frac{2V_1 V_2}{V_1^2 + V_2^2} \in [0, 1], \quad (2.31)$$

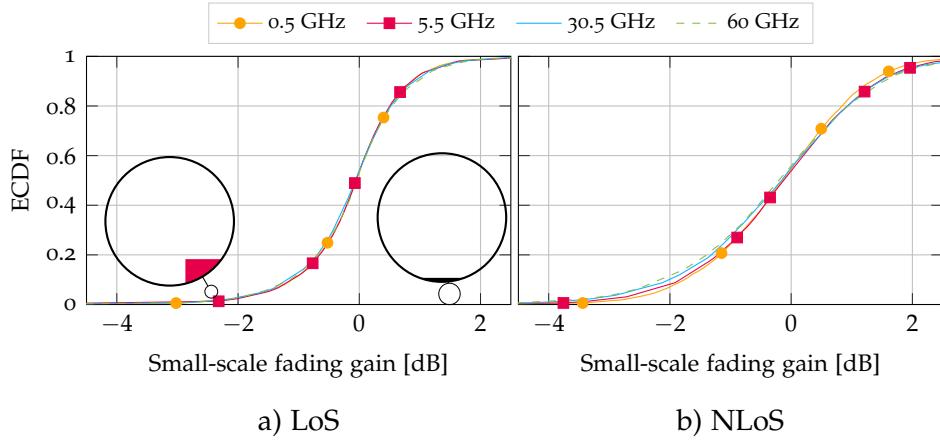


Figure 2.12: Small-scale fading gain statistics for the UMi propagation scenario versus the carrier frequency f_C , for both LoS and NLoS channel conditions.

where K represents the ratio of the power of the specular components with respect to the diffuse ones, while Δ denotes how similar the received powers of the specular components are. By tuning these parameters, a high degree of flexibility can be achieved. Notably, a choice of $\Delta = 0$ effectively yields a Rician-distributed signal amplitude [81].

Calibration. In our work, we calibrated the V_1, V_2 and m parameters of the FTR fading model using the full 3GPP TR 38.901 channel model as a reference. In particular, we first obtained the statistics of the small-scale fading of the 3GPP model, using an ad hoc calibration script (`three-gpp-two-ray--channel-calibration.cc`). The script produces a collection of channel gain samples obtained by using the `ThreeGppSpectrumPropagationLossModel` and the `ThreeGppChannelModel` classes, and neglecting the beamforming gain, path-loss, shadowing and blockages. Accordingly, we isolate the variation around the mean received power caused by the small-scale fading only. A separate set of these samples has been retrieved for both LoS and NLoS channel conditions, the different propagation scenarios of [29], and a set of carrier frequencies ranging from 0.5 to 100 GHz. However, a preliminary evaluation of the obtained data showed a negligible dependence of the small-scale fading with respect to the carrier frequency, as can be observed in Figure 2.12. Therefore, we calibrated the FTR parameters considering only the channel condition and the propagation scenario.

The small-scale fading samples have been used to estimate the Δ , K and m FTR parameters, and then derive analytically the values of V_1 and V_2 yielding the fading realizations that are the closest (in a goodness-of-fit sense) to the TR 38.901 model. To this end, we defined a discrete grid of FTR parameters, spanning their whole domain, and considered the corresponding set of parameterized FTR distributions. To find the best matching one, we measured the distance between each of these distributions and the 3GPP reference curves by using the Anderson-Darling goodness-of-fit test [82]. This test is used to discern whether a sorted collection of n samples $\{Y_1, \dots, Y_n\}$ originates from a specific distribution, by evaluating the test statistic [82]

$$A^2 = -n - S(\mathcal{F}), \quad (2.32)$$

where

$$S(\mathcal{F}) = \sum_{i=1}^n \frac{2i-1}{n} [\ln(\mathcal{F}(Y_i)) + \ln(1 - (\mathcal{F}(Y_{n+1-i}))), \quad (2.33)$$

and \mathcal{F} is the Cumulative Distribution Function (CDF) of the target distribution. In the standard Anderson-Darling test, A^2 is then compared to a pre-defined critical value to validate the hypothesis. Instead, in our work we find the FTR distribution $\mathcal{F}_{m,K,\Delta}$ which yields the lowest S . Specifically, for each combination of propagation scenario, LoS condition and corresponding samples $\{Y_1, \dots, Y_n\}$ we find

$$\mathcal{F}_{m^*,K^*,\Delta^*} \doteq \operatorname{argmin}_{m,K,\Delta} S(\mathcal{F}_{m,K,\Delta}). \quad (2.34)$$

Finally, we exported the calibrated FTR parameters into ns-3, by storing them in `SIM_PARAMS_TO_FTR_PARAMS_TABLE`, i.e., an `std::map` which associates the propagation scenario and condition to the corresponding best fitting FTR parameters. We remark that this calibration process represents a pre-computation step which needs to be done only once. Indeed, when running a simulation with this channel model, the FTR parameters get simply retrieved from the pre-computed lookup table by the `GetFtrParameters` function. Nevertheless, for the sake of reproducibility and maintainability of the code, we provide this functionality in the Python script `two-ray-to-three--gpp-ch-calibration.py`.

2 Simulation tools for future cellular networks

2.3.4 Benchmarks, examples and use cases

In this section, we provide an example on how to use the performance-oriented channel model presented above, in conjunction with the New Radio (NR) [18] module, to simulate 5G MIMO networks. Moreover, we present benchmarks which quantify the simulation time reduction achieved with this work, and we outline some possible use cases.

2.3.4.1 Examples and Benchmarks

We demonstrate how to use the performance-oriented channel model in the `cttc-nr-demo-two-ray` script, i.e., a custom version of the `cttc-nr-demo` example which is included in the NR module. The script deploys N_{gNB} 5G NR base stations, along with N_{UE} users in each cell. Each UE uploads data using two Bandwidth Parts (BWP) operating at 28 and 30 GHz, respectively. Both base stations and user terminals feature UPAs with multiple radiating elements.

Most simulation parameters can be tuned by ns-3 users. Notably, the script provides the possibility to choose whether to use the 3GPP TR 38.901 channel model of [28] or the FTR-based channel model proposed in this work. In such regard, the use of the `TwoRaySpectrumPropagationLossModel`, instead of the TR 38.901 one, is achieved by:

1. Setting the `TypeId` of the `SpectrumPropagationLossModel` factory to `TwoRaySpectrumPropagationLossModel`;
2. Creating an instance of the `TwoRaySpectrumPropagationLossModel` class using the above factory, and setting the corresponding pointer as the `SpectrumPropagationLossModel` of both BWPs;
3. Setting the attribute `Frequency` of the `TwoRaySpectrumPropagationLossModel` instance as the BWP carrier frequency;
4. Specifying the 3GPP propagation scenario by setting the attribute `Scenario` and
5. Creating and setting the `ChannelConditionModel` by using the `TwoRaySpectrumPropagationLossModel` class `ChannelConditionModel` attribute.

2.3 Improving the scalability of wireless channel simulation in ns-3

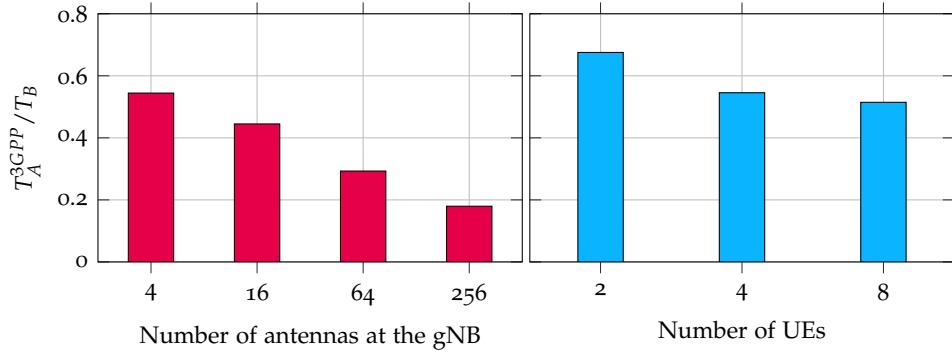


Figure 2.13: Ratio of the median simulation times after the merge of this work with the Eigen integration (T_A^{3GPP}) and as per ns-3.37 (T_B), when using the 3GPP channel model of [29].

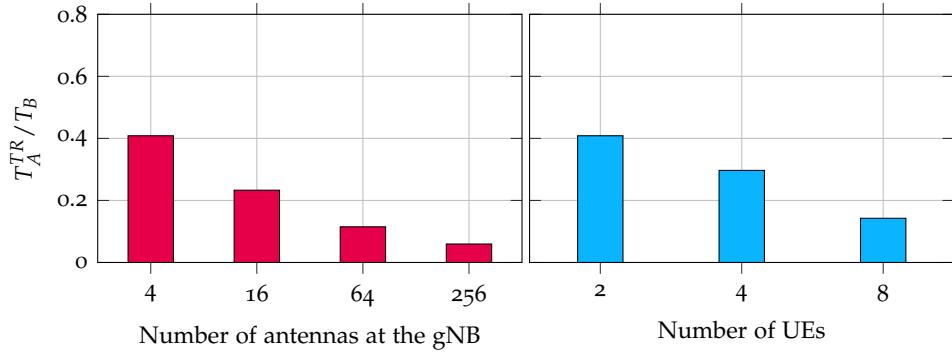


Figure 2.14: Ratio of the median simulation times using the performance-oriented channel model presented in this work (T_A^{TR}) and the 3GPP channel model of [29] after the merge of this work. In this case, Eigen is disabled.

On the other hand, the Eigen optimizations simply require users to have the corresponding library installed in their system, and to enable Eigen when configuring ns-3, using the flag `enable-eigen`.

We validated our contributions by benchmarking the simulation times exhibited by the above simulation script, which depicts a typical MIMO 5G NR scenario. To such end, we varied the number of gNB antennas and UEs deployed, and we timed 100 simulation runs for each parameter combination. Figure 2.13 reports the ratio of the median simulation time achieved when using the Eigen-based optimizations, and of the same metric obtained using the vanilla ns-3.37. It can be seen that the matrix multiplication routines offered by Eigen can significantly reduce simulation times. For instance, a reduction of 5 times in the simulation time is achieved when equipping gNBs

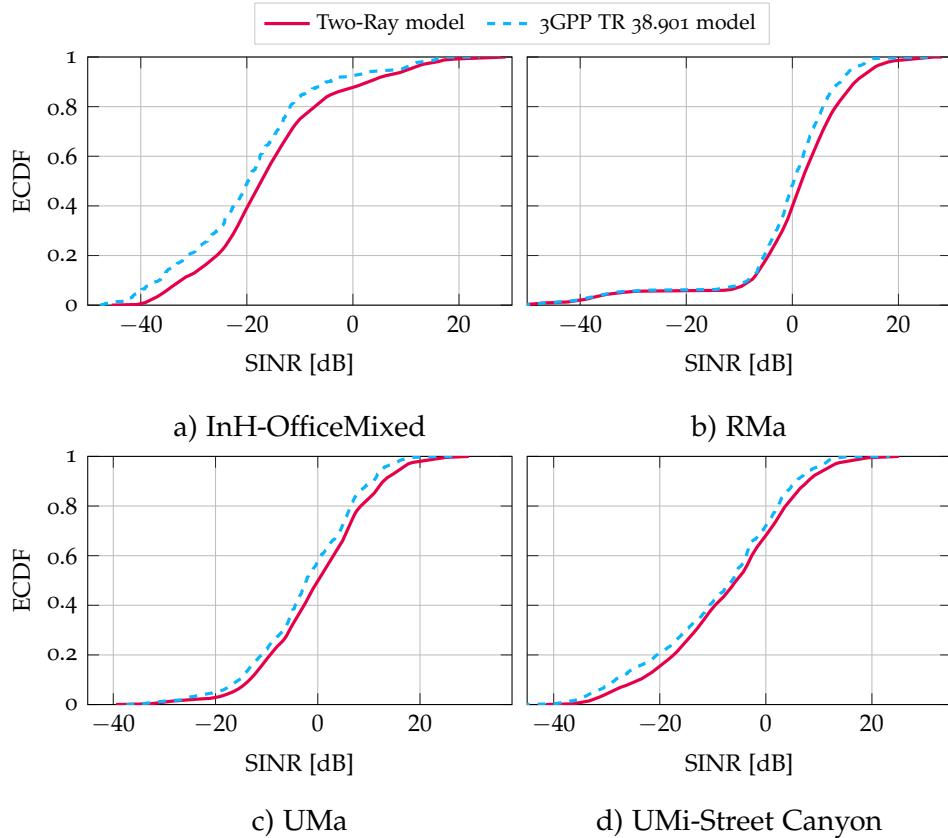


Figure 2.15: ECDF of the SINR obtained using the 3GPP channel model of [29], and the performance-oriented channel model presented in this work, for different propagation scenarios.

with 256 radiating elements. Similarly, Figure 2.14 depicts the ratio of the median simulation time obtained by using the FTR-based channel model, and the 3GPP TR 38.901 with Eigen disabled. In this case the computational complexity improvement is even more dramatic, with simulations taking as low as 6 % of the time to complete, with respect to the 3GPP model implementation of [28]. As a reference, the median simulation time obtained on an Intel[®] i5-6700 processor system, before the merge of this work and for $\{2, 4, 8\}$ users is $\{64.7, 210.5, 666.6\}$ [s], respectively.

Finally, we also computed (using the same simulation script, i.e., `cttc--nr-demo-two-ray`) the SINR statistics achieved by the proposed FTR-based model, and compared them to those obtained using the model of [28]. As can be seen in Figure 2.15, the two models provide similar results. Indeed, a non-

2.4 Conclusions and future work

negligible difference can be found only in the case of the InH-OfficeMixed propagation scenario.

We remark that all the results presented in this section can be reproduced by using the SEM [83] scripts which we provide³.

2.3.4.2 Use Cases

The main goal of both the performance oriented channel model and the optimizations to the 3GPP TR 38.901 model is to enable system-level simulations of large-scale MIMO scenarios for which the implementation of [28] exhibits prohibitive computational complexity. Specifically, our contributions allow ns-3 users to simulate wireless deployments where the devices feature antenna arrays with more than hundreds of radiating elements, and/or the number of communication endpoints is particularly high. For example, the modifications presented in this work can be used in the NR and mmwave [17] modules (which both already support the proposed channel models) to simulate massive MIMO 5G NR networks. Notably, a preliminary version of the Eigen port has been used in conjunction with the mmwave [17] module to simulate 5G networks aided by IRSs, i.e., devices which feature up to 100×100 reflecting elements [84].

Moreover, since the supported frequency range is 0.5 – 100 GHz, this encompasses not only terrestrial 5G and Long Term Evolution (LTE) deployments, but also most non-terrestrial networks and IEEE Radio Access Technologies (RATs). Finally, the proposed TwoRaySpectrumPropagationLossModel can be further extended to support frequencies above 100 GHz using reference fading and path loss statistics.

2.4 Conclusions and future work

In this chapter, we proposed a signal model for IRSs and AF relays based on the 3GPP TR 38.901 channel for 5G NR networks, and explained the methodology we used to perform network-level simulations of 5G and beyond scenarios with IRS and AF relay nodes. Based on this framework, we performed simulations to provide numerical guidelines to dimension IRS/AF-assisted

³<https://gitlab.com/pagmatt/ns-3-dev/-/tree/gsoc-wns3>

2 Simulation tools for future cellular networks

networks. Moreover, we presented an ns-3 implementation of the 3GPP channel model for NTN, developed following the specifications provided in [26]. The code, which is publicly available at [85], well integrates with the rest of the ns-3 framework, and enables full-stack end-to-end simulations in different NTN scenarios. We validated the link-level and end-to-end accuracy of our module against 3GPP calibration results reported in [62]. Finally, we presented a set of optimizations concerning the simulation of MIMO wireless channels in ns-3. These improvements comprise the optimization of the related linear algebra routines, and the design and implementation in ns-3 of a performance-oriented statistical channel model based on the FTR fading model, which further reduces the simulation time of MIMO scenarios.

As part of our future work, we plan to extend our smart relays simulator by considering more sophisticated scenarios in which heterogeneous types of relays are deployed, and compare the numerical performance of IRS/AF relays with that of IAB. Moreover, we will also relax some of our assumptions by considering quantization of the relay phase shifters.

Furthermore, we foresee to further extend our NTN module to simulate end-to-end NTN 5G NR networks. To this end, we will incorporate additional functionalities, such as a delay model, along with the adaptations to the terrestrial 5G NR protocol stack which it entails, and the support for satellite mobility.

Additionally, we plan to further improve the scalability of the wireless channel simulation framework in ns-3 by studying more refined beamforming gain correction factors, and possibly making the estimation of such term scenario-dependent. Moreover, we envision to design more efficient storage/access data structures and linear algebra operations for 3D matrices, by better leveraging Eigen also in this context. Finally, we will consider using Single Instruction, Multiple Data (SIMD) for speeding up the evaluation of trigonometric functions, and caching the beamforming gain in the `TwoRay-SpectrumPropagationLossModel` class to further reduce the simulation time of MIMO scenarios in ns-3.

3 *Towards wireless-backhauled next-generation cellular networks*

Future wireless networks will accommodate data-rate intensive use cases which include untethered Virtual Reality (VR) and mobile metaverse applications. This will further exacerbate the congestion of mobile access networks and backhaul systems [86]. To accommodate this traffic increase, the 3GPP has introduced various technological advancements with the specifications of the 5G Radio Access Network (RAN) and Core Network (CN), namely NR and 5G Core (5GC) [87]. In particular, NR features a user and control plane split, a flexible Orthogonal Frequency Division Multiplexing (OFDM) frame structure, and the support for mmWave communications, while the CN introduces virtualization and slicing [88].

Notably, the use of the mmWave band, with typical deployments in the spectrum around 28 GHz and 39 GHz [89], possibly coupled with sub-terahertz mobile links [8, 90], represents the major technological enabler toward the Gbit/s capacity target. Indeed, these frequencies are characterized by the availability of vast chunks of contiguous and currently unused spectrum, in stark contrast with the crowded sub-6 GHz bands. However, mmWaves and terahertz frequencies exhibit unfavorable propagation characteristics, such as high isotropic losses and a marked susceptibility to blockages and signal attenuation [91, 92]. These issues can be partially mitigated using beamforming through large antenna arrays, thanks to the small wavelengths and advances in low-power Complementary Metal-Oxide Semiconductor (CMOS) RF circuits [93]; nevertheless, their introduction alone is not enough for meeting the high service availability requirement. In fact, wireless networks operating at such high frequencies will be deployed with extremely high density, to improve the probability of LoS coverage and mitigate the impact of the harsh propagation environment. Nonetheless, while the theoretical effectiveness of this technique is well understood [94], achieving dense cellular deployments

is extremely challenging from a practical point of view. Specifically, providing a fiber backhaul among base stations and the CN is deemed economically impractical, even more so in the initial 5G deployments [95].

To make ultra-dense deployments viable, the 3GPP has standardized an extension of 5G NR, i.e., IAB, which exploits the same waveform and protocol stack to provide access to mobile users and wireless backhaul for gNBs (i.e., the IAB nodes) thus limiting the need for fiber drops. The wireless backhaul topology terminates at a gNB with fiber connectivity to the data core, the IAB donor [96–98]. IAB also simplifies the deployment of cellular networks in on-demand or ad hoc contexts, as it removes the need for part of the wired backhaul. Prior research has highlighted that IAB represents a cost-performance trade-off [95, 97], as base stations need to multiplex access and backhaul resources, and as the wireless backhaul at mmWaves is less reliable than a fiber connection. In particular, IAB networks may suffer from excessive buffering (and, consequently, high latency and low throughput) when a suboptimal partition of access and backhaul resources is selected, thus hampering the benefits that the high bandwidth mmWave links introduce [45, 95]. Therefore, it is fundamental to solve these non-trivial challenges to enable a smooth integration of IAB in 5G and beyond deployments.

In this chapter, we introduce several solutions for optimizing routing and backhaul/access resource partitioning in IAB networks. In particular, Section 3.1 describes a semi-centralized resource allocation scheme for IAB networks, designed to be flexible, with low complexity, and compliant with the 3GPP IAB specifications. The proposed solution, which is based on the Maximum Weighted Matching (MWM) problem, is compared with state-of-the-art distributed approaches through end-to-end, full-stack system-level simulations with a 3GPP-compliant channel model, protocol stack, and a diverse set of user applications. Results show that this scheme can increase the throughput of cell-edge users up to 5 times, while decreasing the overall network congestion with an end-to-end delay reduction of up to 25 times. Section 3.2 describes Safehaul, a risk-averse learning-based solution for IAB mmWave networks. Instead of optimizing the average latency performance, Safehaul ensures reliability by minimizing the losses in the tail of the performance distribution. We show via extensive simulations that Safehaul not only reduces the latency by up to 43.2% compared to the benchmarks, but also exhibits significantly more reliable performance, e.g., 71.4% less vari-

ance in latency. Finally, in Section 3.3 we consider the deployment of mixed mmWave and sub-terahertz links to increase the capacity of the backhaul network, and provide the first performance evaluation of the potential of sub-terahertz frequencies for 6G IAB. To do so, we develop a greedy algorithm that allocates frequency bands to the backhaul links (considering constraints on spectrum licenses, sharing, and congestion) and generates the wireless network mesh. Then, we profile the performance through a custom extension of the open-source system-level simulator Sionna that supports Release 17 IAB specifications and channel models up to 140 GHz. Results show that IAB with sub-terahertz links can outperform a mmWave-only deployment with improvements of $4\times$ for average user throughput and a reduction of up to 50% for median latency.

3.1 Semi-centralized framework for resource management in 5G NR Integrated Access and Backhaul

The literature adopts different approaches to model and solve the resource allocation problem in multi-hop wireless networks. The first, discussed in [99–105] is based on conventional optimization techniques. Specifically, the authors of [99] present a simple and thus tractable system model and find the minimal number of gNBs featuring a wired backhaul that are needed to sustain a given traffic load. Their work is further extended in [100], which provides an analysis of the performance benefits introduced by additional, fiber-less gNBs. In [101], the mobile network is modeled as a noise-limited, k -ring deployment. Such model is then used to obtain closed-form expressions for the max-min rates achieved by UEs in the network. Moreover, [102] proposes a system model which leads to an NP-hard optimization problem, even though it considers single-hop backhaul networks only, and uses deep Reinforcement Learning (RL) to reduce its computation complexity. In [103], the joint routing and resource allocation problem is tackled via a Linear Programming (LP) technique. Notably, this work assumes that data can be transmitted (received) toward (from) multiple nodes at the same time. Similarly, the authors of [104] formulate a Time Division Duplexing (TDD), multi-hop resource allocation optimization problem which leverages the directionality of mmWave antennas, albeit in the context of Wireless Personal Area Networks (WPANs). Since such problem is also NP-hard, a sub-optimum solution is found. Finally, [105] focuses on joint link scheduling, routing and power allocation in multi-hop wireless networks. As in previous cases the obtained optimization problem is not tractable: in this instance such obstacle is overcome by studying the dual problem via an iterative approach.

The second approach relies on stochastic geometry to model IAB networks [97, 106]. Specifically, [106] determines the rate coverage probability of IAB networks and compares different access/backhaul resource partitioning strategies. Similarly, [97] provides a comparison of orthogonal and integrated resource allocation policies, although limited to single-hop wireless networks.

Another significant body of literature leverages Markov Chains (MCs) to study IAB networks; some of these works can be interpreted as a direct application of such theory [107, 108], while others [109–112] exploit a more

complex framework. The papers which belong to the former class are based on the pioneering work of [113], which inspects the stability of generic multi-hop wireless networks and formulates a throughput-maximizing algorithm known as *back-pressure*. In particular, [107] focuses on the optimization of the timely-throughput, i.e., takes into account that packets usually have an arrival deadline. Such problem is then addressed by formulating a Markov Decision Process (MDP), leading to a distributed resource allocation algorithm. Similarly, [108] proposes an algorithm that also targets throughput optimality but, contrary to the back-pressure algorithm, manages to avoid the need for per-flow information. On the other hand, the body of literature which belongs to the latter class uses the MC-derived Network Utility Maximization (NUM) framework first introduced in [114] and [115]. Specifically, the authors of [109] focus on satisfying the URLLC Quality of Service (QoS) requirements by jointly optimizing routing and resource allocation. Then, the problem is solved using both convex optimization and RL techniques. In [110], an in-depth analysis of a mmWave, multi-hop wireless system is presented, proposing and comparing three different interference frameworks, under the assumption of a dynamic TDD system. This work is extended in [111] and [112], which consider respectively a Spatial Division Multiple Access (SDMA) and a Multi-User (MU)-MIMO capable system.

Finally, only a small portion of the literature [45, 95, 116] analyzes the end-to-end performance of IAB networks. Specifically, the authors of [45] extend the ns-3 mmWave module, introducing realistic IAB functionalities which are then used to characterize the benefit of deploying wireless relays in mmWave networks. Their work is extended in [116], where path selection policies are formulated and their impact on the system performance is inspected. A further end-to-end analysis of IAB networks is carried out in [95], providing insights into the potentials of this technology and the related open research challenges.

Concluding, the literature exhibits the presence of algorithms relying on a varying degree of assumptions on the network topology and the knowledge of system. Furthermore, most of the aforementioned studies lack an end-to-end, full-stack system-level analysis of the proposed solution. To fill these gaps, this section proposes a semi-centralized resource allocation scheme which exhibits low complexity, both computationally and in terms of required feedback. Moreover, we provide considerations on how our proposed

solution can be implemented and deployed in standard-compliant 3GPP IAB networks, and compare such solution to the state of the art with an end-to-end, realistic performance analysis

3.1.1 Contributions

This remainder of this section tackles the access and backhaul partitioning problem by proposing an optimal, semi-centralized resource allocation scheme for 3GPP IAB networks, based on the MWM problem on graphs. It receives periodic L₁ and/or L₃ measurements from the nodes of the IAB deployment, a possibility which is explicitly mentioned by 3GPP in [117, Section 7.3.3], constructs a spanning tree that represents the deployment, and uses a simplified, low-complexity version of the MWM to partition the links between access and backhaul. After a feedback step, each node can then schedule the resources at a subframe-level among the connected devices.

To the best of our knowledge, this is the first MWM-based resource allocation framework for 3GPP IAB networks at mmWaves. As such, it exhibits the following benefits: (i) no constraints on the number of hops in the IAB-network are introduced, and, more in general, it is 3GPP-compliant; (ii) a globally optimum is computed; (iii) generic network utility functions can be used; (iv) it features a computational complexity which is linear in the number of gNBs which are connected to the same IAB-donor; and (v) a very limited communication overhead is required.

In particular, the flexibility makes it possible to easily adapt the resource allocation strategy to different requirements, use cases, and classes of traffic for 5G networks. We achieve this by developing a generic optimization algorithm, which identifies with a configurable periodicity the access and backhaul partition that optimizes a certain utility function. The selection of the utility function prioritizes the optimization of different metrics, e.g., throughput or latency, which in turn can be mapped to different classes of traffic.

Moreover, to achieve the compliance with the 3GPP IAB specifications, the resource allocation framework relies only on information that can be actually exchanged and reported in a 3GPP deployment. In this regard, we also review the latest updates related to the 3GPP IAB standardization activities. Nevertheless, our solution can be easily extended to consider other types of

feedback information. Finally, the algorithm operates with a low complexity, i.e., we propose a version of the MWM algorithm that can be applied on spanning trees with linear complexity in the number of nodes in the network infrastructure, and demonstrate its equivalence to the generic (and more complex) MWM. Additionally, the proposed framework also relies on a feedback exchange that is linear in the number of base stations, and is thus decoupled from the number of users. Along this line, the semi-centralized nature of the proposed solution combines the benefit of a centralized point of view for the allocation of inter-dependent IAB links and a limited complexity.

Furthermore, we evaluate the performance of the proposed scheme with an end-to-end, full-stack system-level simulation, using the ns-3 mmWave module [17] and its IAB extension [45]. This represents the first evaluation of an optimized resource allocation scheme for IAB with a simulator that is based on a 3GPP-compliant protocol stack, uses 3GPP channel models, and integrates realistic applications and transport protocols. The extended performance evaluation highlights how the proposed scheme improves the throughput of a diverse set of applications, with a 5-fold increase for the worst case users, with different packet sizes and transport protocols, while decreasing the latency and buffering at intermediate nodes by up to 25 times for the smallest packet sizes.

The remainder of this section is organized as follows. Section 3.1.2 describes our assumptions and the system model. Then, Section 3.1.3 presents a novel scheme for resource partitioning in mmWave IAB networks, along with considerations on how it can be implemented in 3GPP NR. Finally, Section 3.1.6 describes the performance evaluation results.

3.1.2 IAB networks

The following paragraphs identify the characteristics and constraints of mmWave IAB, according to the 3GPP design guidelines presented in [117] and the specifications of [98].

3.1.2.1 Network topology

In general, an IAB network is a deployment where a percentage of gNBs (i.e., the IAB-nodes) use wireless backhaul connections to connect to a few gNBs (i.e., the IAB-donors) which feature a wired connection to the core network,

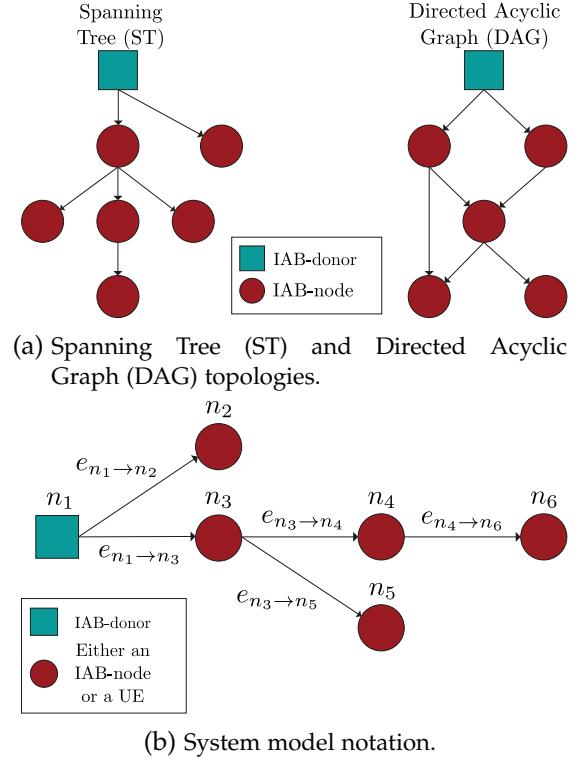


Figure 3.1: Comparison of the IAB network topologies analyzed in [117] and related notation.

as can be seen in Fig. 3.1. Moreover, these deployments exhibit a *multi-hop* topology where a strict parent-child relation is present. The former can be represented by the IAB-donor itself or an IAB-node; the latter by either UEs or downstream IAB-nodes. In [117], no a priori limit on the number of backhaul hops is introduced. As a consequence, 3GPP argues that IAB protocols should provide sufficient flexibility with respect to the number of backhaul hops. Moreover, the Study Item (SI) on IAB [117] highlights the support for both the topologies depicted in Fig. 3.1a, i.e., ST and DAG IAB. Clearly, the former exhibits less complexity but, at the same time, poses some limits in terms of network performance: the possible presence of obstacles may result in a service interruption, due to the unique backhaul route established by the UEs. On the other hand, a DAG topology offers routing redundancy, which can be used not only to decrease the probability of experiencing a “topological blockage,” but also for load balancing purposes.

3.1.2.2 Multiple access schemes and scheduling

An in-band, dynamic partitioning of the access and backhaul spectrum resources is currently preferred by 3GPP [98, 117], together with half-duplex operations of the IAB-nodes. Moreover, most of the literature suggests that 5G mmWave systems will operate in a TDD fashion [91, 118]. This choice is mainly driven by the stringent latency requirements which the next generation of mobile networks will be required to support, and by the usage of analog or hybrid beamforming. The usage of Frequency Division Duplexing (FDD), in conjunction with the presence of large chunks of bandwidth, would lead to severe resource under-utilization and make channel estimation more difficult. Based on these considerations, the system model exhibits a TDD, TDMA-based scheduling where the access/backhaul interfaces are multiplexed in a half-duplex manner. Coupled with mmWaves directionality, this means that self and inter-cell-interference are both limited, as reported by [99]. Furthermore, at any given time instant, each node of the IAB network cannot be simultaneously involved in more than one transmission or reception. In particular, IAB-nodes cannot schedule time and frequency resources which are already allocated by their parent for backhaul communications which involve them. Moreover, the backhaul links of a given gNB might also carry data which is destined to (and/or generated by) UEs which are connected to different base stations. As a consequence, an IAB-network exhibits a marked and peculiar inter-dependence between the resource allocations of the various base stations, which is the major motivation for the introduction of a semi-centralized framework.

Finally, the introduction of resource coordination mechanisms and related signaling is explicitly supported in the IAB specification drafts [98, 117]. Nevertheless, these solutions must reuse as much as possible the available NR specifications and require at most minimal changes to the Rel.15 5GC and NR.

3.1.2.3 System model

According to these assumptions and referring to Fig. 3.1b., a generic IAB network can be modeled as a directed graph $\mathcal{G} = \{\mathcal{N}, \mathcal{E}\}$, where the set of nodes $\mathcal{N} \triangleq \{n_1, n_2, \dots, n_{|\mathcal{N}|}\}$ comprises the IAB-donor, the various IAB-nodes and the UEs. Accordingly, the set of directed edges $\mathcal{E} \triangleq \{e_1, e_2, \dots\}$

$e_{|\mathcal{E}|} \} \equiv \{e_{n_j \rightarrow n_k}\}_{j,k}$, where the edge $e_{n_j \rightarrow n_k}$ originates at the parent node n_j and terminates at the children n_k , comprises in all the active cell attachments, either of mobile terminals to a gNB or from IAB-nodes towards their parent node. Since the goal of this section is to study backhaul/access resource partitioning policies, this generic model can be actually simplified: in fact, all the UEs connected to a given gNB can be represented by a single node in \mathcal{G} without any loss of generality. Similarly, the same holds true for their links toward the serving gNB, which can then be represented by a single edge. Furthermore, this work focuses on ST topologies only.

We define as *feasible schedule* any set of links $\mathcal{E}' \subseteq \mathcal{E}$ such that none of them share a common vertex, i.e., $\forall e_{n_j \rightarrow n_k} \neq e_{n_l \rightarrow n_m} \in \mathcal{E}'$ it holds that $n_j \neq n_m$ and $n_l \neq n_k$. Let then f_u be a utility *additive map*, namely, a function such that the overall utility experienced by the system when scheduling edges e_1 and e_2 satisfies $f_u(e_1, e_2) = f_u(e_1) + f_u(e_2)$. Let also $\mathcal{W} \stackrel{\Delta}{=} \{w_1, w_2, \dots, w_{|\mathcal{E}|}\}$ be the set of positive weights whose generic entry w_j represents the utility which is obtained when scheduling the j -th edge, namely, $w_j \stackrel{\Delta}{=} f_u(e_j)$. Then, the overall utility of the system is $\mathcal{U} \stackrel{\Delta}{=} \sum_{e_k \in \mathcal{E}'} f_u(e_k) = \sum_{e_k \in \mathcal{E}'} w_k$. The goal is to find the feasible set \mathcal{E}'^* which maximizes the overall utility, i.e., $\underset{\mathcal{E}'}{\operatorname{argmax}} \mathcal{U}$. In computer science, this task is typically referred to as the *Maximum Weighted Matching* problem [119].

Finding the MWM of a given graph, in the general case, is not trivial from a computational point of view. In fact, the fastest known MWM algorithm for generic graphs has a complexity of $\mathcal{O}(|V||E| + |V|^2 \log |V|)$ [120], posing serious limitations to the suitability of such algorithm to 5G and beyond networks, which target a connection density of 1 million devices per km². However, we argue that under the aforementioned assumptions on the system model, which restrict the network to an ST topology, it is possible to design an MWM-based semi-centralized resource partitioning framework which exhibits linear complexity with respect to the network size and which, as a result, is able to satisfy the scalability requirements highlighted by 3GPP in [117]. Nevertheless, the proposed framework can be easily extended to the case of a DAG IAB network. In such regard, a sub-optimal strategy is to periodically discard, during each centralized allocation, the redundant edges of each node. In such a way, the input which is fed to the T-MWM algorithm is, effectively, an ST. A second, optimal extension can be obtained by computing

at the controller the MWM of the network via a generic MWM algorithm, instead of using the ST-specific T-MWM as in the proposed framework. However, this strategy would feature a higher computational complexity.

3.1.3 Semi-centralized resource allocation scheme for IAB networks

This section presents an MWM algorithm for ST topologies (Section 3.1.3.1), an efficient and MWM-based semi-centralized resource partitioning framework for IAB networks (Section 3.1.4) and some considerations about its implementation (Section 3.1.5). Specifically, the proposed scheme collects at a controller installed on the IAB-donor L₁ and/or L₃ measurements from the various gNBs. Then, it uses such information to build a weighted ST which represents the IAB-network. In particular, the network topology is inferred by examining the incoming parent-child associations. The edge weights are also computed from the received measurements, based on the specific policy (hence, of target Key Performance Indicators (KPIs)) of choice. Finally, the resource partitioning is optimized by computing an MWM of the network and then prioritizing the links which comprise it. A high level diagram of is provided in Fig. 3.2b.

3.1.3.1 MWM for ST graphs

As the first of our contributions, we present an algorithm, hereby called T-MWM, which computes the MWM of an ST in linear time. In particular, T-MWM is a bottom-up algorithm which, upon receiving as input a weighted ST \mathcal{G} described by its edge map \mathbf{E} and the corresponding weight map \mathbf{W} , produces as output a set of active edges \mathbf{E}^* which are an MWM of \mathcal{G} . That is to say, \mathbf{E}^* is a matching of \mathcal{G} which yields the globally maximum utility. Furthermore, \mathbf{E} is from now on assumed to exhibit the following invariant: each IAB parent precedes its children in the map, hence avoiding the need for a recursion. This is automatically obtained as each IAB child connects after its parent, and is thus added to the map in a subsequent position. Nevertheless, this assumption can be easily relaxed, albeit at the cost of losing as a side-effect the bottom-up design.

The proposed algorithm is designed starting from the observation that, given the generic node $n_k \in \mathcal{G}$ and a matching $\bar{\mathbf{E}}$ of \mathcal{G} , we can identify the following mutually exclusive and collectively exhaustive cases: $\bar{\mathbf{E}}$ can contain

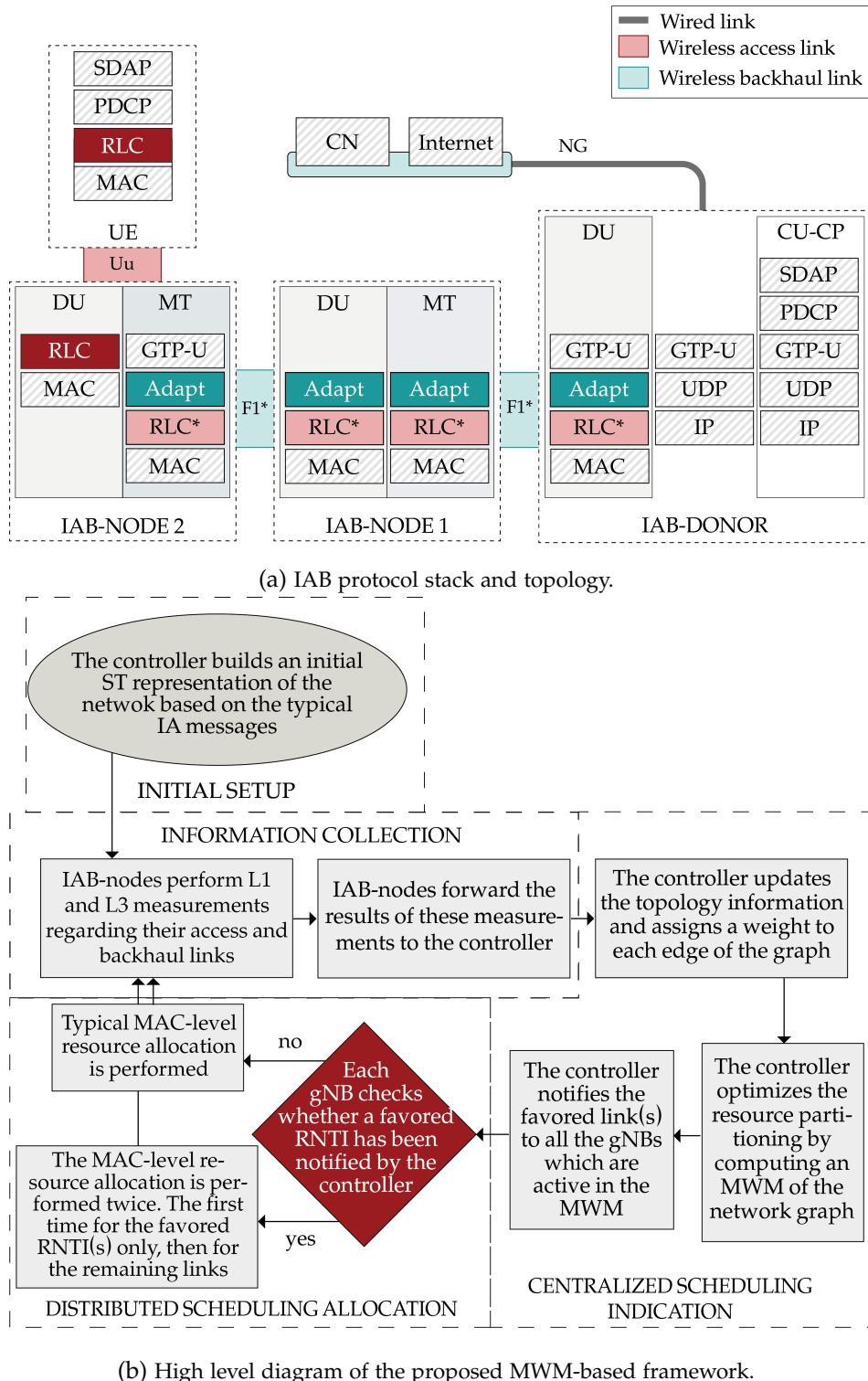


Figure 3.2: IAB topology and proposed MWM-based framework.

Algorithm 1 Tree-Maximum Weighted Matching

Input: A weighted ST \mathcal{G} encoded by a map \mathbf{E} , which associates each node in \mathcal{G} to its edges, and the corresponding weights map \mathbf{W} .

Output: An MWM \mathbf{E}^* of \mathcal{G} .

```

1: procedure T-MWM( $\mathbf{E}, \mathbf{W}$ )
2:    $\mathbf{F} \leftarrow \mathbf{0}; \mathbf{G} \leftarrow \mathbf{0}$                                  $\triangleright$  Initialize the utility vectors to zero vectors
3:    $\mathbf{E}^* \leftarrow \{\}$                                           $\triangleright$  Initialize the set of active edges as empty
4:   for each internal node  $n_k \in \mathbf{E}$  do           $\triangleright$  In ascending order w.r.t. to their
      depth in  $\mathcal{G}$ 
5:      $maxUtil \leftarrow -\infty; \mathbf{maxUtilChild}(n_k) \leftarrow \{\}$ 
6:     for each edge  $e_{n_k \rightarrow n_j} \in \mathbf{E}(n_k)$  do            $\triangleright$  Iterate over its edges
7:        $\mathbf{G}(n_k) \leftarrow \mathbf{G}(n_k) + \max\{\mathbf{F}(n_j), \mathbf{G}(n_j)\}$ 
8:        $currUtil \leftarrow \mathbf{W}(e_{n_k \rightarrow n_j}) - [\mathbf{F}(n_j) - \mathbf{G}(n_j)]^+$ 
9:       if  $currUtil > maxUtil$  then
10:         $maxUtil \leftarrow currUtil; \mathbf{maxUtilChild}(n_k) \leftarrow n_j$ 
11:       end if
12:     end for
13:      $\mathbf{F}(n_k) \leftarrow \mathbf{G}(n_k) + maxUtil$ 
14:   end for
15:   for each internal node  $n_k \in \mathbf{E}$  do           $\triangleright$  In ascending order w.r.t. to their
      depth in  $\mathcal{G}$ 
16:     if  $\mathbf{F}(n_k) \geq \mathbf{G}(n_k)$  then
17:        $\mathbf{E}^* \leftarrow \mathbf{E}^* \cup e_{n_k \rightarrow \mathbf{maxUtilChild}(n_k)}$ 
18:        $\mathbf{F}(\mathbf{maxUtilChild}(n_k)) \leftarrow -\infty$             $\triangleright$  Ensure child does not
          get activated multiple
          times
19:     end if
20:   end for
21:   return  $\mathbf{E}^*$ 
22: end procedure

```

either one or zero edges which originate from n_k . Based on this fact, we then discern the optimal utilities which can be obtained in each of these cases. Specifically, we define the maximum utilities yielded by a matching of n_k 's sub-tree which either contains a link originating from n_k or not as $\mathbf{F}(n_k)$ and $\mathbf{G}(n_k)$, respectively. Then, as can be seen in Alg. 1, the T-MWM algorithm basically consists in two traversals of the network graph. During the first one we compute the \mathbf{G} and \mathbf{F} functions for all the nodes in \mathcal{G} using the recursive formulas provided by Lemma 1. Finally, during the second traversal, this knowledge is used for computing an MWM of the network; the correctness of this last phase is proved by Lemma 2.

Lemma 1. *Given an ST \mathcal{G} , consider its generic internal node n_k . Let then $\mathbf{F}(n_k)$ be the maximum utility yielded by a matching of n_k 's sub-tree which activates a link*

originating from n_k , and $\mathbf{G}(n_k)$, conversely, the utility provided when such matching contains no links which feature n_k as parent. Then, we have that:

$$\begin{cases} \mathbf{G}(n_k) = \sum_{\{n_j\}_k} \max\{\mathbf{F}(n_j), \mathbf{G}(n_j)\} \\ \mathbf{F}(n_k) = \mathbf{G}(n_k) + \max_{\{n_j\}_k} \{\mathbf{W}(e_{n_k \rightarrow n_j}) \\ \quad - [\mathbf{F}(n_j) - \mathbf{G}(n_j)]^+\} \end{cases}$$

where the set $\{n_j\}_k$ comprises all the children of n_k and $[x]^+ = \max\{x, 0\}$ is the positive part of x . Conversely, for leaf nodes n_l it holds that $\mathbf{F}(n_l) \equiv \mathbf{G}(n_l) \equiv 0$.

Proof. This lemma can be proved by induction over the height h_k of the sub-tree corresponding to node n_k . The base case is $h_k = 0$, i.e., when n_k is a leaf node; in this case, trivially, both $\mathbf{F}(n_k)$ and $\mathbf{G}(n_k)$ are zero since no links exhibit n_k as parent node and the sub-tree of \mathcal{G} which originates in n_k consists of n_k only, respectively.

Then, assume that n_k 's sub-tree exhibits a generic height $h_k > 0$, and that the above formulas hold for each of its children sub-trees, which exhibit a height $h_j < h_k$. If we do not activate any edge which originates from n_k , then no added constraints are introduced concerning the edges which can be activated in its children sub-trees. Therefore, $\mathbf{G}(n_k)$ is simply the sum of the utilities achieved by any MWM computed on its children sub-trees, i.e., $\mathbf{G}(n_k) = \sum_{\{n_j\}_k} \max\{\mathbf{F}(n_j), \mathbf{G}(n_j)\}$. The remaining option is to activate exactly one edge, hereby called $e_{n_k \rightarrow n_m}$, which originates from n_k . In this case, no additional edges which feature n_m as parent can be added to the matching. As a consequence, the contribution of n_m 's sub-tree on $\mathbf{F}(n_k)$ reads $\mathbf{G}(n_m)$. Conversely, no additional constraints are introduced regarding the other nodes. It follows that the utility obtained in this instance reads:

$$\sum_{\{n_j \neq n_m\}_k} \max\{\mathbf{F}(n_j), \mathbf{G}(n_j)\} + \mathbf{W}(e_{n_k \rightarrow n_m}) + \mathbf{G}(n_m)$$

and can be rewritten as:

$$\mathbf{G}(n_k) + \mathbf{W}(e_{n_k \rightarrow n_m}) - [\mathbf{F}(n_m) - \mathbf{G}(n_m)]^+$$

Finally, such utility is clearly maximized when n_m is chosen as $\underset{\{n_j\}_k}{\operatorname{argmax}} \{ \mathbf{W}(e_{n_k \rightarrow n_j}) - [\mathbf{F}(n_j) - \mathbf{G}(n_j)]^+ \}$, yielding:

$$\mathbf{F}(n_k) = \mathbf{G}(n_k) + \underset{\{n_j\}_k}{\max} \{ \mathbf{W}(e_{n_k \rightarrow n_j}) - [\mathbf{F}(n_j) - \mathbf{G}(n_j)]^+ \} \quad \blacksquare$$

Lemma 2. Given an ST \mathcal{G} of root n_r and the \mathbf{F} and \mathbf{G} functions computed as per Lemma 1, an MWM \mathbf{E}^* of \mathcal{G} can be computed by performing the following procedure:

1. If $\mathbf{F}(n_r) \geq \mathbf{G}(n_r)$, add to \mathbf{E}^* the edge from n_r to n_m , where the latter is defined as $n_m \stackrel{\Delta}{=} \underset{\{n_j\}_r}{\operatorname{argmax}} \{ \mathbf{W}(e_{n_r \rightarrow n_j}) - [\mathbf{F}(n_j) - \mathbf{G}(n_j)]^+ \}$. Then, repeat recursively on all the sub-trees corresponding to n_r 's children $\{n_j\}_r | n_j \neq n_m$ and on the children of n_m itself.
2. If $\mathbf{F}(n_r) < \mathbf{G}(n_r)$, repeat recursively on all the sub-trees corresponding to n_r 's children.

Proof. The above procedure always yields a feasible activation, i.e., a matching of \mathcal{G} . In particular, in either options we never recurse on a node which has already been activated, hence no pair of edges $\in \mathbf{E}^*$ can share any vertices. Furthermore, due to the properties of \mathbf{F} and \mathbf{G} , whenever $\mathbf{F}(n_r) \geq \mathbf{G}(n_r)$ a matching yielding maximal utility can be obtained by activating the edge $e_{n_r \rightarrow n_m}$, where $n_m \stackrel{\Delta}{=} \underset{\{n_j\}_r}{\operatorname{argmax}} \{ \mathbf{W}(e_{n_r \rightarrow n_j}) - [\mathbf{F}(n_j) - \mathbf{G}(n_j)]^+ \}$. Since the procedure is then recursively repeated on n_r 's children and the validity of \mathbf{F} and \mathbf{G} properties holds for each sub-tree in \mathcal{G} , the set of edges \mathbf{E}^* produced by the above procedure comprise a *maximal* matching, i.e., they yield the maximum possible utility among all the feasible schedules. \blacksquare

Regarding the computational complexity of the proposed algorithm, it can be observed that during the first phase the main loop effectively scans each edge of \mathcal{G} , hence exhibiting a complexity $\mathcal{O}(|\mathbf{E}|)$. Moreover, the second phase of T-MWM has complexity $\mathcal{O}(|\mathbf{V}|)$, since it loops through all the network nodes. Therefore, we can conclude that the overall asymptotic complexity of the algorithm is $\mathcal{O}(|\mathbf{V}| + |\mathbf{E}|)$, or, equivalently, $\mathcal{O}(|\mathbf{V}|)$ since in an ST the number of edges equals $|\mathbf{V}| - 1$.

3.1.4 Semi-centralized resource partitioning scheme

Based on the system model introduced in Section 3.1.2, and the T-MWM algorithm, we present a generic optimization framework which partially centralizes the backhaul/access resource partitioning process, in compliance with the guidelines of [117]. The goal of this framework is to aid the distributed schedulers, adapting the number of OFDM symbols allocated to the backhaul and access interfaces to the phenomena which exhibit a sufficiently slow evolution over time, i.e., large scale fading and local congestion. This optimization is undertaken with respect to a generic additive utility function f_u . An IAB network of arbitrary size is considered, composed of a single IAB-donor, multiple IAB-nodes and a (possibly time-varying) number of UEs which connect to both types of gNBs. Furthermore, assume that a central controller is installed on the IAB-donor.

The proposed framework can be subdivided into the following phases, which are periodically repeated every T_{alloc} subframes:

1. **Initial setup.** This step, which is depicted in Fig. 3.3a, consists in the computation of the simplified IAB network graph $\mathcal{G} \equiv \{\mathcal{V}, \mathcal{E}\}$. Specifically, after this phase \mathcal{V} comprises the donor and the various IAB-nodes. Accordingly, \mathcal{E} contains their active cell associations.
2. **Information collection.** During this phase, the various IAB-nodes send to the central controller a pre-established set of information for each of their children in \mathcal{G} . For instance, this feedback may consist in their congestion status and/or information regarding their channel quality. To such end, the implementation presented in this section uses modified versions of pre-existing NR Release 16 Control Elements (CEs), as strongly recommended in the IAB SI [117]. However, the scheme does not actually impose any limitations in such regard.
3. **Centralized scheduling indication.** Upon reception of the feedback information, the central controller updates \mathcal{G} by inspecting the received node-parent associations. Then, the set of weights \mathcal{W} is calculated and an MWM of \mathcal{G} is computed, using the T-MWM algorithm. The output of this procedure is the activation set \mathbf{E}^* , which yields a globally optimum solution with respect to the chosen utility function. Subsequently, \mathbf{E}^* is used as to create a set of *favored* downstream nodes, i.e., of children

which will be served with the highest priority by their parent, as depicted in Fig. 3.3b. Finally, these scheduling indications are forwarded to the various IAB-nodes which act as parents in the edges of E^* .

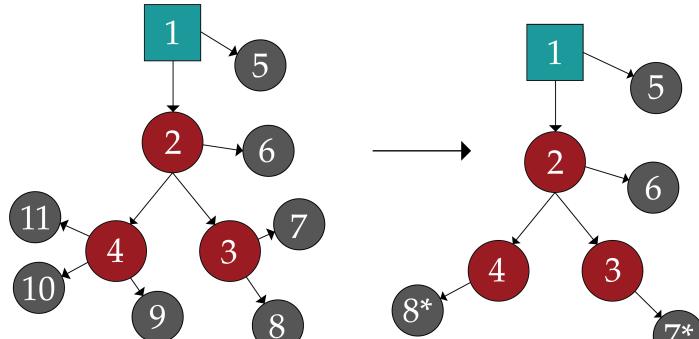
4. **Distributed scheduling allocation.** During this phase, the various IAB-nodes make use of the indications received by the central controller, if available, in order to perform the actual scheduling (which is, therefore, predominantly distributed). Specifically, the favored nodes are served with the highest priority, while the remaining downstream nodes are scheduled if and only if the resource allocation of the former does not exhaust the available OFDM symbols.

It is important to note that since \mathcal{G} contains only the IAB-nodes, the donor and at most one “representative” UE per gNB, the proposed scheme effectively performs only the backhaul/access resource partitioning in a centralized manner. On the other hand, the actual MAC-level scheduling is still undertaken in a distributed fashion, albeit leveraging the indications produced by the central controller. The major advantages which this two-tier design exhibits, compared to a completely centralized solution, are the presence of a relatively light signaling overhead and the ability to promptly react to fast channel variations, for instance caused by small scale fading.

3.1.5 Implementation of semi-centralized allocation schemes in mmWave IAB networks

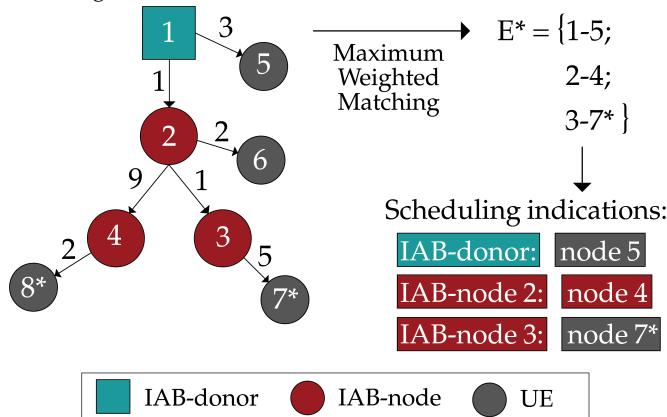
The remainder of this section discusses how the proposed scheme can be implemented in IAB deployments, with references to how the 3GPP specifications can support it. Moreover, an in-depth analysis of the framework’s communication overhead and computational complexity is provided. To such end, let $\mathcal{G} = \{V, E\}$ be the reduced network graph, computed as per Section II-C, and, conversely, let $\bar{\mathcal{G}} = \{\bar{V}, \bar{E}\}$ comprise all the nodes in the IAB network.

In general, the resource allocation framework requires (i) a central controller, which is installed on the IAB-donor, or could be deployed in a RAN Intelligent Controller (RIC) following the O-RAN architecture [37]; and (ii) a scheduler which exchanges resource coordination information with the former. In particular, and referring to the aforementioned phases of the proposed scheme, the following additional considerations can be made.



IAB-donor IAB-node UE

- (a) The original topology, exhibiting the actual cell attachments, is depicted on the left. Conversely, the reduced one is shown on the right.



IAB-donor IAB-node UE

- (b) Computation of the MWM and of the corresponding scheduling indications.

Figure 3.3: High level scheme of the initial setup and centralized scheduling indication phases.

Initial setup During this phase, which takes place when the IAB-nodes perform their first connection to the network, the controller acquires preliminary topology information by leveraging the configuration messages which are already exchanged during the typical Rel.16 Initial Access (IA) procedure [117, Section 9.6]. Therefore, no additional overhead is introduced. Specifically, a map which associates each IAB-node in the network to a list of its edges, identified by global identifiers (which from now on will be referred to as "IDs"), is computed. As a consequence, $\mathcal{O}(|V|)$ insertions in a sorted map are performed and this one-time setup exhibits a computational complexity of $\mathcal{O}(|V| \log(|V|))$.

Information collection The generation of the feedback information is performed in a distributed manner by the gNBs. To such end, the current implementation features the forwarding of information on the channel quality and buffer status, in the form of Channel Quality Informations (CQIs) and Buffer Status Reports (BSRs) respectively. This choice is driven by both the will of maximizing the re-utilization of the NR Rel.16 specifications and the goal of making use of MAC-level CEs only, hence avoiding the introduction of any constraint regarding the supported IAB-relaying architecture. In particular, the CQI and BSR information is generated by analyzing the corresponding CEs, which are already received by the scheduler of each gNB, and checking whether the source Radio Network Temporary Identifier (RNTI) belongs to an IAB-node or to a UE. In the first case, the corresponding ID is retrieved and an entry carrying such identifier along with its CQI/BSR value is generated. The feedback information concerning the UEs, instead, is averaged in the case of the CQIs and added up for the BSRs, to obtain a single value for each gNBs.

Referring to the 3GPP specifications of [121], the buffers occupancy can then be forwarded to the IAB-donor by introducing a Short BSR, which carries a single Logical Channel Group (LCG) ID and its respective buffer size. This is motivated by the fact that we do not keep track of per-flow information, i.e., we aggregate all the different RLC bearers into a single measurement report. Similarly, the channel qualities can be reported by the various IAB-nodes via an additional CQI-only Channel State Information (CSI) report, based on a Wideband (WB) measurement. Therefore, we can upper bound the size of these CEs as 11 [121] and 7 bits [122] respectively. Regard-

ing the computational complexity, in this phase we generate, at each gNB, one CQI and one BSR for each backhaul link, and (possibly) compute one cumulative CQI and BSR for the UEs. Therefore, the asymptotic complexity of this phase can be identified as $\mathcal{O}(|V|)$.

Centralized scheduling indication During this phase, the controller makes use of the feedback received from the gNBs to update the topology information, compute the weights of the various network links and to generate the centralized scheduling indications.

Regarding the former, no additional control information is required. In fact, the periodic feedback received from the various IAB-nodes, which carries a list of ID-value pairs, can be used in such regard. In particular, the controller checks the child-parent associations for discrepancies with its local knowledge, and, if so, updates the stored associations. Discrepancies can arise under two circumstances: the connection of the first UE to an IAB-node and the handover to a different parent of any IAB-node. In the first case, just the corresponding “cumulative access node” needs to be added to the aforementioned map. On the other hand, whenever a backhaul link changes, the topological information for the whole subtending tree must be updated. Since in the worst case this might require an update of the whole map, the asymptotic complexity of the topology information update is $\mathcal{O}(|V|)$. Thanks to this periodic update, our framework is robust with respect to Radio Link Failures (RLFs) and handovers, which may occur due to blockages or mobility of UEs and, possibly, gNBs.

With respect to the computation of weights for the MWM problem, we propose the following policies:

1. **Max Sum-Rate (MSR).** This policy maximizes the overall PHY-layer throughput, i.e., the utility function is

$$f_u^{\text{MSR}} \triangleq \sum_{e_{i \rightarrow k} \in E^*} c_{i,k},$$

and the weight assigned to the edge from node i to node k reads $w_{i,k} \triangleq c_{i,k}$, where $c_{i,k}$ is the capacity of the link $e_{i \rightarrow k}$.

2. **Backlog Avoidance (BA).** This resource partitioning strategy aims at avoiding congestion. Therefore, the system utility is:

$$f_u^{\text{BA}} \triangleq \sum_{e_{i \rightarrow k} \in E^*} q_{i,k},$$

where the weight $w_{i,k}$ reads $q_{i,k}$, namely, the amount of buffered data which would reach its next hop in the IAB network by crossing the link $e_{i \rightarrow k}$.

3. **Max-Rate Backlog Avoidance (MRBA).** This represents the most balanced option among the three, since it exploits favorable channel conditions while also preventing network congestion and favoring network fairness. The weight assigned to link $e_{i \rightarrow k}$ is:

$$w_{i,k} \triangleq c_{i,k} + \eta \cdot q_{i,k} \cdot \left(\frac{\mu}{\mu_{thr}} \right)^k,$$

where η , μ_{thr} and k are arbitrary parameters and μ represents the number of subframes which have elapsed since the last time edge $e_{i \rightarrow k}$ has been marked as favored.

Regardless of the specific policy used, the computation of the weights exhibits a complexity which is linear in the number of edges $|E|$.

Once the weights are computed, the controller obtains an MWM of the network via an implementation of the aforementioned T-MWM. The algorithm outputs the activation set E^* , i.e., a map associating the ID of the parent gNBs to the one of their favored downstream node. Moreover, E^* is also used by the controller in order to keep track of which link has not been favored and for how long; this information may then be used to introduce a weight prediction mechanism, improving the robustness of the scheme with respect to the information collection period. In terms of overhead, the reporting of E^* to the gNBs would feature as payload just one C-RNTI per IAB-node (at most, since some nodes might not receive any whenever they are not active in the specific MWM solution). In fact, by exploiting the Backhaul Adaptation Protocol (BAP), we can encapsulate this payload as part of a BAP message, while the destination node is already included as part of the BAP header in the “BAP destination” field. Therefore, the payload size of the scheduling indications is 16 bits.

Finally, based on the previous considerations and the analysis of Section 3.1.3.1, the overall complexity of this phase is $\mathcal{O}(|E| + |E| + |V|)$, which, when considering ST topologies, is perfectly equivalent to $\mathcal{O}(|V|)$.

Distributed scheduling allocation The last phase of the resource allocation procedure consists in the distributed MAC-level scheduling. Before assigning the available resources, the various schedulers check whether any indication has been received from the controller. Based on this condition, the buffer occupancy information is then split into two groups. The first contains the BSRs related to the favored RNTI (if any), with the caveat that if the latter indicates the cumulative access link, then this set contains the BSRs of all the UEs attached to the host gNB, while the other comprises the remaining control information. The resource allocation process is then undertaken twice: first considering the set of favored BSRs only, then the remainder of these CEs. Thanks to this repeated allocation, the favored link(s) is (are) scheduled with the highest priority, while the rest of the network only gets the remaining resources. In such a way, the information received by the controller is actually used as an **indication** and not as the eventual **resource allocation**. For instance, the gNBs are free to override these indications whenever the buffer of the favored child is actually empty, due to discrepancies between its actual status and the related information available to the controller. Moreover, the actual Downlink Control Informations (DCIs) can then be generated by the various gNBs themselves (instead of being generated only by the controller and then forwarded to the IAB-nodes), hence making use of the most updated information on the channel quality and buffer status as well. In fact, the exchange of information between the IAB-nodes and the IAB-donor introduces an inevitable delay, proportional to their distance in terms of wireless hops, between the generation of the control information at a given node (BSRs and/or CQIs) and the reception of the corresponding scheduling indications computed by the controller. Thanks to the aforementioned architecture, we limit quite significantly the performance degradation caused by these possible discrepancies between the actual nodes statuses and the (slightly outdated) information which the controller holds about them.

The computational complexity of this last phase is different from the baseline, since it requires an additional MAC-level resource allocation. However, the specific impact of this modification is difficult to determine, since

the choice of the scheduling algorithm is not part of the NR specifications. Anyhow, it is reasonable to assume such algorithm to exhibit an asymptotic complexity which is at least linear in the number of users N to be scheduled, i.e., the number of computational steps is $\mathcal{O}(N^\alpha) | N \in \mathbb{N}^+; \alpha \in \mathbb{R}, \alpha \geq 1$. Furthermore, it can be observed that in our framework, the two allocations receive as input disjoint subsets of the links; let $\hat{N}, \bar{N} | \hat{N}, \bar{N} \in \mathbb{N}; \hat{N} + \bar{N} = N$ be their respective sizes. Therefore, the number of operations required for the scheduling can be estimated as $\mathcal{O}((\hat{N} + \bar{N})^\alpha)$ for the typical network operation and $\mathcal{O}(\hat{N}^\alpha + \bar{N}^\alpha)$ when using our framework. Since the following holds:

$$(\hat{N} + \bar{N})^\alpha \geq \hat{N}^\alpha + \bar{N}^\alpha \quad \forall \alpha \in \mathbb{N}^+$$

we can claim that, under the aforementioned assumptions, the last phase of the proposed framework introduces no computational overhead with respect to the typical network operation.

In addition of the previous considerations, we also need to take into account that, if no modifications to the Rel. 16 NR specifications are introduced, a set of MAC and BAP headers would also be added to the aforementioned payload estimates; their respective sizes can be estimated as 16 [121] and 46 [123] bits, respectively. Accordingly, the worst-case *overall* network overhead can be estimated as follows. During phase 2, for each backhaul link in the network and towards the controller, up to two BSRs and CQIs are exchanged, originating from the link's parent and child respectively. Moreover, for each IAB-node in the network, one BSR and one CQI are exchanged for the (possible) "cumulative" access link. Then, in phase 3 the controller sends up to one scheduling indication per IAB-node. Letting then N be the number of IAB-nodes which are connected to the same IAB-donor, the communication overhead can be upper bounded by $N \cdot (2 + 1) \cdot (65 + 69) = 402 \cdot N$ [bits] in the UL and $76 \cdot N$ [bits] in the DL.

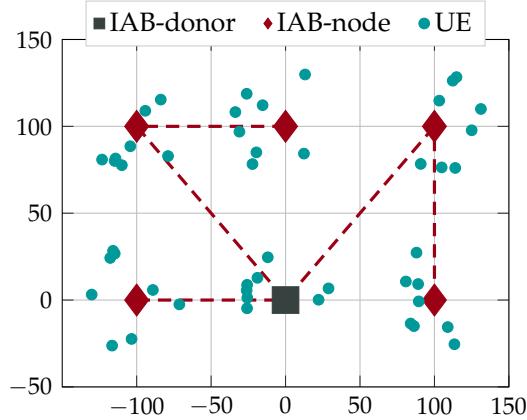
Notably, the 3GPP also considered the possibility of realizing heterogeneous IAB deployments [117], in which IAB-nodes hold an additional connection with a macro cell (ideally co-located with the IAB-donor) to handle the control plane. In this context, our framework can be enhanced by carrying feedback information (i.e., CQIs and BSRs) and scheduling indications

over the additional connection, reducing the overhead and avoiding the need to travel through multiple hops before reaching the IAB-donor.

We implemented the proposed resource allocation scheme in the popular open source simulator ns-3, exploiting the mmWave module [17] and its IAB extension [45], to characterize the system-level performance of the proposed solution with realistic protocol stacks, scenarios, and user applications.

The ns-3 mmWave module is based on [124] and features highly customizable PHY and MAC layer implementations, with an NR-like flexible OFDM numerology and frame structure. It also includes accurate interference and error models, as well as a detailed channel model, which is compliant with the 3GPP specifications [22] and accounts for large and small scale fading phenomena, as well as for interference. Additionally, the IAB module [45] models wireless relaying functionalities which mimic the specifications presented in [117]. Specifically, this module supports both single and multi-hop deployment scenarios, auto-configuration (within the network) of the IAB-nodes and a detailed 3GPP protocol stack, allowing wireless researchers to perform system-level analyses of IAB systems in ns-3.

It is of particular relevance to understand how the scheduling operations are implemented in the IAB module, since they offer not only the baseline for the proposed scheme, but also valid guidelines for real-world deployments. The current ns-3 IAB schedulers exhibit a TDMA-based multiplexing between the access and backhaul interfaces. Moreover, scheduling decisions are undertaken in a distributed manner across the IAB network, i.e., each gNB allocates the resources which its access interface offers (to both UEs and IAB-nodes) independently of the other gNBs in the network. In fact, in an IAB network these scheduling decisions are *almost* independent of one another: if a parent node schedules the backhaul interface of a downstream node, clearly the latter will be constrained in its own scheduling decisions, as it will not be allowed to allocate the time resources which have already been scheduled for backhaul transmissions by its parent. Therefore, in a tree-based, multi-hop wireless network the various gNBs need to know in advance the scheduling decisions performed by their upstream nodes: to solve this problem, the authors of the IAB module for ns-3 introduced a "*look-ahead backhaul-aware scheduling mechanism*" [45]. Such mechanism features an exchange of DCI between the access and backhaul interfaces: in such a way, any time resources already scheduled by the parent for backhaul commu-



(a) A realization of the simulation scenario; the dotted lines represent the cell-attachments of the IAB-nodes.

SIMULATION PARAMETERS	
PARAMETER	VALUE
Number of runs N_{runs}	25
Simulation time T_{sim}	3 s
MWM period T_{alloc}	{1, 2, 4} subframes
Layer 4 protocol	{UDP, TCP}
UDP packet size s_{UDP}	{50, 100, 200, 500} B
Weight policy f_u	{MSR, BA, MRBA}

(b) Simulation parameters.

Figure 3.4: Simulation configuration.

nifications can be marked as such by the corresponding downstream node, preventing any overlap with other transmissions. Furthermore, the *look-ahead* mechanism requires the schedulers of the various gNBs to commit to their resource allocation for a given time T at a time $T - k$, where $k - 1$ is the maximum distance (in terms of wireless hops) of any node from the donor. In such a way, the DCIs will have time to propagate across the IAB network and reach the farthest node at time $T - 1$, thus allowing its scheduler to perform the resource allocation process at least one radio subframe in advance.

3.1.5.1 Simulation scenario and parameters

The purpose of these simulations is to understand the performance of the proposed resource partitioning framework in the context of its target deployment, i.e., a multi-hop IAB network. As a consequence, the reference scenario

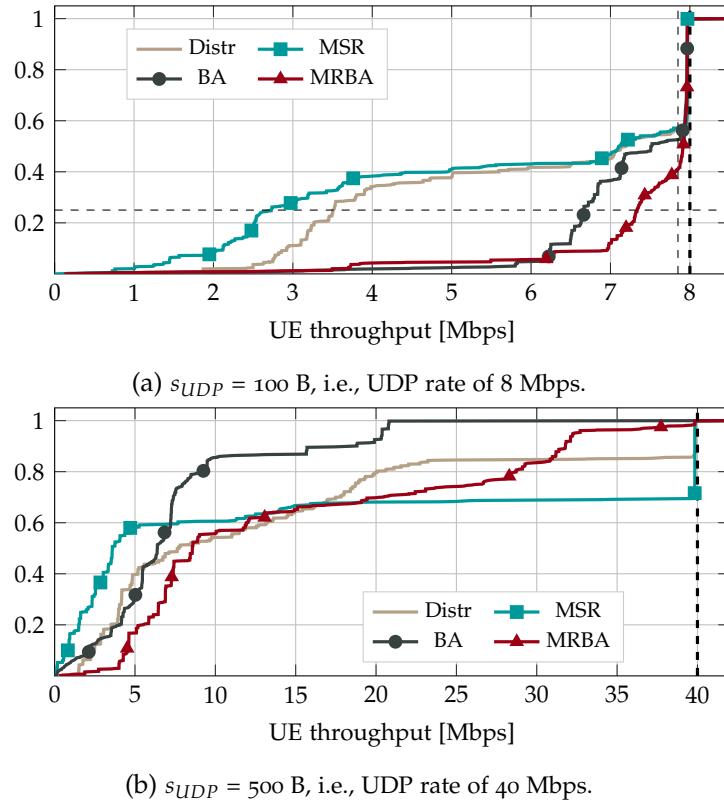


Figure 3.5: Per-UE end-to-end throughput Empirical Cumulative Distribution Functions (ECDFs). The thick dashed line represents the rate of the UDP sources.

consists of a dense urban deployment with a single IAB-donor and multiple IAB-nodes, as depicted in Fig. 3.4a. In particular, the various gNBs are distributed along an urban grid where the donor is located at the origin while the IAB-nodes are deployed along the street intersections, with a minimum inter-site distance of 100 m. The IAB-nodes attachments are computed using the so-called HQF policy presented in [116]; however, this choice does not introduce any loss of generality since such parameter is fixed for all the runs. A given number of UEs are deployed within the surroundings of these base stations, with an initial position which is randomly sampled from circles of radius ρ and whose centers are the various gNBs. A summary of the simulation parameters in provided in Tab. 3.4b.

3.1.6 Performance evaluation

Both the IAB-donor and the IAB-nodes are equipped with a phased array featuring 64 antenna elements, and transmit with a power of 33 dBm; conversely UEs are equipped with 16 antenna elements and their transmission power is restricted to 23 dBm. Notably, the presence of additional antenna elements at the gNBs is a key (but reasonable) assumption, as it allows base stations to achieve a high beamforming gain. In turn, it is possible to achieve a high capacity, which is fundamental to avoid performance bottlenecks, given the absence of a fiber backhaul. The UEs download data which originates from sources that are installed on a remote host; both the UDP and the TCP are used. For the UDP simulations, the rate of the sources is varied from 4 to 40 Mbps to introduce different degrees of saturation in the network. Therefore, in these simulations only DL traffic is considered. Finally, the performance of the proposed policies is hereby compared with the baseline of [45], indicated as “Dist.” by examining end-to-end throughput, latency, and a network congestion metric.

3.1.6.1 Throughput

The first metric which is inspected in this analysis is the end-to-end throughput at the application layer. As a consequence, only the packets which are correctly received at the uppermost layer of the destination node in the network are taken into account. In particular, for each UE and each simulation run, the long-term average throughput is computed as follows:

$$S_{k,n}^{\text{APP}} \triangleq \frac{B(T_{\text{sim}}, k, n)}{T_{\text{sim}}}$$

where $B(T, k, n)$ is the cumulative number of bits received up to time T by the k -th UE, during the n -th simulation run. Then, the distribution of \mathbf{S}^{APP} , namely, the vector containing the collection of the $S_{k,n}^{\text{APP}}$ values across the different runs and UEs, is analyzed.

Figs. 3.5a and 3.5b report the ECDF of \mathbf{S}^{APP} , for a UDP packet size of 100 and 500 bytes, respectively, and the policies introduced in Section 3.1.5. In the former, we can notice that the introduction of the semi-centralized framework increases by up to 15% the percentage of UEs whose throughput almost matches the rate of the UDP sources, i.e., achieving approximately

7.9 Mbps. Moreover, by focusing on the leftmost portion of Fig. 3.5a we can observe another interesting result, concerning the throughput experienced by the UEs which do not fulfill their QoS requirements. In fact, with respect to the first quartile the distributed scheduler and the MSR policy achieve the worse performance. On the other hand, the MRBA and BA policies significantly improve these results, even though the extent of such improvements varies quite dramatically across the two.

In particular, compared with the distributed case the BA and MRBA policies introduce a 2 and 3-fold increase of the worst case throughput respectively, coupled with a significantly lower variance in both cases.

These results can be explained as follows: since a UDP packet size of 100 bytes does not saturate the capacity of the access links, the main performance bottleneck of this configuration is represented by the buffering of the aggregated traffic on the intermediate backhaul links. Therefore, the MSR policy provides no improvements compared to the performance of the distributed scheduler, since it simply favors the links which exhibit a higher SINR. Conversely, the prioritization of the most congested links which is introduced by the other two strategies successfully tackles the former problem. In particular, the BA policy exhibits the highest worst case throughput, while also satisfying the QoS requirements of approximately 40% of the UEs. Moreover, the bias towards high SINR channels introduced by the MRBA strategy further improves the higher percentiles, compared to the BA policy, and dramatically outperforms MSR and the baseline across all percentiles.

By increasing the UDP packet size to 500 bytes, the network becomes noticeably saturated, as depicted by Fig. 3.5b; in fact, in this instance only a minority of the UEs achieves a throughput which is comparable to the source rate. With this configuration, the BA strategy achieves the worst performance, providing a significantly lower throughput across most percentiles. On the other hand, the remaining strategies both introduce significant improvements, although with different trade-offs. In particular, compared to distributed case, the MSR policy exhibits an increase of approximately 20% of the number of UEs which satisfy their QoS requirements, albeit at the cost of worse lower percentiles. The MRBA, conversely, introduces performance benefits which mostly affect the bottom percentiles only. However, with this strategy only a limited portion of the UEs achieves the target throughput of 40 Mbps. As a consequence, we can conclude that with the configuration

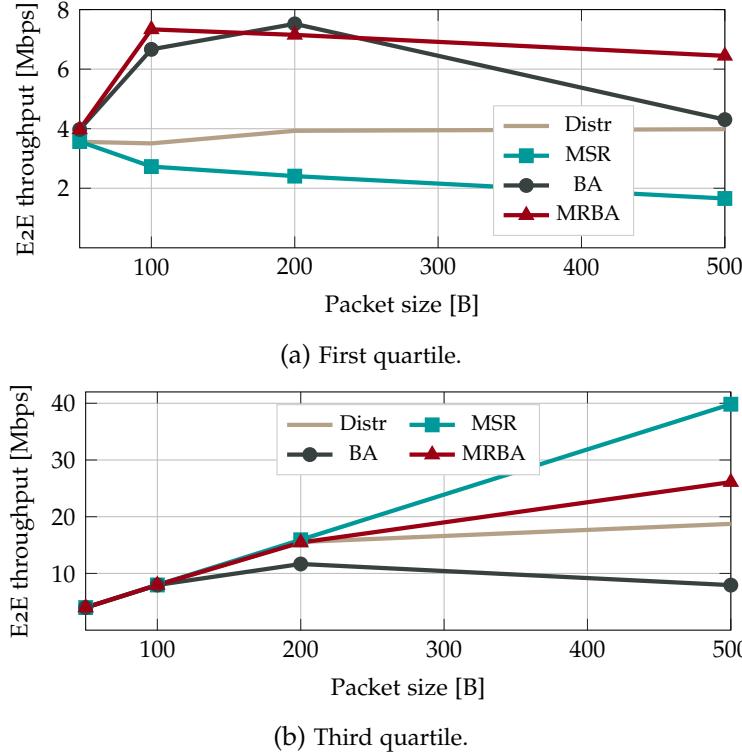


Figure 3.6: End-to-end throughput quartiles, for $s_{UDP} \in \{50, 100, 200, 500\}$ B.

depicted in Fig. 3.5b the network is approaching the capacity of the mmWave channels. Therefore, buffering phenomena are likely occurring at each intermediate IAB-node. Moreover, we can say that in a saturated network the congestion is so severe that prioritizing the bottleneck links is not enough: we also need to take into account the channel conditions and prioritize the links which not only are congested, but also have the “biggest chance” of getting rid of the buffered data due to the temporary better channel quality.

Finally, Fig. 3.6 presents the first and third quartiles of \mathbf{S}^{APP} as a function of the UDP packet size s_{UDP} . It can be noted that, with respect to the first quartile, the MRBA outperforms all the other policies by delivering a throughput which is up to 90% higher than the one obtained by the distributed scheduler. On the other hand, Fig. 3.6b shows how the best third quartile is achieved by MSR, with up to a 2-fold improvement over the distributed solution. Furthermore, we can observe how the positive impact of the BA strategy is inversely proportional to the saturation in the network. We can then conclude that the

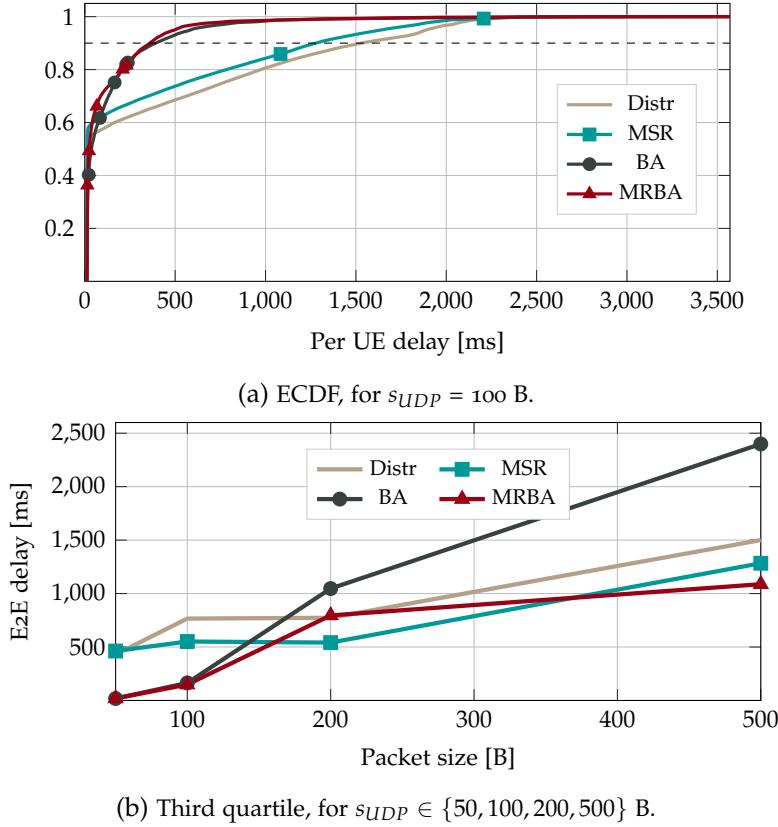

 (b) Third quartile, for $s_{UDP} \in \{50, 100, 200, 500\}$ B.

Figure 3.7: Per-UE end-to-end delay statistics.

bias it introduces loses its effectiveness as the buffering phenomena start to affect the majority of the IAB-nodes.

3.1.6.2 Latency

Just like the aforementioned metric, the latency is measured end-to-end at the application layer. Thanks to this choice, the resulting delay accurately represents the system-level performance, as it includes the latency which is introduced at each hop in the IAB network.

In particular, for each packet correctly received at the uppermost layer of its final destination, the following quantity is traced:

$$D_i^{\text{APP}} \triangleq \sum_{l_k \in \mathcal{E}_i} D_i^{l_k}$$

where \mathcal{E}_i comprises the links in the IAB network that are crossed by the i -th packet, while the term $D_i^{l_k}$ indicates its point-to-point latency over the path

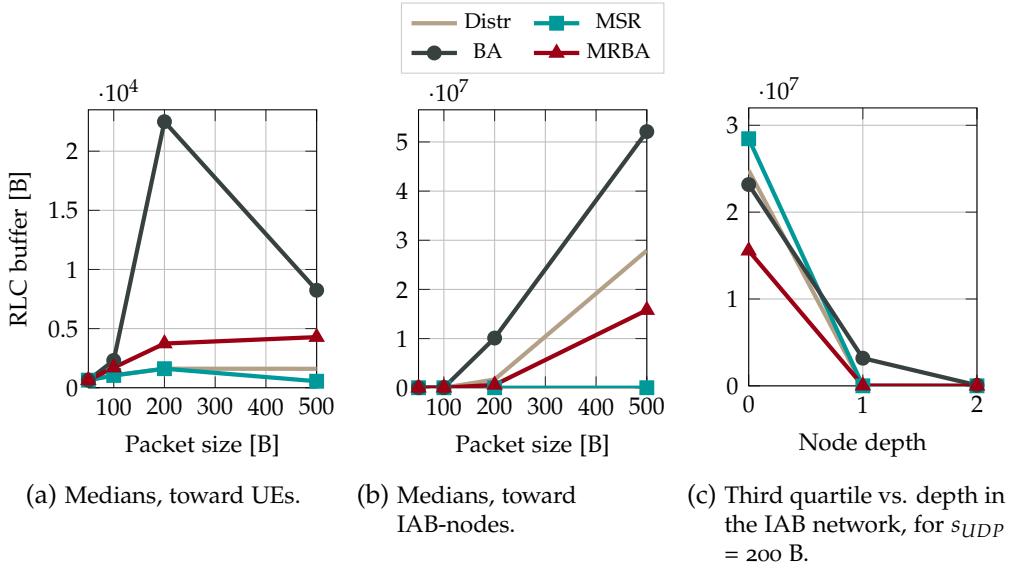


Figure 3.8: Buffer occupancy statistics, for $s_{UDP} \in \{50, 100, 200, 500\}$ B.

link l_i . Finally, these values are collected for each of the various runs into the vector \mathbf{D}^{APP} and its statistical properties are inspected.

Fig. 3.7a shows the empirical ECDF of \mathbf{D}^{APP} for a packet size of 100 bytes. It can be noticed that, in this case, the 90th percentile achieved by the BA and the MRBA policies are approximately 20 % smaller than the one obtained by the distributed scheduler. Moreover, these strategies manage to dramatically reduce the number of packets received with extremely high delay, i.e., in the order of seconds, showing the dramatic impact of buffering in the baseline configuration. Conversely, the MSR policy provides the best performance with respect to the best case delay only, although it still outperforms quite significantly the distributed strategy.

These trends are exacerbated by Fig. 3.7b, which shows the third quartile of \mathbf{D}^{APP} as a function of the UDP packet size s_{UDP} . In fact, we can notice that the effectiveness of the BA policy is inversely proportional to the network saturation; the opposite holds true with respect to the MSR strategy. It follows that, for UDP rates in the order of 5 to 10 Mbps, the network is mainly plagued by local congestion which causes the insurgence of buffering in some of the nodes. Conversely, as the rate of the UDP sources increases the system shifts to a capacity-limited regime, a phenomenon which explains the dominance of the MSR and MRBA policies.

3.1.6.3 Network congestion

The network congestion is measured by collecting, every T_{alloc} subframes, the RLC buffers status of the various nodes into the vector \mathbf{B}^{RLC} . It must be noted that, since RLC Acknowledged Mode (AM) is used, these values will indicate data which is related to both new packets and possible retransmissions.

Figs. 3.8a and 3.8b show the median of \mathbf{B}^{RLC} , for traffic flows whose next hop in the network is represented by either UEs or IAB-nodes respectively. Specifically, the BA strategy achieves the worst performance in this metric, leading to unstable systems in the cases of $s_{UDP} = \{200, 500\}$ B. A reason for this behavior can be found in the “locality” of the BA policy criteria and the lack of influence of the past allocations on the weights. These characteristics may lead to favoring the same link in a repeated manner, hence offering little remedy to the end-to-end congestion.

On the other hand, the buffer occupancy achieved by the MSR strategy depicts an effectiveness which, in accordance with previous observations, is proportional with respect to the source rate. In particular, a dramatic decrease of up to 4 orders of magnitude is achieved for $s_{UDP} = 500$ B. Finally, when compared to the distributed scheduler, the MRBA policy also achieves a lower median RLC buffer occupancy towards the backhaul links, albeit the difference is less striking than in the case of the MSR policy, and at the cost of slightly more congested UE buffers.

Additionally, Fig. 3.8c depicts the third quartiles of \mathbf{B}^{RLC} as a function of the depth of the corresponding gNB in the IAB network. It is possible to notice that, regardless of the policy in use, the amount of buffering at the various gNBs generally decreases as their distance to the donor increases. This follows from the fact that nodes which have a lower depth exhibit, on average, a bigger subtending tree; therefore the amount of traffic which makes use of their backhaul links is significantly higher.

3.1.6.4 Performance with TCP traffic

This subsection extends the aforementioned analysis by inspecting the performance of the proposed scheme in the case of TCP traffic. Specifically, a TCP full-buffer source model is used, and the various semi-centralized resource allocation policies are compared against the distributed scheduler.

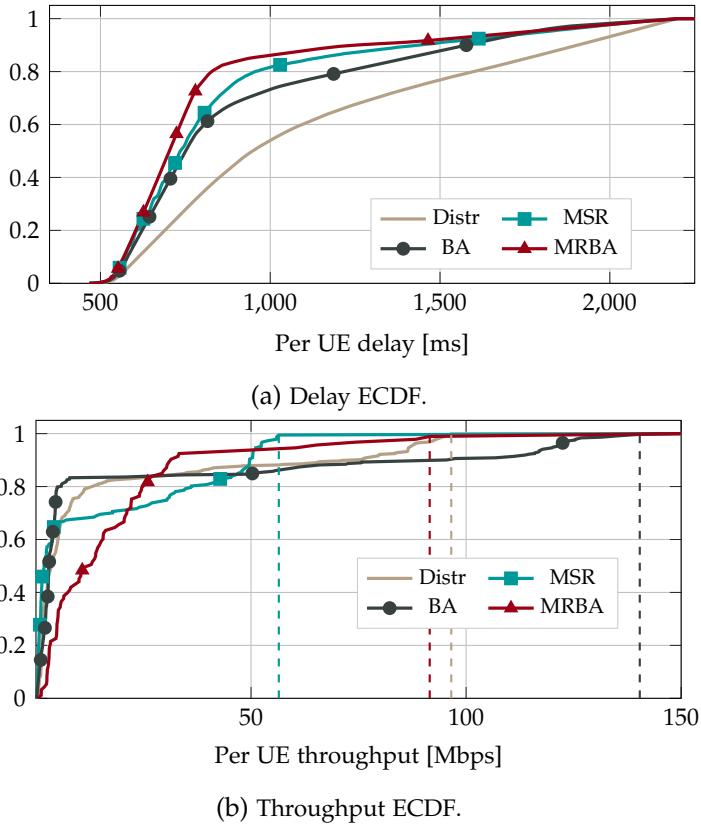


Figure 3.9: End-to-end delay and throughput statistics, for TCP layer 4 protocol.

Fig. 3.9a shows the ECDF of the end-to-end delay experienced by the successfully received packets. Similarly to the UDP case, the distributed scheduler exhibits the worst results in this regard. In fact, the performance benefits introduced by the semi-centralized policies are noticeable across all percentiles. In particular, with this configuration the MRBA policy provides the best results, followed quite closely by the BA and MSR strategies. Fig. 3.9b, which depicts the statistics of the end-to-end throughput achieved by the various UEs, further explains the effect on the system of the various semi-centralized policies. In particular, the BA policy achieves, approximately, a 45% increase of the peak throughput. Conversely, the MRBA strategy causes a redistribution of the achieved data rate, massively improving the lower quartiles (up to the 80-th), albeit at the expense of the maximum throughput. Finally, MSR also causes a redistribution of the throughput across the different percentiles, but the net benefit is less noticeable.

Therefore, we can conclude that regardless of the specific policies used, the proposed scheme improves the system performance also with this configuration, by limiting the insurgence of local buffering and aiding the end-to-end congestion control mechanism offered by TCP. Furthermore, it can be noted that both a prioritization of the most congested links and of the channels featuring a higher quality results in performance benefits in the average case, although it also causes a decrease of the network fairness. On the other hand, the MRBA policy manages to optimize the backhaul/access resource partitioning, while introducing an increase in the throughput fairness at the same time.

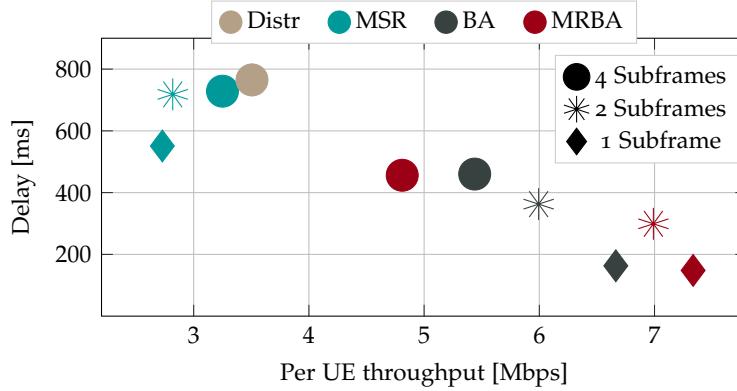
3.1.6.5 Further considerations

It is of particular relevance to analyze the performance of the semi-centralized policies when relaxing the most restrictive hypothesis, i.e., the capability of exchanging feedback information in a timely manner.

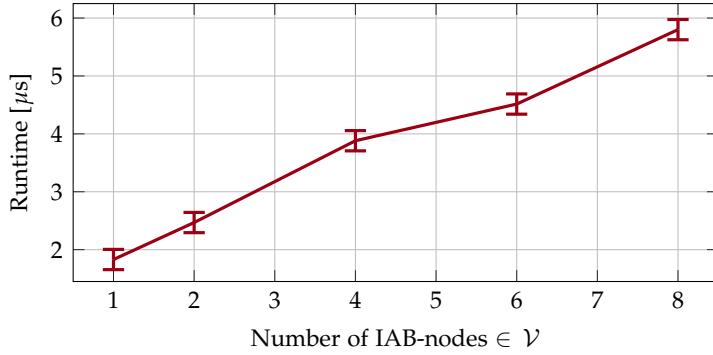
Actually, such analysis provides also insights regarding the effects of errors and/or crashes in the control messages. Indeed, both control and data channels implement error detection mechanisms, hence we deem the likelihood of undetected errors in the feedback information to be negligible. As a consequence, the errors would be detected at the receiver, and lost information would be either retransmitted by the source or simply discarded, waiting for the following periodic update; in both cases, the net effect would be a delay in the reception of the message.

To such end, Fig. 3.10a shows the performance of the proposed framework as a function of the semi-centralized allocation period T_{alloc} . In particular, each of the depicted points represents the joint end-to-end throughput and delay achieved with the different configurations.

As expected, in general the effectiveness of the various semi-centralized policies progressively deteriorates as the frequency of the scheduling indications decreases. Interestingly, the BA policy exhibits the lowest performance degradation with respect to an increase of the allocation period, which suggests that this phenomenon has a slower evolution over time compared to the one exhibited by the channels quality. Nevertheless, the key takeaway is that all of the proposed allocation strategies except MSR outperform the distributed solution, across both metrics. In fact, the latter exhibits the lowest throughput first quartile, but only because it introduces a strong bias



(a) Combined per UE end-to-end throughput first quartile and delay third quartile, as a function of the semi-centralized allocation period T_{alloc} .



(b) MWM runtime as a function of the number of IAB-nodes in the network.

Figure 3.10: Considerations on the formulated assumptions.

towards high SINR channels, as discussed in Section 3.1.6.1. However, the trend depicted by Fig. 3.10a also suggests that there exists a threshold value of T_{alloc} after which the performance of the proposed frameworks brings only marginal performance benefits.

Additionally, the running time of the MWM algorithm presented in Section 3.1.3.1 was analyzed, in order to understand whether it may partially invalidate the timely feedback assumption. Specifically, Fig. 3.10b presents the statistics of the various MWM execution times, obtained on a machine equipped with an i7-6700 4-core processor clocked at 3.4 GHz. The first observation which can be made is that this empirical analysis confirms the previously estimated asymptotic complexity, depicting a running time which exhibits a linear dependence on the number of gNBs in the network. Furthermore, it can be noted that the runtime of the MWM algorithm does not exceed 6 μ s, even for a significant number of IAB-nodes connected to the same

3 Towards wireless-backhauled next-generation cellular networks

IAB-donor. As a consequence, we can conclude that the execution times of the semi-centralized allocation process do not pose any threat to the timely feedback assumption, since they are reasonably smaller than the duration of the minimum semi-centralized allocation period.

3.2 Risk-averse learning for latency minimization in mmWave integrated access and backhaul networks

Regardless of the methodology, prior works mostly aim at maximizing the network capacity [101, 103, 125–130], minimizing latency [109, 131] and improving throughput fairness [84, 126]. Although these approaches successfully improve the network performance, Mobile Network Operators (MNOs) are often more *concerned about their reliability*. For this reason many commercial products rely on *simplified* but reliable algorithms for resource allocation, despite their sub-optimal performance. In this section, we address these limitations by proposing Safehaul, a reinforcement learning-based solution for scheduling and path selection in IAB mmWave systems which reaps the benefits of learning-based algorithms, while guaranteeing reliable network performance. To this end, we use the concept of risk aversion, commonly used in economics [132, 133], to measure and enhance the reliability of Safehaul. The following summarizes our contributions:

- We model the scheduling and path selection problem in IAB mmWave networks as a multi-agent multi-armed bandit problem (Section 3.2.2). We consider multiple fiber base stations, simultaneously supporting many self-backhauled mmWave base stations. In our model, the self-backhauled base stations independently decide the links to activate. The consensus among the base stations is reached via standard-defined procedures (Section 3.2.3.3).
- We present the first solution to provide reliable performance in IAB-enabled networks (Section 3.2.3). Specifically, we investigate the joint minimization of the average end-to-end latency and its expected tail loss. To this aim, we propose Safehaul, a learning approach that leverages the coherent risk measure CVaR[132]. CVaR measures the tail average of the end-to-end latency distribution that exceeds the maximum permitted latency, thus ensuring the network’s reliability.
- We analytically bound the regret of Safehaul, i.e., we bound the loss of Safehaul compared to the case when the delays associated to all end-to-end paths between self-backhauled base stations and fiber base stations are known a priori. We show that, for the case when there are no

conflicts between the decisions of the self-backhauled base stations, the average regret of Safehaul tends to zero as the time increases. This regret bound characterizes the learning speed and proves that Safehaul converges to the optimal scheduling and path selection solution that jointly minimizes the average end-to-end latency and its expected tail loss.

- We provide a new means of simulating multi-hop IAB networks by extending NVIDIA’s GPU-accelerated simulator Sionna [134] (Section 3.2.4). Specifically, we add codebook-based analog beamforming capabilities for both uplink and downlink communications. In addition, we add internal RT of Sionna in order to generate Channel Impulse Response (CIR). Further, we extend Sionna by implementing system-level components such as layer-2 schedulers and buffers and BAP-like routing across the IAB network. We believe our IAB extensions will be instrumental for the open-source evaluation of future research on self-backhauled mmWave networks.
- Exploiting the above simulator, we evaluate and benchmark Safehaul against two state-of-the-art algorithms [131, 135] based on deployment in two different locations (Manhattan and Padova). The results confirm that Safehaul is significantly more reliable than the considered benchmarks, as it exhibits much tighter variance in terms of both latency (up to 71.4% smaller) and packet drop rate (at least 39.1% lower). Further, Safehaul achieves up to 43.2% lower average latency and 11.7% higher average throughput than the reference schemes.

3.2.1 System Model

We consider a cellular system with N base stations capable of self-backhauling and D base stations with a fiber connection to the core network. Following 3GPP terminology, we refer to self-backhauled base stations as IAB-nodes (BS-nodes) and to fiber base stations as IAB-donors (BS-donors)¹. Each BS-node connects to the core network via a (multi-hop) wireless link to a BS-donor. The sets of all BS-nodes and BS-donors are denoted by $\mathcal{N} =$

¹Please note that throughout the section we will use interchangeably BS-nodes and IAB-nodes (and similarly for BS-donors and IAB-donors)

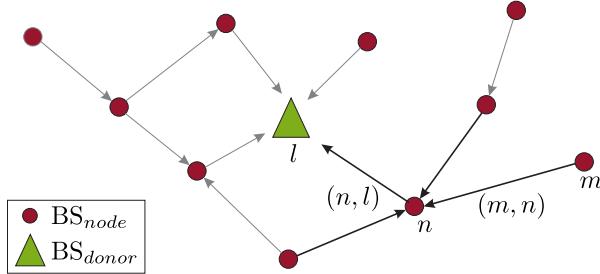


Figure 3.11: Example of a graph \mathcal{G}_i

$\{1, \dots, N\}$ and $\mathcal{D} = \{N + 1, \dots, N + D\}$, respectively. The system works in a time-slotted fashion starting from time slot $i = 1$ until a finite time horizon I . All the time slots $i = 1, \dots, I$ have the same duration. The BS-nodes are equipped with two RF chains. One RF chain is used exclusively for the communication with cellular users (access network), while the second RF chain is used for self-backhauling. In line with the 3GPP specification [117], we assume half-duplex self-backhauling, i.e., in each time slot i a BS-node can either transmit, receive or remain idle.

We model the connections between the base stations in slot i as a graph $\mathcal{G}_i = \{\mathcal{V}, \mathcal{E}_i\}$, see Figure 3.11. The set $\mathcal{V} = \mathcal{N} \cup \mathcal{D}$ of vertices is formed by all the BS-nodes and BS-donors in the system. The set \mathcal{E}_i of edges is composed of the available wireless links (n, l) between a BS-node $n \in \mathcal{N}$ and any BS (BS-donor or BS-node) $l \in \mathcal{V}$ in time slot i . Note that \mathcal{G}_i is not static. In a given time slot i , some links may be unavailable due to failure, blockage, or interference. Thus, only feasible wireless links are considered in the set \mathcal{E}_i . The path $X_{n,d}$ from BS-node n to any BS-donor d is a sequence of intermediate links (n, l) . $X_{n,d}$ changes over time according to the traffic loads of the intermediate BS-nodes and to the channel conditions. We model the activation of link (n, l) with the binary variable $x_{n,l,i}$. When $x_{n,l,i} = 1$, the link is activated and BS-node n transmits to BS $l \in \mathcal{V}$ in time slot i , whereas $x_{n,l,i} = 0$ indicates that the link is deactivated. $x_{n,n,i} = 1$ indicates that BS-node n does not transmit nor receives backhaul data in time slot i .

Each BS-node n has a finite data buffer with capacity B_n^{\max} to store the backhaul data to be transmitted to any of the BS-donors. In each time slot i , BS-node n is characterized by its load and average queuing time. The load, denoted by $B_{n,i} \in \mathbb{N}$, indicates the number of data packets stored in the buffer at the beginning of time slot i . The average queuing time $t_{n,i}^q \in \mathbb{R}^+$ is the average number of time slots the current packets in the data buffer

have been stored. Additionally, we denote by $M_{n,l,i} \in \mathbb{N}$ the number of data packets transmitted from n and successfully received at l in time slot i (i.e., when $x_{n,l,i} = 1$), and with $t_{n,l,i}^{\text{tx}} \in \mathbb{R}^+$ the transmission time needed to send these packets. Note that $M_{n,l,i} \leq B_{n,i}$ as only packets stored in the data buffer can be transmitted. At the receiver BS-node l , the load $B_{l,i+1}$ of its data buffer is updated at the beginning of the next time slot $i + 1$ such that $B_{l,i} + M_{n,l,i} \leq B_l^{\max}$ holds. That is to say, packets exceeding the buffer capacity are dropped. Finally, when $x_{n,l,i} = 0$ both $M_{n,l,i}$ and $t_{n,l,i}^{\text{tx}}$ are equal to zero.

We define the maximum tolerable latency T_{\max} as the maximum time a packet can take from its source BS-node to any BS-donor. Any packet that is not delivered before T_{\max} milliseconds will be dropped. The average maximum end-to-end latency $\bar{T}_{n,d}$ from BS-node n to BS-donor d is the average, over the complete time horizon I , of the maximum delay a packet originating from BS-node n takes to reach any BS-donor d in time slot i . This is calculated as $\bar{T}_{n,d} = \frac{1}{I} \sum_{i=1}^I T_{n,d,i}$, where $T_{n,d,i}$ is the maximum end-to-end latency among all the packets originating in BS-node n which reach BS-donor d in time slot i . $T_{n,d,i}$ is a sample of the random variable $T_{n,d}$ drawn from an unknown stationary probability distribution P that depends on the links $x_{n,l,i'}$, $n \in \mathcal{N}$, $l \in \mathcal{V}$, $i' = 1, \dots, i$, activated up to time i , the user's mobility, the location of the BS-node n , the interference in the system, and the queue dynamics. Accordingly, we define the average maximum end-to-end latency in the system \bar{T} as

$$\bar{T} = \frac{1}{ND} \sum_{n=1}^N \sum_{d=1}^D \bar{T}_{n,N+d}. \quad (3.1)$$

3.2.2 Problem Formulation

The joint minimization of the average maximum end-to-end latency and the expected value of its tail loss in IAB-enabled networks is formulated in this section. We first introduce CVaR, the risk metric accounting for minimizing the events in which the end-to-end latency is higher than T_{\max} . Next, we formulate the optimization problem in the complete network.

3.2.2.1 Preliminaries on CVaR

Traditionally, latency minimization in IAB-enabled networks has focused on optimizing the expected value of a latency function [109, 131]. However, such an approach fails to capture the time variability of the latency distribution, thus potentially leading to unreliable systems in which $T_{n,d,i} > T_{\max}$, for any $i = 1, \dots, I$, $n \in \mathcal{N}$ and $d \in \mathcal{D}$. For this purpose, we consider not only the average end-to-end latency \bar{T} in the system, but also its expected tail loss based on the CVaR[132, 136].

Having in mind that $T_{n,d}$ is a random variable, we assume it has a bounded mean on a probability space (Ω, \mathcal{F}, P) , with Ω and \mathcal{F} being the sample and event space, respectively. Using a risk level $\alpha \in (0, 1]$, the $\text{CVaR}_\alpha(T_{n,d})$ of $T_{n,d}$ at risk level α quantifies the losses that might be encountered in the α -tail. More specifically, it is the expected value of $T_{n,d}$ in its α -tail distribution [136]. Formally, $\text{CVaR}_\alpha(T_{n,d})$ is defined as [132]

$$\text{CVaR}_\alpha(T_{n,d}) = \min_{q \in \mathbb{R}} \left\{ q + \frac{1}{\alpha} \mathbb{E}[\max\{T_{n,d} - q, 0\}] \right\}, \quad (3.2)$$

where the expectation in (3.2) is taken over the probability distribution P . Note that lower $\text{CVaR}_\alpha(T_{n,d})$ results in higher system reliability because the expected end-to-end latency in the α -worst cases is low. Moreover, note that α is a risk aversion parameter. For $\alpha = 1$, $\text{CVaR}_\alpha(T_{n,d}) = \mathbb{E}[T_{n,d}]$ which represents the traditional risk-neutral case. Conversely, $\lim_{\alpha \rightarrow 0} \text{CVaR}_\alpha(T_{n,d}) = \sup\{T_{n,d}\}$. CVaR has been shown to be a coherent risk measure, i.e., it fulfills monotonicity, subadditivity, translation invariance, and positive homogeneity properties [137].

3.2.2.2 Optimization Problem

We jointly minimize the average maximum end-to-end latency and its expected tail loss for each BS-node. For this purpose, we decide which of the (n, l) links to activate in each time slot i during the finite time horizon I . In the following, we formulate the optimization problem from the network perspective and consider the sum over all BS-nodes in the system. The latency minimization problem should consider three different aspects: (i) link activation is constrained by the half-duplex nature of self-backhauling, (ii) only

data stored in the data buffers can be transmitted, and (iii) packet drop due to buffer overflow should be avoided. Formally, the problem is written as:

$$\underset{\{x_{n,l,i}\}}{\text{minimize}} \quad \sum_{n \in \mathcal{N}} \left(\sum_{d \in \mathcal{D}} \left(\frac{1}{I} \sum_{i=1}^I T_{n,d,i} \right) + \eta \text{CVaR}_\alpha(T_{n,f}) \right) \quad (3.3a)$$

subject to

$$\sum_{l \in \mathcal{V}, l \neq n} x_{n,l,i} + \sum_{l \in \mathcal{N}} x_{l,n,i} = 1, \quad n \in \mathcal{N}, i = 1, \dots, I, \quad (3.3b)$$

$$B_{n,i} \geq M_{n,l,i}, \quad n \in \mathcal{N}, l \in \mathcal{V}, i = 1, \dots, I, \quad (3.3c)$$

$$B_{l,j} + M_{n,l,j} \leq B_l^{\max}, \quad n \in \mathcal{N}, l \in \mathcal{V}, i = 1, \dots, I, \quad (3.3d)$$

$$x_{n,l,i} \in \{0, 1\}, \quad n \in \mathcal{N}, l \in \mathcal{V}, i = 1, \dots, I. \quad (3.3e)$$

In (3.3a), $\eta \in [0, 1]$ is a weighing parameter to control the trade-off between minimizing the average maximum end-to-end latency $\bar{T}_{n,d}$ and the expected loss of its tail. The constraint in (3.3b) considers half-duplex transmissions by ensuring that, in each time slot i , every IAB-node communicates with up to one of its neighbors by either receiving or transmitting backhaul data. (3.3c) considers data causality, i.e., only data already stored in the data buffers can be transmitted, and (3.3d) prevents buffer exhaustion. As the considered scenario is not static, solving (3.3) would require complete non-causal knowledge of the system dynamics during the complete time horizon I . However, in practical scenarios, knowledge about the underlying random processes is not available in advance. For example, the BS-node's loads $B_{n,i}$ depend not only on the transmitted and received backhaul data, but also on the randomly arriving data from its users. Similarly, the amounts of transmitted data $M_{n,l,i}$ depend on the varying channel conditions of both BS n and l . As a result, the exact values of $T_{n,l,i}$, $B_{n,i}$ and $M_{n,l,i}$ are not known beforehand. For this reason, we present in Sec. 3.2.3 Safehaul, a multi-agent learning approach to minimize in each BS-node the average maximum end-to-end latency and the expected value of the tail of its loss.

3.2.3 Our proposed solution: Safehaul

In this section, we describe Safehaul, a multi-agent learning approach for the joint minimization of the average maximum end-to-end latency and its expected tail loss in IAB mmWave networks. In Safehaul, each BS-node inde-

pendently decides which links (n, l) to activate in every time slot i by leveraging a multi-armed bandit formulation. The consensus among the BS-nodes is reached by exploiting the centralized resource coordination and topology management role of IAB-donors [138, Sec. 4.7.1].

3.2.3.1 Multi-Armed Bandit Formulation

Multi-armed bandit is a tool well suited to problems in which an agent makes sequential decisions in an unknown environment[139]. In our scenario, each BS-node n decides, in each time slot i , which of the links (n, l) to activate without requiring prior knowledge about the system dynamics. The multi-armed bandit problem at BS-node n can be characterized by a set \mathcal{A}_n of actions and a set \mathcal{R}_n of possible rewards. The rewards $r_{n,i} \in \mathcal{R}_n$ are obtained in each time slot i as a response to the selected action $a_{n,i} \in \mathcal{A}_n$ and the observed latency. Since every BS-node n selects only one action during each time slot, we enforce the half-duplex constraint in (3.3b) by defining the set of possible actions as the set of feasible links for BS-node n . In particular, we define \mathcal{A}_n for $n \in \mathcal{N}$ as $\mathcal{A}_n = \{(n, l), (m, n) | m \in \mathcal{N}, l \in \mathcal{V}\}$, where link (n, n) indicates that BS-node n remains idle. As blockages, overloads, or failures might render certain links (n, l) temporarily unavailable, we define the set $\mathcal{A}_{n,i} \subseteq \mathcal{A}_n$ of available actions in time slot i as $\mathcal{A}_{n,i} = \{(n, l), (l, n) | (n, l), (l, n) \in \mathcal{E}_i\}$. Selecting action $a_i = (n, l)$ in time slot i implies $x_{n,l,i} = 1$.

The rewards $r_{n,i}$ are a function of the end-to-end latencies $T_{n,d,i}$ and depend on whether at BS-node n a link (n, l) or (l, n) is activated. BS-node n is connected to the BS-donor via multi-hop wireless links. Consequently, $T_{n,d,i}$ cannot be immediately observed when a link (n, l) , with $l \notin \mathcal{D}$ is activated. In fact, the destination BS-donor d might not even be known to BS-node n in time slot i . To overcome this limitation, we define the rewards $r_{n,i}$ as a function of the next-hop's estimated end-to-end latency $\hat{T}_{l,d,i}$ as

$$r_{n,i} = \begin{cases} t_{l,i}^q + t_{n,l,i}^{\text{tx}} + \hat{T}_{l,d,i}, & \text{for link } (n, l) \\ t_{n,i}^q + \hat{T}_{n,d,i}, & \text{for link } (l, n), \end{cases} \quad (3.4)$$

where $\hat{T}_{l,d,i}$ is calculated as $\hat{T}_{l,d,i} = \min_{(l,m) \in \mathcal{E}_i} \hat{T}_{l,m,i}$ and $t_{n,l,i}^{\text{tx}}$ is calculated based on $M_{n,l,i}$ to ensure the causality constraint in (3c) is fulfilled. Note that the

constraint in (3d) cannot be enforced, since multi-armed bandit algorithms learn from the activation of both optimal and suboptimal links.

3.2.3.2 Latency and CVaR Estimation

As given in (3.4), BS-node n learns which links (n, l) to activate by building estimates of the expected latency $\hat{T}_{n,l}$ associated to each of them. Let $K_{n,l,i} = \sum_{j=1}^i x_{n,l,j}$ be the number of times link (n, l) has been activated up to time slot i . The estimated $\hat{T}_{n,l}$ is updated using the sample mean as

$$\hat{T}_{n,l,i+1} = \frac{K_{n,l,i} \hat{T}_{n,l,i} + r_{n,i}}{K_{n,l,i} + 1}, \quad (3.5)$$

where the subindex i is introduced to emphasize that the estimate is built over time.

The CVaR definition given in (3.2) requires $T_{n,d}$ which, as discussed before, is not known a priori. Hence, we leverage the CVaR estimator derived in [140] to calculate the estimated CVaR of a link (n, l) . Let $\tilde{r}_n^1, \dots, \tilde{r}_n^{K_{n,l,i}}$ be all the rewards received up to time i . The estimated $\widehat{\text{CVaR}}_i(n, l)$ in time slot i is calculated as [140]

$$\widehat{\text{CVaR}}_i(n, l) := \inf_{t \in \mathbb{R}} \left(t + \frac{1}{\alpha \cdot K_{n,l,i}} \sum_{k=1}^{K_{n,l,i}} [\tilde{r}_n^k - t]^+ \right). \quad (3.6)$$

Using the estimates in (3.5) and (3.6), BS-node n computes the value $Q_n(a_{n,i} = (n, l))$ associated to the selected action $a_n \in \mathcal{A}_n$, and defined as

$$Q_n(a_{n,i}) = \hat{T}_{n,l,i} + \eta \widehat{\text{CVaR}}_i(n, l). \quad (3.7)$$

Note that (3.7) is aligned with the objective function in (3.3a). Actions with an associated low value $Q_n(a_{n,i})$ lead to lower end-to-end latency and a low expected value on its tail.

3.2.3.3 Consensus

All the BS-nodes independently decide which links to activate based on their estimates of the end-to-end latency. As a consequence, conflicting actions may be encountered. A conflict occurs when two or more BS-nodes n and m aim at activating a link to a common BS l , $l \in \mathcal{V}$, i.e., $x_{n,l,i} = x_{m,l,i} = 1$.

Algorithm 2 Safehaul algorithm at each BS-node

Input: $\alpha, \eta, \mathcal{A}_n$

- 1: Initialize $\hat{T}_{n,l}, \widehat{\text{CVaR}}(n, l)$, and Q_n for all $(n, l) \in \mathcal{E}_1$
 - 2: Set counters $K_{n,l} = 0$ and initial action $a_{n,1} = (n, n)$
 - 3: **for** every time slot $i = 1, \dots, I$ **do**
 - 4: perform action $a_{n,i}$, observe reward $r_{n,i}$ and increase counter $K_{n,l}$ by one ▷ Eq. (3.4)
 - 5: update latency estimate $\hat{T}_{n,l}$ ▷ Eq: (3.5)
 - 6: update CVaR estimate $\widehat{\text{CVaR}}(n, l)$ ▷ Eq: (3.6)
 - 7: update $Q_n(a_{n,i})$ ▷ Eq: (3.7)
 - 8: select next action $a_{n,i+1}$ using ϵ -greedy ▷ Eq. (3.8)
 - 9: share $a_{n,i+1}, t_{n,i}^q$ and $B_{n,i}$ with the other BS-nodes
 - 10: if required, update $a_{n,i+1}$ to reach consensus ▷ Sec. 3.2.3.3
 - 11: **end for**
-

We reach consensus by first retrieving the buffer and congestion status of the various IAB-nodes, leveraging the related BAP layer functionality [138, Sec. 4.7.3]. With this information at hand, conflicts are resolved by prioritizing the transmission of the BS-node with the larger queuing times $t_{n,i}^q$ and loads $B_{n,i}$. Then, we let the IAB-donor mark as *unavailable* the time resources of the remaining base stations with conflicting scheduling decisions [138, Sec. 10.9]. Note that as the learning is performed at each BS-node, only the link activation decision and the weighted sum of $t_{n,i}^q$ and $B_{n,i}$ are transmitted. Thus, low communication overhead is achieved.

3.2.3.4 Implementation of Safehaul

Here, we describe how the above-mentioned solution can be implemented in a real system. Specifically, we elaborate on the required inputs and the interactions among the different entities as well as the pseudo-code of Safehaul, see Alg. 2.

Safehaul is executed at each BS-node n . For its implementation, the MNO provides α, η and \mathcal{A}_n as an input. α is the risk level parameter that influences the level of reliability achieved in the system. Similarly, η controls the impact of the minimization of the latency in the α -worst cases on the overall performance. Both parameters, α and η , are set by the MNO depending on its own reliability requirements. The set \mathcal{A}_n depends on the considered network topology, which is perfectly known by the MNO. \mathcal{A}_n includes all links (n, l) and (l, n) to and from the first-hop neighbors of BS-node n .

The execution of Safehaul begins with the initialization of the latency and CVaR estimates, and the values Q of the actions in \mathcal{A}_n . Additionally, the

counters $K_{n,l}$, that support the calculations of $\hat{T}_{n,l}$ and $\widehat{\text{CVaR}}(n,l)$, are initialized for all links in \mathcal{A}_n (lines 1-2). These parameters are updated and learnt throughout the execution of Safehaul. At time slot $t = 0$, no transmission has occurred and $B_{n,0} = 0$. Hence, BS-node n remains idle for the first time slot $i = 1$, i.e., $a_{n,1} = (n, n)$ (line 2). Next, and in each of the subsequent time slots $i \in \{1, \dots, I\}$, the selected action is performed and the corresponding reward is obtained (line 4). If BS-node n transmits in time slot i , i.e., $a_{n,i} = (n, l)$, the reward $r_{n,i}$ is sent by the receiving BS l through the control channel. If $a_{n,i} = (l, n)$, the reward $r_{n,i}$ depends, as given in (3.4), only on the current estimates at BS-node n and the status of its buffer $B_{n,i}$. With the observed reward $r_{n,i}$, the counter for action $a_{n,i}$ is increased and the latency and CVaR estimates are updated (lines 4-6). Using the new estimates (lines 5 and 6), the value $Q(a_{n,i})$ of the performed action $a_{n,i}$ is updated (line 7). The next action $a_{n,i+1}$ is then selected according to ϵ -greedy (line 8), which is a well-known method to balance the exploitation of links with estimated low latency, and the exploration of unknown but potentially better ones. In ϵ -greedy, a random action $a_{n,i+1}$ from the set $\mathcal{A}_{n,i+1}$ is selected with probability $\epsilon \in [0, 1]$. With probability $(1 - \epsilon)$, instead, the action that yields the estimated lowest value is chosen, i.e.,

$$a_{n,i+1} = \begin{cases} \text{randomly selected action from } \mathcal{A}_{n,i+1}, & \text{if } x \leq \epsilon \\ \underset{b_n \in \mathcal{A}_{n,i+1}}{\operatorname{argmax}} Q_n(b_n), & \text{if } x > \epsilon, \end{cases} \quad (3.8)$$

where x is a sample taken from a uniform distribution in the interval $[0, 1]$. Once the action $a_{n,i+1}$ is selected, it is shared with other BS-nodes in the network along with $t_{n,i}^q$ and $B_{n,i}$ (line 9). As described in Section 3.2-3.3, this goes through the control channel. If conflicts arise, consensus is reached by prioritizing the transmission of the BS-node with the largest loads and queuing times (line 10).

3.2.3.5 Regret Analysis

The regret ζ is defined as the expected loss caused by the fact that the optimal action is not always selected [141]. Let \bar{T}^* and \bar{T}_{a_n} be the expected delay associated to the optimal action $a^* \in \mathcal{A}_n$ and the non-optimal action $a_n \in \mathcal{A}_n$, respectively. Similarly, let CVaR^* and CVaR_{a_n} be the CVaR of the optimal

action $a^* \in \mathcal{A}_n$ and the non-optimal action $a_n \in \mathcal{A}_n$, respectively. Formally, the regret ζ_i after i time slots is defined as

$$\begin{aligned}\zeta_i &= \sum_{a_n \in \mathcal{A}_n} ((\bar{T}_{a_n} + \eta \text{CVaR}_{a_n}) - (\bar{T}^* + \eta \text{CVaR}^*)) \mathbb{E}[K_{a_n, i}] \\ &= \sum_{a_n \in \mathcal{A}_n} \Delta_{a_n} \mathbb{E}[K_{a_n, i}],\end{aligned}\quad (3.9)$$

where $K_{a_n, i}$ is the number of times action a_n has been selected up to time slot i .

Proposition 1. *For a network \mathcal{G} in which the independent decisions of the BS-nodes do not lead to conflicts, let $A_n = |\mathcal{A}_n|$ be the number of available actions for BS-node n . Additionally, let $c > 0$, $0 < d \leq 1$, and $\epsilon_i := \min(1, \frac{cA_n}{d^2i})$. Then, there exists a positive constant $C > 1$, such that the probability that Safehaul chooses a non-optimal action $a_n \neq a^*$ after $i \geq cA_n/d$ time slots is upper bounded as*

$$\begin{aligned}\mathbb{P}[a_{n,i} = a_n] &\leq \frac{c}{d^2i} + \frac{4e}{d^2} B_i^{\frac{c}{2}} + \frac{2Cd^2}{c \ln\left(\frac{(i-1)d^2e^{0.5}}{cA_n}\right)} \\ &\quad + 4C\left(\frac{c}{d^2} \ln\left(\frac{(i-1)d^2e^{0.5}}{cA_n}\right)\right) B_i^{\frac{c}{5d^2}},\end{aligned}$$

with $B_i = \frac{cA_n}{(i-1)d^2e^{0.5}}$.

Proof. See the Appendix. ■

Theorem 1. *For a network \mathcal{G} in which the independent decisions of the BS-nodes do not lead to conflicts, the regret ζ_i of Safehaul after i time slots is upper bounded by*

$$\begin{aligned}\zeta_i &\leq \sum_{a_n \in \mathcal{A}_n} \Delta_{a_n} \left(1 + \sum_{i'=2}^i \left[\frac{c}{d^2i'} + \frac{4e}{d^2} B_{i'}^{\frac{c}{2}} + \frac{2Cd^2}{c \ln\left(\frac{(i'-1)d^2e^{0.5}}{cA_n}\right)} \right. \right. \\ &\quad \left. \left. + 4C\left(\frac{c}{d^2} \ln\left(\frac{(i'-1)d^2e^{0.5}}{cA_n}\right)\right) B_{i'}^{\frac{c}{5d^2}} \right], \right),\end{aligned}$$

where $c > 0$ and $0 < d \leq 1$.

Proof. From the definition in (3.9), the regret can be upper bounded as

$$\zeta_i \leq \sum_{a_n \in \mathcal{A}_n} \Delta_{a_n} \left(1 + \sum_{i'=2}^i \mathbb{P}[a_{n,i'} = a_n] \right), \quad (3.10)$$

by considering that $\mathbb{E}[K_{a_n,i}] \leq 1 + \sum_{i'=2}^i \mathbb{P}[a_{n,i'} = a_n]$. The bound is obtained by including the result of Proposition 1 in (3.10) as

$$\begin{aligned} \zeta_i &\leq \sum_{a_n \in \mathcal{A}_n} \Delta_{a_n} \left(1 + \sum_{i'=2}^i \left[\frac{c}{d^2 i'} + \frac{4e}{d^2} B_{i'}^{\frac{c}{2}} + \frac{2Cd^2}{c \ln\left(\frac{(i'-1)d^2 e^{0.5}}{c A_n}\right)} \right. \right. \\ &\quad \left. \left. + 4C \left(\frac{c}{d^2} \ln\left(\frac{(i'-1)d^2 e^{0.5}}{c A_n}\right) \right) B_{i'}^{\frac{c}{5d^2}} \right] \right), \end{aligned} \quad (3.11)$$

As every term in square brackets decreases monotonically in i' , the regret ζ_i grows sub-linearly. ■

3.2.4 Simulation setup

Given the lack of access to actual 5G (and beyond) network deployments, prior works mostly rely on *home-grown* simulators for performance evaluation. Although this is a valid approach, these simulators often cannot fully capture the real network dynamics, introducing strong assumptions in the physical and/or the upper layers of the protocol stack. Until very recently, the most complete simulator for IAB networks was a system-level simulator [142] developed as an extension of the ns-3 *mmWave* module [143]. However, despite accurate modeling of the IAB protocol stack, it is currently behind the latest IAB specifications². Moreover, the ns-3 IAB extension is unsuitable for large simulations with hundreds of nodes due to reliance on an older version of the *mmWave* module. Therefore, in our work we opt for Sionna [134], which is an open-source GPU-accelerated toolkit based on TensorFlow.

However, unlike the aforementioned ns-3 module, Sionna is a physical layer-focused simulator that does not explicitly model 5G networks, thus lacking the characterization of the 5G-NR upper-layer protocol stack. Hence, we extend Sionna by including the system-level functionalities such as MAC-level scheduling and RLC-level buffering. Furthermore, since Sionna exhibits slight differences compared to the 5G-NR physical layer, we extend Sionna's

²For instance due to the assumption of L-3 (instead of L-2) relaying at the IAB-nodes which was based on a draft version of TR 38.874 [144].

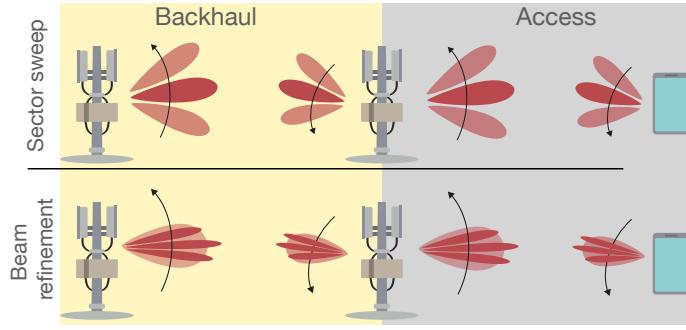


Figure 3.12: Schematic of the hierarchical beam management procedure. First, the general direction is estimated using wide beams (top). Then, the search is refined using the narrow beams codebook.

physical layer model [134] with the 5G-NR procedures. In the following, we describe the details of our extensions, which are publicly available³.

3.2.4.1 Extensions to Sionna's physical layer module

In this section, we describe the physical layer modification that were necessary to evaluate IAB scenarios using Sionna.

Codebook-based Beamforming Sionna's native beamforming only supports Zero-Forcing (ZF) pre-coding in downlink. Therefore, as a first step, we extend Sionna by implementing an NR-like codebook-based analog beamforming both at the transmitter and at the receiver. Specifically, we assume that the beamforming vectors at the transmitter w_{tx} and at the receiver w_{rx} are a pair of codewords selected from a predefined codebook. The codebook is computed by defining a set of beam directions $\{\omega_{p,q}\}$ which scans a given angular sector with a fixed beamwidth. The steering vector $a_{p,q}$ corresponding to direction $\omega_{p,q}$ can be computed as:

$$a_{p,q} = \left[1, \dots, e^{j \frac{2\pi}{\lambda} d (i_H \sin \alpha_p \sin \beta_q + i_V \cos \beta_q)}, \dots, e^{j \frac{2\pi}{\lambda} d ((N_H - 1) \sin \alpha_p \sin \beta_q + (N_V - 1) \cos \beta_q)} \right]^T, \quad (3.12)$$

where N_H and N_V are the number of horizontal and vertical antenna elements, respectively. The horizontal and vertical indices of a radiating element are denoted by $i_H \in \{0, \dots, N_H - 1\}$ and $i_V \in \{0, \dots, N_V - 1\}$, respec-

³https://github.com/TUDA-wise/safehaul_infocom2023

tively. α_p and β_q represent the azimuth and elevation angles of $\omega_{p,q}$. Next, we define the codebook as the set $\{(\sqrt{N_H N_V})^{-1} w_{p,q}\}$.

In line with the 5G-NR beam management procedure [49], we assume the lack of complete channel knowledge, i.e., the communication endpoints do not know the corresponding channel matrix. Accordingly, an exhaustive search is conducted to identify the best pair of codewords resulting in the highest SINR. We leverage a hierarchical search, in which the communication pairs first perform a wide-beam search in which the transmitter and the receiver approximate the direction of communication, see Figure 3.12. Next, the beamforming direction is fine-tuned through a beam refinement procedure going through a codebook with narrow beams. Consequently, we employ two types of codebooks, one with wide beams for sector sweep and another with narrow beams for beam refinement.

SINR Computations Since Sionna does not natively calculate the SINR, we add this functionality to the simulator to better model the impact of interference in our simulations. We compute the SINR experienced by TBs by combining the power of the intended signal with that of the interferers and of the thermal noise. Specifically, we first compute the power $P_n(i, f)$ of the intended signal at receiver n over frequency f and in time slot i . Then, we obtain the overall interference power by leveraging the superposition principle and summing the received power from all other interfering base stations $P_m(i, f)$ where $m \neq n$. For the purposes of this computation, we assume that each interferer employs the beamforming vector yielding the highest SNR towards its intended destination. Similarly, the transmitter and the receiver use the beamforming configuration estimated via the hierarchical search procedure. Finally, the SINR is

$$\gamma_n(i, f) = \frac{P_n(i, f)}{\sum_{m \neq n} P_m(i, f) + \sigma^2(i, f)}, \quad (3.13)$$

where $\sigma^2(i, f)$ is the thermal noise power at the receiver.

As mentioned, Sionna is mainly a physical layer simulator. However, to get closer to IAB networks as specified in Rel. 17, we have extended Sionna by implementing a selection of system-level features. To such end, we introduced a discrete-event network simulator for modeling IAB networks. This

3.2 Risk-averse learning for latency minimization in mmWave integrated access and backhaul networks

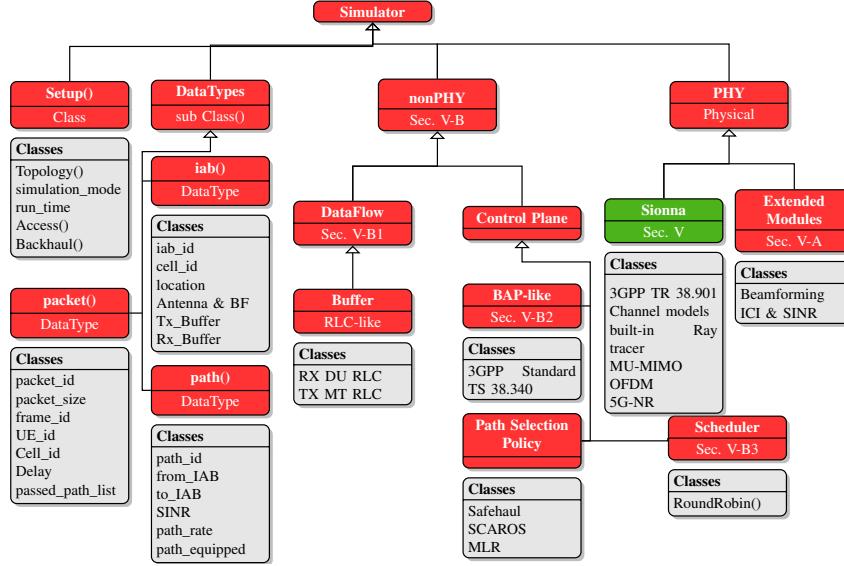


Figure 3.13: Overall design of our Sionna’s extension. The red blocks represent our additions to the baseline simulator, i.e., Sionna [134].

system-level extension operates on top of Sionna and provides basic functionalities such as a MAC-level scheduler, layer-2 buffers, and data flow and path selection mechanisms. Our simulator, depicted in Figure 3.13, generates a variety of system-level KPIs such as latency, throughput, and packet drop rate.

Data Flow and buffer 3GPP has opted for a layer-2 relaying architecture for BS-nodes where hop-by-hop RLC channels are established. This enables re-transmissions to take place on the affected hops only, thus preventing the need for traversing again the whole route from the BS-donor whenever a physical layer TB cannot be successfully decoded. This design results in a more efficient recovery from transmission failures and reduces buffering at the communication endpoints [145]. To mimic this architecture, we have implemented RLC-like buffers at each base station. Specifically, each BS-node features layer-2 buffers for both received and transmitted packets. For instance, the data flow for an uplink packet is the following. The UE generates packets and sends a transmission request to the base station. Consequently, the scheduler allocates OFDM symbols for this transmission, which is eventually received and stored at the RX buffer of its Distributed Unit (DU). Next, the packet is placed into the TX buffer to be forwarded to the suitable next

hop BS-node. This procedure is repeated until the packet crosses all the wireless-backhaul hops and reaches the BS-donor. Note that the packet can be dropped due to latency constraints or to interference.

BAP To manage routing within the wireless-backhauled network, the 3GPP introduced BAP, i.e., an adaptation layer above RLC which is responsible for packet forwarding between the BS-donor and the access BS-nodes [123]. Our simulator mimics this by associating each BS-node to a unique BAP ID. Moreover, we append a BAP routing ID to each packet at its entry point in the RAN (i.e., the BS-donor and the UEs for DL and UL data, respectively). Then, this identifier is used to discern the (possibly multiple) routes toward the packet’s intended destination [123]. The choice of the specific route is managed by Safehaul.

Scheduler Finally, we implemented a MAC-level scheduler which operates in a TDMA mode. The scheduler periodically allocates the time resources to backhaul or access transmissions in a Round-Robin fashion⁴. Specifically, each cell first estimates the number of OFDM symbols needed by each data flow by examining the corresponding buffer. Then, the subframe’s OFDM symbols are equally allocated to the users. If a user requires fewer symbols to transmit its complete buffer, the excess symbols (the difference between the available slot length and the needed slot length) are distributed to the other active users.

3.2.5 Performance Evaluation

In our simulations, we consider realistic cellular base station deployments in Manhattan, New York City⁵ and in the historical city center of Padova. Specifically, for the former we collect the locations of $N = 223$ 5G-NR base stations in an area of 15 Km^2 as depicted in Figure 3.14a. On the other hand, in the Padova topology we combine locations of $N = 100$ 4G-LTE Base Station (BS) of different MNOs (WINDTRE, TIM, and Vodafone) in an area

⁴The choice of the specific scheduling algorithm is outside of the scope of the 3GPP NR specifications, and is thus left to the MNOs. Accordingly, a Round-Robin scheduling policy represents a typical baseline assumption.

⁵The locations correspond to the network of T-Mobile, which has the largest deployment among the MNOs.

3.2 Risk-averse learning for latency minimization in mmWave integrated access and backhaul networks

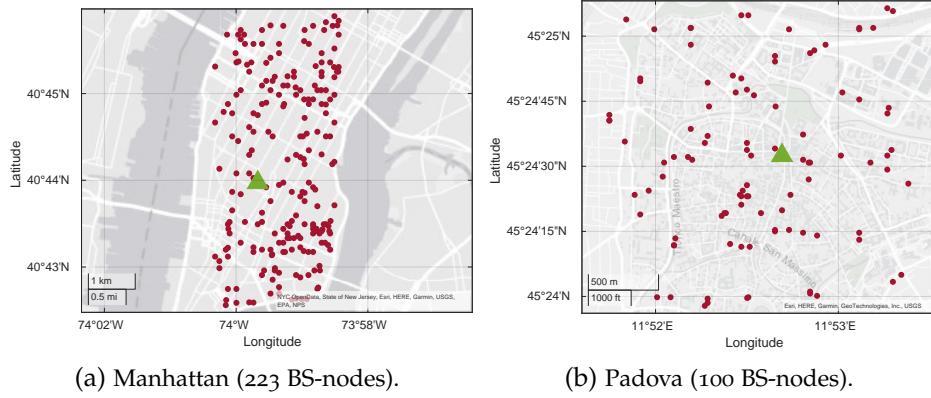


Figure 3.14: Locations of BS-nodes (red dots) and of the BS-donor (green triangle) in the Manhattan (left) and Padova (right) topologies.

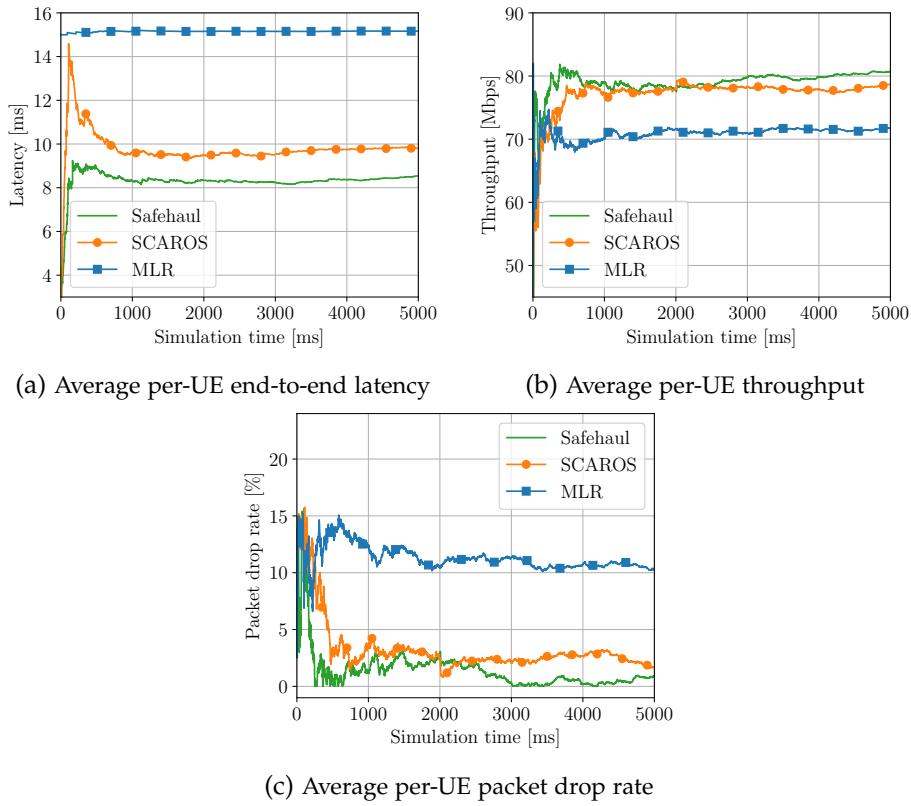


Figure 3.15: Average network performance for 50 UEs and 80 Mbps per-UE source rate (Scenario 1).

of 10 Km^2 as depicted in Figure 3.14b, due to the lack of 5G-NR base station deployment at the time of writing of this thesis. The detailed simulation

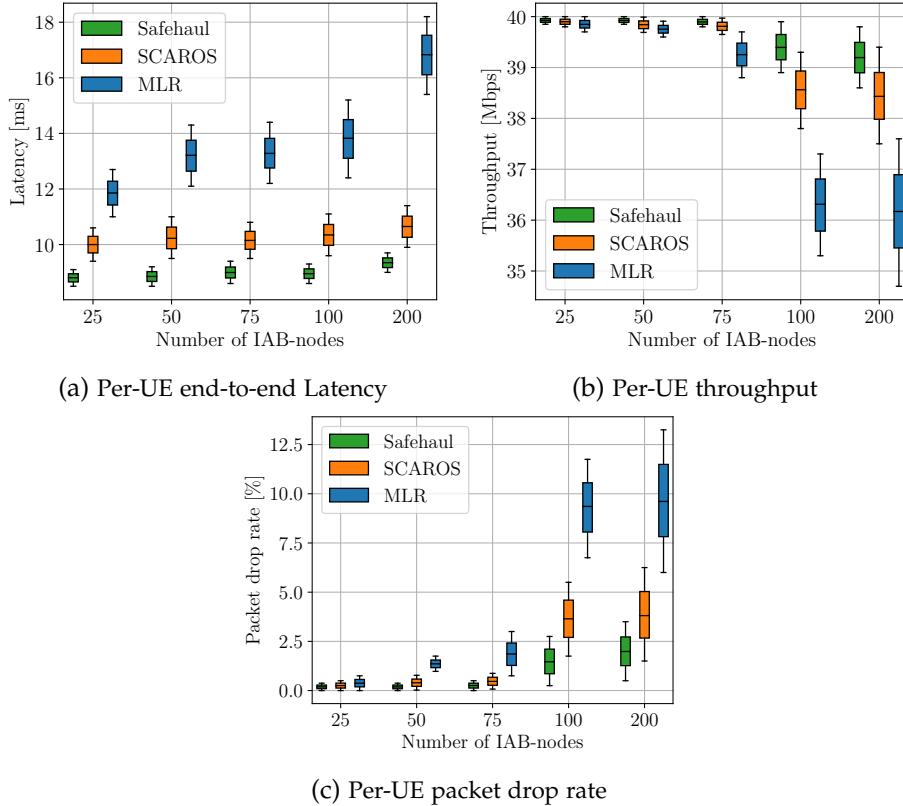


Figure 3.16: Network performance for $\{25, 50, 75, 100, 200\}$ BS-node, 2 UEs per BS-nodes on average, and 40 Mbps per-UE source rate (Scenario 2).

parameters are provided in Table 3.1. We used the channel model outlined by 3GPP in TR 38.901 [23], which provides a statistical channel model for 0.5-100 GHz, and analyzed the "Urban Micro (UMi)-StreetCanyon" scenario.

Benchmarks. To provide better insights on the performance of Safehaul, we replicate two approaches from the state of the art: (i) Scalable and Robust Self-backhauling Solution (SCAROS), a learning-based approach that minimizes the average latency in the network [131], and (ii) Maximum-local-rate (MLR), a greedy approach aiming to maximize throughput by selecting the links with the highest data rate.

Our evaluation consists of six scenarios, in which we study the convergence of the algorithms to a steady state, the number of BS-nodes, the number of BS-donors, and the impact of risk aversion. When demonstrating the results, we show the average throughput, latency, and packet drop rate per UE. Furthermore, we show the statistical variance of the obtained results using

Table 3.1: Simulation parameters.

Parameter	Value
Carrier frequency and bandwidth	28 GHz and 400 MHz
IAB RF chains	2 (1 access + 1 backhaul)
Pathloss model	UMi-Street Canyon [23]
Number of BS-nodes N	{223 NY, 100 Padova}
Source rate	{40, 80} Mbps
IAB Backhaul and access antenna array	8H×8V and 4H×4V
UE antenna array	4H×4V
IAB and UE height	15 m and 1.5 m
IAB antenna gain	33 dB
Noise Figure	10 dB
Risk level α	0.1
Reliability weight factor η	1

candlesticks which include the corresponding max, min, mean, and 10 and 90 percentiles.

3.2.5.1 Scenario 1: Average Network Performance

Analyzing the performance of the algorithms as a function of time is crucial to determine the convergence speed of the learning-based techniques, i.e., Safehaul and SCAROS. Hence, in Figure 3.15 we show the average network performance over time for three metrics: latency, throughput, and packet drop rate.

In Figure 3.15a, we can observe that Safehaul rapidly converges to an average latency of approximately 8.6 ms which is 12.2% and 43.4% lower than the latency of SCAROS and MLR, respectively. The high performance of Safehaul stems from the joint minimization of the average latency and the expected value of its tail loss, which results in avoiding risky situations where latency goes beyond T_{\max} . This is not the case for SCAROS where we observe a high peak in the latency before convergence, i.e., between zero and 1000 ms. *It is exactly the avoidance of such transients in Safehaul that leads to higher reliability in the system.* The reliability offered by Safehaul allows MNOs to deploy self-backhauling in an online fashion and without disrupting the network operation. The performance of MLR is constant throughout the simulation, as it is not designed as an adaptive algorithm.

Figure 3.15b shows that the risk-aversion capabilities of Safehaul have no negative impact on the average throughput of the network. The performance of Safehaul is comparable to that of SCAROS, approximately 79.3 Mbps, and 11.7% larger than the performance of MLR. The performance shown in Figure 3.15c is consistent with the behavior observed in Figure 3.15a. As Safehaul additionally minimizes the α -worst latency, it achieves the lowest packet drop rate compared to the reference schemes, namely, 30.1% (84.0%) lower than SCAROS (MLR).

3.2.5.2 Scenario 2: Impact of the Network Size

In Figure 3.16 we evaluate the reliability of the three considered approaches for different network sizes. Specifically, we vary the number of BS-nodes from 25 to 200. At the same time, we increase the load in the network by increasing the number of UEs. From the figures, we can clearly see that Safehaul consistently achieves a lower variation compared to the reference schemes. This verifies that Safehaul achieves the intended optimization goal, i.e., the joint minimization of the average end-to-end delay and its expected tail loss.

Figure 3.16a shows that Safehaul is able to maintain an almost constant latency as the number of BS-nodes increases. Specifically, the variation of latency with Safehaul is 56.1% and 71.4% less than with SCAROS and MLR, respectively. Furthermore, Safehaul achieves 11.1% and 43.2% lower latency compared to SCAROS and MLR, where the high variance exhibited by the latter is due to a lack of adaptation capabilities.

As shown in Figure 3.16b, the average throughput of the learning-based approaches Safehaul and SCAROS remains constant for the different values of the network size. However, the lowest variation in the throughput is achieved by Safehaul, i.e., only 0.90 compared to 1.9 and 2.8 in the benchmark schemes. Such behavior corroborates Safehaul's reliability capabilities.

The packet drop rate for different numbers of BS-nodes is shown in Figure 3.16c. Safehaul not only consistently outperforms the reference schemes, but also has the minimum variation in the results (at least 47.3% lower compared to the benchmarks). Considering the largest network size and load, i.e., 200 BS-nodes and 400 UEs, Safehaul achieves 49.3% and 81.2% lower packet drop rate compared to SCAROS and MLR, respectively.

3.2 Risk-averse learning for latency minimization in mmWave integrated access and backhaul networks

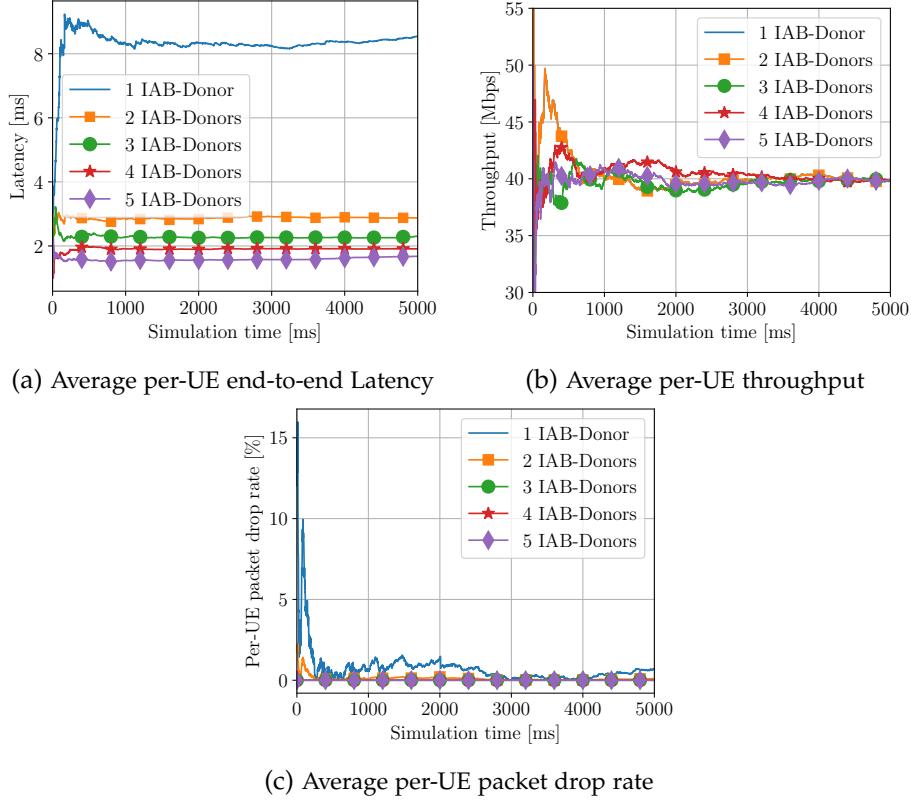


Figure 3.17: Network performance for 50 UEs and 40 Mbps per-UE source rate, versus the number of BS-donors (Scenario 3).

3.2.5.3 Scenario 3: Impact of the number of BS-donors

Although the benchmark schemes do not support multiple BS-donors, Safehaul is designed to accommodate such scenarios. In Figure 3.17, we investigate the impact of the number of BS-nodes on Safehaul. To this end, we keep the number of UEs and their data rate constant. We observe in Figure 3.17a that the highest latency is experienced when only one BS-donor is present in the network. This stems from the tributary effect of self-backhauling where the traffic flows towards a central entity which itself can become a bottleneck. As the number of BS-donors increases, the traffic is more evenly distributed, resulting in lower latency. Specifically, the average latency decreases from 8.2 ms for $D = 1$ to 1.7 ms when $D = 5$. As mentioned, since the load is constant in this scenario, the average throughput also remains constant for all different numbers of BS-donors, see Figure 3.17b. Notably, Safehaul's learning speed is maintained for the different values of D . This is an impor-

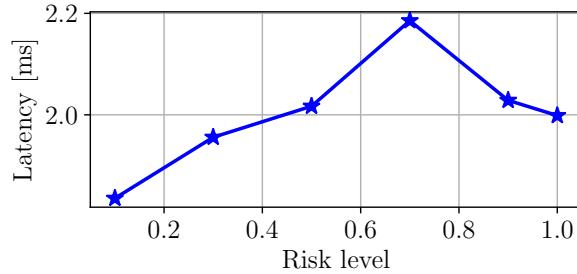


Figure 3.18: Average latency for 50 UEs and 20 Mbps per-UE source rate, versus the risk level α (Scenario 4)

tant design feature of Safehaul because having more BS-donors means that the number of paths a BS-node has to the core network increases exponentially. From a learning perspective, such increment implies a larger action set and a lower learning speed. Safehaul avoids this problem by learning the average latency based on the estimates of its neighbors and not on the complete paths to the BS-donors. Finally, Figure 3.17c shows that increasing the number of BS-donors significantly reduces the packet drops, which also stems from a better distribution of traffic flows in the network, as observed in Figure 3.17a.

3.2.5.4 Scenario 4: Impact of the risk parameter α

The definition of losses in the tail of the latency distribution is controlled by the risk level parameter α . Its impact on the average latency is shown in Figure 3.18, where an increasing behavior is observed for $\alpha \leq 0.7$. The lowest latency is achieved for $\alpha = 0.1$, which corresponds to the most risk-averse, and therefore the most reliable, case out of all the considered ones. The non-monotonic behavior of the average latency versus α can be explained by the so-called exploration-exploitation trade-off: the higher α , the higher the level of risk, which in turn leads Safehaul to learn more about the environment and choose a more reliable action. Eventually, as α grows beyond approximately 0.7, the performance of Safehaul tends to that of the risk-neutral case. As a consequence, the algorithm undertakes excessive exploration, which causes a degradation of the average latency performance.

3.2 Risk-averse learning for latency minimization in mmWave integrated access and backhaul networks

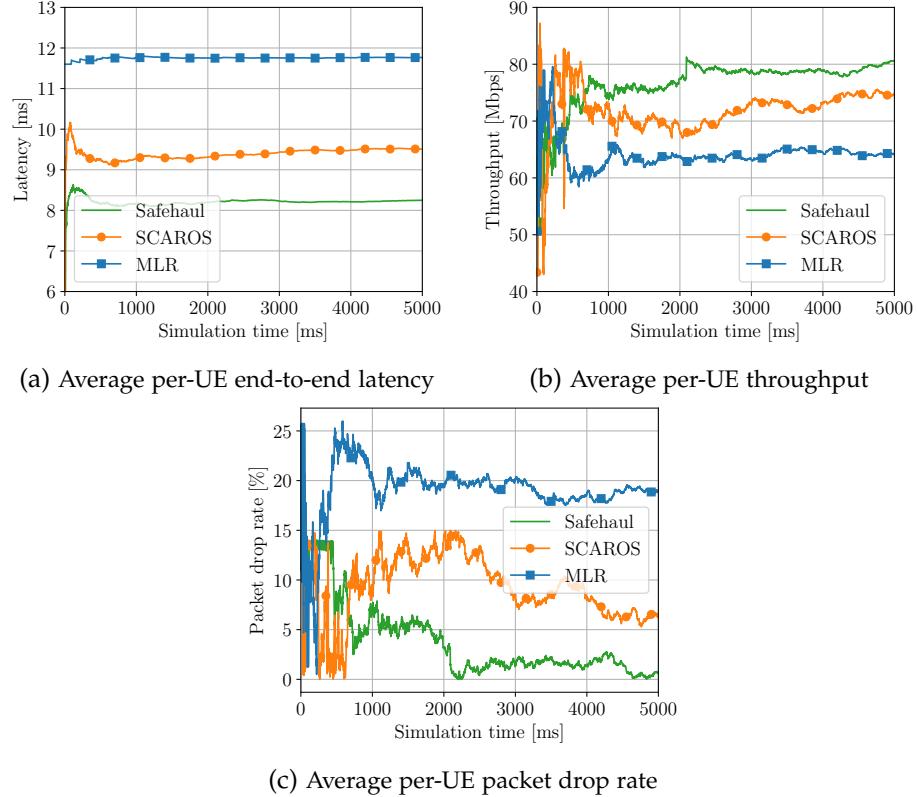


Figure 3.19: Average network performance for 50 UEs and 80 Mbps per-UE source rate (Scenario 1) in Padova.

3.2.5.5 Scenario 5: Performance in different topologies

To verify the generality of the proposed algorithms, it is essential to examine how they perform in different topologies, and consider both typical network performance metrics (i.e., along the lines of Scenario 1) and their stability with respect to the number of BS-nodes and BS-donors (Scenarios 2 and 3). To this end, we ran additional simulations in the deployment depicted in Figure 3.14b, which mimics the BS-nodes locations of the historic center of Padova. We report the average network performance over time, in terms of end-to-end packet drop rate, throughput, and latency in Figure 3.19. Overall, the outcomes of this simulation campaign are in line with those obtained in Scenario 1. Specifically, as seen in Figure 3.19a, Safehaul quickly converges to an average latency of approximately 8 ms, which is 14% and 31% lower than SCAROS and MLR's latency. Figure 3.19b shows the average per-UE throughput, for which Safehaul achieves about 4% and 17% better perfor-

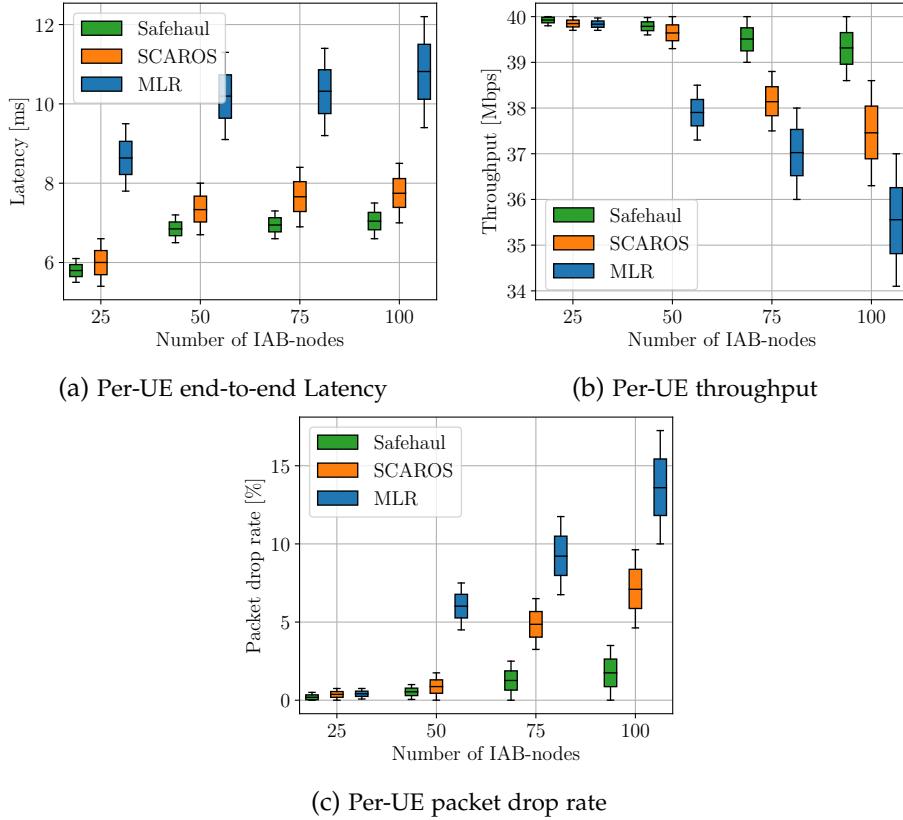


Figure 3.20: Network performance for $\{25, 50, 75, 100\}$ BS-nodes, 2 UEs per BS-node on average, and 40 Mbps per-UE source rate (Scenario 2) in Padova.

mance than SCAROS and MLR, respectively. Similarly, the performance depicted in Figure 3.19c is in line with that reported in Figs. 3.19a and 3.19b, with Safehaul achieving approximately a 24% and 38% smaller packet drop rate than SCAROS and MLR, respectively.

In Figure 3.20, we compare the consistency of the performance of the three algorithms with respect to the network size. In particular, we change the number of BS-nodes from 25 to 100, keeping fixed the number of UEs per BS-node and thus effectively increasing the network load on the BS-donor. Results show that Safehaul, when compared to other schemes, exhibits minimal performance degradation when introducing additional BS-nodes and UEs. As can be seen in Figure 3.20a, the latency achieved by Safehaul increases by at most 16% in the case of 100 BS-nodes, while SCAROS and MLR lead to a latency which is consistently higher and increases up to 27% and 25% when deploying additional nodes, respectively. Similar trends can be

3.2 Risk-averse learning for latency minimization in mmWave integrated access and backhaul networks

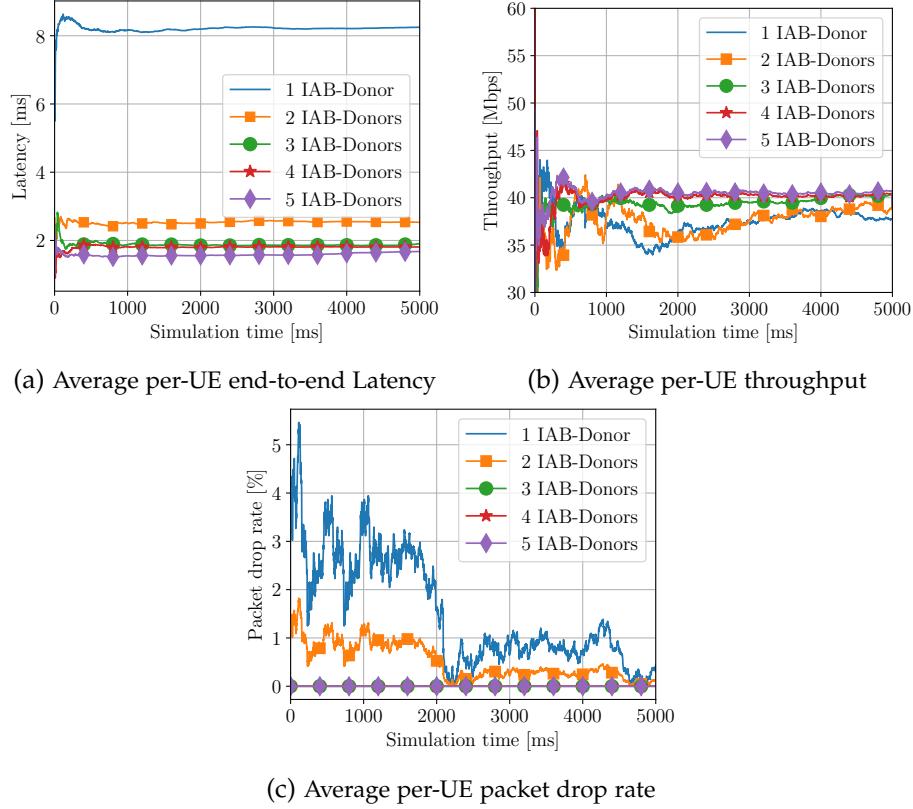


Figure 3.21: Network performance for 50 UEs and 40 Mbps per-UE source rate, versus the number of BS-donors in Padova (Scenario 3).

observed in Figs. 3.20b and 3.20c, which report throughput and packet loss versus the network size, respectively. Indeed, Safehaul is the best performer across the whole range of BS-nodes which have been considered. Furthermore, Safehaul loses 20% more packets with the denser network deployment (i.e., 100 BS-nodes), while reference schemes exhibit an increase in packet loss of up to 33%.

We complete this analysis by examining how the number of donors affects the performance achieved by Safehaul in the Padova-like topology. As can be seen in Figure 3.21, increasing the number of fiber-backhauled base stations progressively reduces the latency. Similarly, and in line with the results obtained in Scenario 3 and reported in Figure 3.21c, the packet drop rate varies from approximately 0.08% when considering a single BS-donor, to approximately 0.003% in the presence of five BS-donors. The performance

improvements introduced by additional fiber links saturate after 3 donors, thanks to the efficient routing and scheduling performed by Safehaul.

In summary, the results obtained in the additional topology mimicking the historical center of Padova are well aligned with those obtained in the Manhattan topology. Although the specific values of the network metrics achieved by the considered schemes in the two topologies are different (for instance, SCAROS achieves a 66% lower packet drop rate in Scenario 1 compared to Scenario 5), the trends among the various schemes are the same. Specifically, we observed that Safehaul consistently achieves the best performance in comparison to SCAROS and MLR across different metrics, which supports the claim that the proposed scheduler is capable of learning how to optimize arbitrary deployment topologies.

3.2.5.6 Scenario 6: Network resilience

In networking, resilience refers to the ability of a network to recover in a quick and effective fashion from disruptions, thus providing reliable and high-quality communication services to its users. Specifically, the ability to recover from link failures is particularly important in IAB networks, where backhaul links are susceptible to the typical disruptions which plague the RAN due to its mobile and wireless nature. For instance, the links among BS-nodes can be degraded by adverse environmental conditions such as heavy rain and monsoons, physical obstacles and network congestion. These disruptions can cause temporary or permanent communication failures, which in turn result in degraded performance and/or loss of connectivity for the end users. To prevent and/or recover from these undesired events, a backhaul scheduler must detect, mitigate, and recover from various types of disruptions and failures, and must maintain the required level of service availability and performance despite the time-varying channel conditions.

We benchmark the resilience of the proposed algorithm by mimicking radio link failures, which we simulate by stopping BS-nodes at a fixed time instant (2000 s), and inspecting the resulting performance degradation. Since the failed node(s) is (are) chosen at random, we run multiple simulations to estimate the average network performance, as shown in Figs. 3.22 and 3.23 for the case of one and three link failures, respectively.

3.2 Risk-averse learning for latency minimization in mmWave integrated access and backhaul networks

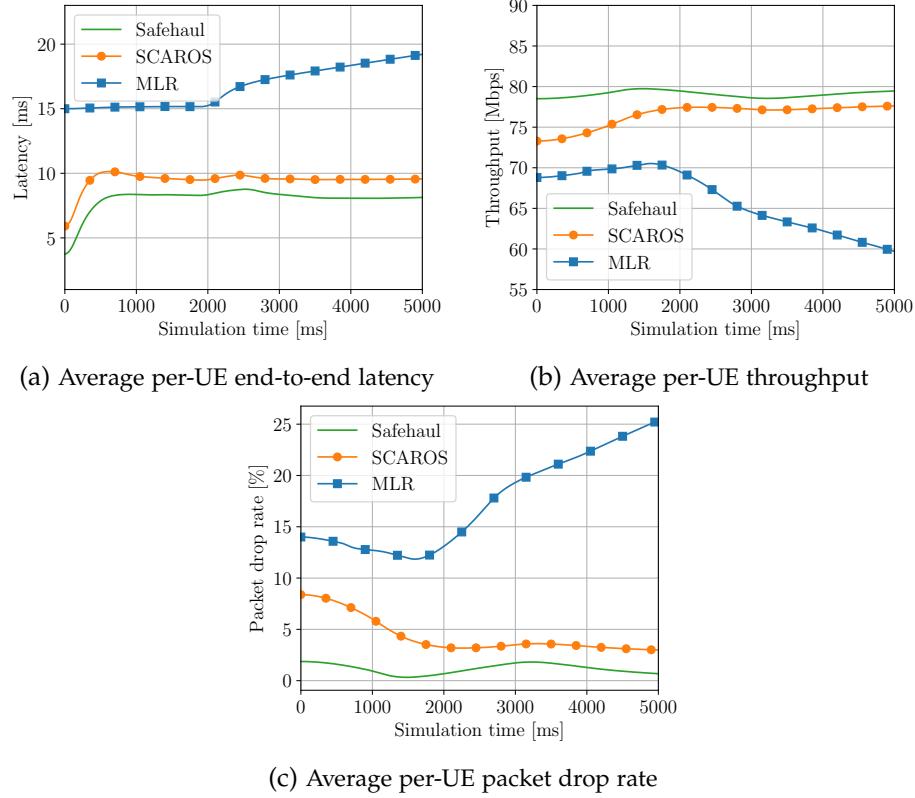


Figure 3.22: Average network performance for 50 UEs and 80 Mbps per-UE source rate where 1 random BS-node is shut down.

Results show that MLR is unable to react to the link failure(s) due to its static and myopic policy. Specifically, the disruption causes an increase of 33% (60%) in latency, and a decrease of up to 15% (23%) in throughput when considering one (three) link failure(s). On the other hand, both Safehaul and SCAROS are capable of adapting the scheduling to the new topology. Indeed, both schemes show a transient region where the performance is slightly degraded since the algorithms are learning new routes and resource partitions to account for the lost link. Nevertheless, Safehaul and SCAROS eventually converge to a solution which provides approximately the same network performance as before the failures, in both cases of one and three lost links.

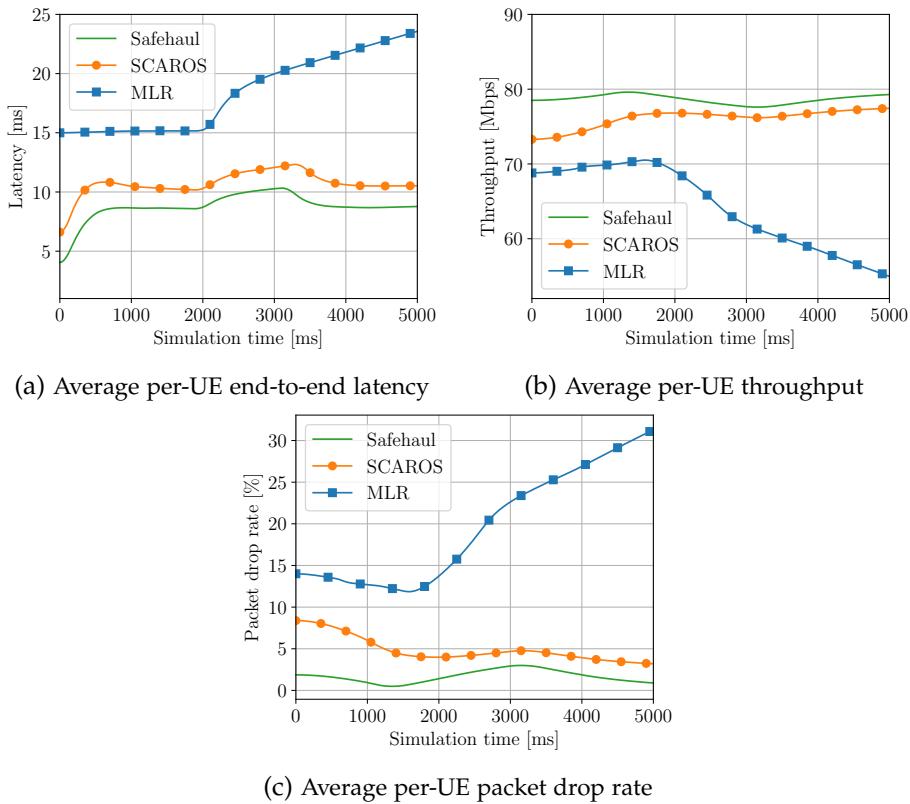


Figure 3.23: Average network performance for 50 UEs and 80 Mbps per-UE source rate where 3 random BS-nodes are shut down.

3.2.6 Related work

Self-backhauling wireless networks have been studied in different contexts. Ranging from the so-called Heterogeneous Networks (HetNets) and IAB 5G NR systems, to Cloud Radio Access Networks (RANs), each has considered a different set of premises and optimization goals. In this section, we review the related work in the context of basic assumptions and their optimization goals.

Ideal backhaul links. Numerous works assume either an *infinite or fixed capacity backhaul link*. This is often motivated by the presence of a wired fiber link between the Small Base Stations (SBSs) and the Macro Base Station (MBS) [103, 125, 127, 128]. Indeed, most of these works consider a scenario where a centralized Baseband Unit (BBU) is connected to several Remote Radio Heads (RRHs), i.e., radios which lack signal processing capabilities [125, 127, 128]. In particular, the authors of [128] consider an even

more complex C-RAN scenario where RRHs feature caching and signal processing capabilities. However, in an IAB context it is fundamental to consider *limited-rate, time-varying backhaul channels* and to study the impact of such limitations on the performance of the RAN.

Constrained topologies. It is often assumed that self-backhauled networks have a *specific topology*. This assumption usually simplifies the problem and makes it tractable and/or solvable in polynomial time. For instance, the authors of [102, 129, 130] assume a single-hop network where each SBS is directly connected to the MBS. In [101], a k -ring deployment is considered, i.e., a topology where a single IAB-donor provides backhaul connectivity to k rings of IAB-nodes. Even though this topology can be used to model networks with arbitrary depth, it maintains a symmetric load for each node, an assumption which generally does not hold in real networks. In fact, the 3GPP does not impose any limits on the number of IAB-nodes which can be connected to a given IAB-donor, nor does it set an upper bound on the number of wireless hops from the latter to other wireless-backhauled base stations [117]. Accordingly, in our problem formulation we consider IAB networks with an *arbitrary number of nodes and an arbitrary maximum number of wireless hops* between MBSs and SBSs.

Simplistic traffic models. Some works either assume a *full buffer traffic model and/or impose flow conservation constraints*. In particular, the authors of [103, 146] consider systems where the capacity of each link can always be fully exploited thanks to the presence of *infinite data to transmit at each node*. However, in actual IAB deployments the presence of packets at the MBSs and SBSs is conditioned on the *status of their RLC buffers and, in turn, on the previous scheduling decisions*. Moreover, *packets can actually be buffered at the intermediate nodes*, thus preventing the need for transmitting a given packet in consecutive time instants along the whole route from the IAB-donor to the UEs (or vice versa).

Optimization goals. The works in the literature focus on different optimization goals. Therefore, they prioritize different network metrics. For instance, the authors of [147] aim to optimize the beam alignment between MBSs and SBSs. Instead, the work of [126] aims to compute the optimal user-to-base-station association. However, they neglect backhaul associations and focus on the access only. In [126, 129, 146, 148] the objective function is a function of the users data-rate. In particular, the authors of [146] optimize

the max-min user throughput, arguing that such a metric better captures the performance of the bottleneck links. In [109], the average rate of each link is maximized under bounded delay constraints. In our work, we focus on reliability by minimizing not only the average end-to-end delay, but also the expected value of the worst-case performance. The work closest to this article is SCAROS [131], a learning-based latency-aware scheme for resource allocation and path selection in self-backhauled networks. Assuming a single IAB-donor, the authors study arbitrary multi-tier IAB networks considering the impact of interference and network dynamics. In contrast, we aim at enhancing the reliability of the IAB-network by jointly minimizing the average end-to-end delay and its expected tail loss.

Appendix

For the proof of Proposition 1, Theorem 3 in [140] is needed. For completeness, we first present the theorem in [140] for the special case in which the considered random variables are independent. Next, we present the proof of Proposition 1.

Theorem 2. *Let $T_{a_n,i}$ be independent random variables where $\max_{1 \leq j \leq i} T_{a_n,j} = T_{\max}$, with $i \in \{1, 2, \dots\}$. Then, for any $0 < \delta \leq 1/2$, $\xi > 0$ and $\gamma > 0$, there exists a positive constant C which only depends on ξ and γ , such that the probability of the event $|\widehat{\text{CVaR}}_{a_n,i} - \text{CVaR}_{a_n,i}| \geq 2\xi\alpha^{-1}T_{\max}i^{-\delta}(\ln \ln i)^{1/2} \ln i$ is smaller than or equal to $Ce^{-(1+\gamma) \ln i}$.*

Proof. See Theorem 3 in [140]. ■

Proof of Proposition 1

In this proof, we use the result of the regret bound for the risk-neutral case without CVaR, shown in [141, Theorem 3], as a basis. Additionally, we use the bound for the terms related to the CVaR formulated in [140, Theorem 3]. Using both these results, we first bound the probability that Safehaul chooses a suboptimal arm in the exploitation phase. Then, we combine the latter with the probability of choosing a suboptimal arm in the exploration phase to derive the bound given in Proposition 1.

From the system model and Proposition 1, we have that $c > 0$, $0 < d \leq 1$, and $\epsilon_n := \min(1, \frac{cA_n}{d^2i})$. Moreover, $a_{n,i}$ is the action chosen by ϵ -greedy in time

slot i and $K_{a_n,i}$ is the number of times, up to time slot i , in which Safehaul chose action a_n at random. Similarly, we use K_i^* for the counter of the optimal action. $T_{a_n,i}$ are independent random variables distributed according to the rewards linked to action a_n . We use T_i^* for the optimal action, and $\hat{T}_{a_n,i}$ is the estimated mean of the probability distribution of the rewards linked to action a_n using $K_{a_n,i}$ samples. As before, we use \hat{T}_i^* for the optimal action. $\widehat{\text{CVaR}}_{a_n,i}$ is the estimated CVaR of action a_n up to time slot i and $\widehat{\text{CVaR}}_i^*$ is the estimated CVaR of the optimal action up to time slot i . Then, the probability that action a_n is chosen in time slot i is upper bounded as

$$\mathbb{P}[a_{n,i} = a_n] \leq \mathbb{P}[\delta_{a_n,i-1} \leq \delta_{i-1}^*] \left(1 - \frac{\epsilon_i}{A_n}\right) + \frac{\epsilon_i}{A_n}, \quad (3.14)$$

with $\delta_{a_n,i-1} = \hat{T}_{a_n,i-1} + \eta \widehat{\text{CVaR}}_{a_n,i-1}$ and $\delta_{i-1}^* = \hat{T}_{i-1}^* + \eta \widehat{\text{CVaR}}_{i-1}^*$. The first term in (3.14) is the probability of exploitation and the second term to the probability of exploration.

Using the mean \bar{T}_{a_n} and CVaR_{a_n} of action a_n , and the likewise defined \bar{T}^* and CVaR^* for the optimal action, we set $\Delta_{a_n}^{\text{mean}} := \bar{T}_{a_n} - \bar{T}^*$ and $\Delta_{a_n}^{\text{cvvar}} := \text{CVaR}_{a_n} - \text{CVaR}^*$. Using these definitions in (3.14) we conclude

$$\begin{aligned} & \mathbb{P}[\delta_{a_n,i-1} \leq \delta_{i-1}^*] \leq \\ & \mathbb{P}\left[\delta_{a_n,i-1} \leq \eta \text{CVaR}_{a_n} - \frac{\Delta_{a_n}^{\text{mean}}}{2} + \bar{T}_{a_n} - \eta \frac{\Delta_{a_n}^{\text{cvvar}}}{2}\right] + \\ & \mathbb{P}\left[\bar{T}^* + \frac{\Delta_{a_n}^{\text{mean}}}{2} + \eta \text{CVaR}^* + \eta \frac{\Delta_{a_n}^{\text{cvvar}}}{2} \leq \delta_{i-1}^*\right] \\ & \mathbb{P}\left[\hat{T}_{a_n,i-1} \leq \bar{T}_{a_n} - \frac{\Delta_{a_n}^{\text{mean}}}{2}\right] + \mathbb{P}\left[\bar{T}^* + \frac{\Delta_{a_n}^{\text{mean}}}{2} \leq \hat{T}_{i-1}^*\right] + \\ & \mathbb{P}\left[\widehat{\text{CVaR}}_{a_n,i-1} \leq \text{CVaR}_{a_n} - \frac{\Delta_{a_n}^{\text{cvvar}}}{2}\right] + \\ & \mathbb{P}\left[\text{CVaR}^* + \frac{\Delta_{a_n}^{\text{cvvar}}}{2} \leq \widehat{\text{CVaR}}_{i-1}^*\right]. \end{aligned} \quad (3.15)$$

Similar to [141], we use the Chernoff-Hoeffding bound for the first two terms in (3.15). For the last two summands, there remains to find a bound for the difference between the CVaR and its estimate $\widehat{\text{CVaR}}$. From Theorem 2, we set $\xi := \Delta_{a_n}^{\text{cvvar}} \alpha / 4T_{\max}$, $\delta = 0.5$ and by using the limit $\gamma \rightarrow 0$, we obtain

$$\mathbb{P}\left[|\widehat{\text{CVaR}}_{a_n,i} - \text{CVaR}_{a_n,i}| \geq \frac{\Delta_{a_n}^{\text{cvvar}}}{2} i^{-0.5} (\ln \ln i)^{0.5} \ln i\right] \leq \frac{C}{i}. \quad (3.16)$$

As $\max_i i^{-0.5}(\ln \ln i)^{0.5} \ln i < 1$, the condition $(\Delta_{a_n}^{\text{cvar}}/2)i^{-0.5}(\ln \ln i)^{0.5} \ln i \leq \frac{\Delta_{a_n}^{\text{cvar}}}{2}$ holds for all i . Therefore, considering the last two summands in (3.15), we conclude that there exists a positive constant C that satisfies

$$\mathbb{P}\left[|\widehat{\text{CVaR}}_{a_n,i} - \text{CVaR}_{a_n,i}| \geq \frac{\Delta_{a_n}^{\text{cvar}}}{2}\right] \leq \frac{C}{i}. \quad (3.17)$$

The number of times action a_n has been selected up to time slot i is smaller than or equal to i , i.e., $K_{a_n,i} \leq i$. Using (3.17) we write the last two summands in (3.15) as

$$\mathbb{P}\left[\widehat{\text{CVaR}}_{a_n,i-1} \leq \text{CVaR}_{a_n} - \frac{\Delta_{a_n}^{\text{cvar}}}{2}\right] \leq \frac{C}{K_{a_n,i-1}}, \quad (3.18)$$

and

$$\mathbb{P}\left[\text{CVaR}^* + \frac{\Delta_{a_n}^{\text{cvar}}}{2} \leq \widehat{\text{CVaR}}_{i-1}^*\right] \leq \frac{C}{K_{i-1}^*}. \quad (3.19)$$

As in [141], we use Bernstein's inequality to get an estimate for $K_{a_n,i-1}$. Defining $x_0 := 1/2A_n \sum_{j=1}^{i-1} \epsilon_j$ for $i-1 \geq \frac{cA_n}{d^2}$ we get $P(K_{a_n,i-1} \leq x_0) \leq e^{-\frac{x_0}{5}}$. Additionally, from [141]:

$$x_0 \geq \frac{c}{d^2} \ln\left(\frac{(i-1)d^2e^{0.5}}{cA_n}\right) =: C'(i). \quad (3.20)$$

The same holds for the optimal action and K_{i-1}^* . Using these estimations for x_0 , we can conclude that for $i-1 \geq cA_n/d^2$

$$\mathbb{P}\left[\widehat{\text{CVaR}}_{a_n,i-1} \leq \text{CVaR}_{a_n} - \frac{\Delta_{a_n}^{\text{cvar}}}{2}\right] \leq \sum_{j=1}^{i-1} \mathbb{P}[K_{a_n,i-1} = j] \frac{C}{j} \quad (3.21)$$

$$\begin{aligned} &= \sum_{j=1}^{\lfloor x_0 \rfloor} \mathbb{P}[K_{a_n,i-1} = j] \frac{C}{j} + \sum_{j=\lfloor x_0 \rfloor + 1}^{i-1} \mathbb{P}[K_{a_n,i-1} = j] \frac{C}{j} \\ &\leq Cx_0 e^{-\frac{x_0}{5}} + \frac{C}{x_0} \leq Cx_0 e^{-\frac{x_0}{5}} + \frac{C}{C'(i)}. \end{aligned} \quad (3.22)$$

The same holds again for the optimal action

$$\mathbb{P}\left[\text{CVaR}^* + \frac{\Delta_{a_n}^{\text{cvar}}}{2} \leq \widehat{\text{CVaR}}_{i-1}^*\right] \leq Cx_0 e^{-\frac{x_0}{5}} + \frac{C}{C'(i)}. \quad (3.23)$$

3.2 Risk-averse learning for latency minimization in mmWave integrated access and backhaul networks

Together with the bounds from Theorem 3 in [141] it follows that for $C \geq 1$:

$$\begin{aligned}
& \mathbb{P}[a_{n,i} = a_n] \\
& \leq \frac{\epsilon_i}{A_n} + 4Cx_0e^{-\frac{x_0}{5}} + \frac{4}{(\Delta_{a_n}^{\text{mean}})^2} e^{-\frac{(\Delta_{a_n}^{\text{mean}})^2 \lfloor x_0 \rfloor}{2}} + 2\frac{C}{C'(n)} \\
& \leq \frac{c}{d^2 i} + 2\frac{Cd^2}{c \ln\left(\frac{(i-1)d^2e^{0.5}}{cA_n}\right)} + \frac{4e}{d^2} \left(\frac{cA_n}{(i-1)d^2e^{0.5}}\right)^{\frac{c}{2}} + \\
& \quad 4C\frac{c}{d^2} \ln\left(\frac{(i-1)d^2e^{0.5}}{cA_n}\right) \left(\frac{cA_n}{(i-1)d^2e^{0.5}}\right)^{\frac{c}{5d^2}}.
\end{aligned}$$

■

3.3 High-capacity integrated access and backhaul networks using sub-terahertz links

The spectrum above 100 GHz has several sub-bands that could provide bandwidths wider than 10 GHz, thus potentially data rates in the excess of tens of Gbps [149]. Backhaul—a static deployment—is a promising use case for sub-terahertz links, which need pencil-sharp beams to close the link budget and are thus less resilient to mobility compared to traditional sub-6 GHz or mmWave frequencies.

In recent years, the literature has closed several gaps in terms of circuit, antenna design [150] and physical and MAC layer solutions for sub-terahertz systems [151]. When it comes to IAB with mixed sub-terahertz and mmWave links,⁶ however, there are still several open questions in terms of network design and path selection. In this work, we consider the problem of identifying a viable topology between IAB nodes and the IAB donors, including the carrier frequency of the backhaul links, and profile the performance that network planners can expect when mixing sub-terahertz and mmWave IAB links.

To this end, we develop a greedy path generation algorithm that automatically selects the frequency band of an IAB link (between 28 GHz and 140 GHz) and assigns routes so that each IAB node can reach the IAB donor. The frequency selection aims at avoiding bottlenecks, i.e., the algorithm selects the band that provides the highest capacity when accounting for the congestion that may arise in the proximity of the IAB donor. In addition, we consider and compare different ratios of sub-terahertz and mmWave links, which can be mapped to licensing constraints for out-of-band backhaul, and two different bandwidths for the sub-terahertz links (10 GHz and 32 GHz), which consider exclusive licensing or sharing with other services, respectively [152].

We model the IAB network in a custom-developed 3GPP Release 17 simulator based on the open-source tool Sionna [134], with 3GPP and state-of-the-art mmWave and sub-terahertz channel models, and realistic and detailed 3GPP-based physical and MAC layers. Our results quantify for the first time the performance improvement that sub-terahertz links can introduce in IAB networks, which can push beyond the limits of the in-band mmWave backhaul

⁶In this section, we consider the FR2 range of 3GPP NR (24.25 GHz to 71 GHz) as mmWaves.

and support more than 50 users with 120 Mbps streams and a single donor without congestion (compared to about 33 Mbps for in-band mmWaves).

This is the first work that provides a numerical evaluation of the potential associated with sub-terahertz links for IAB. Notably, [153] evaluates the sub-THz potential in backhaul networks from a physical layer perspective. This research demonstrates that sub-THz spectrum links can achieve multi-Gbps ratios in outdoor backhaul scenarios. [154] proposed UAV-assisted backhaul solution to improve network coverage and data rate in heterogeneous networks with multiple tiers composed of sub-6 GHz, THz and UAV layers. In addition, the authors of [155] successfully adopted concurrent scheduling to increase system throughput in dense THz backhaul scheduling. Finally, [156] considers a multi-band IAB deployment, but with a bandwidth that is more limited than those considered in future 6G scenarios.

The rest of the section is organized as follows. Sec. 3.3.1 introduces the system model. Sec. 3.3.2 describes the algorithm for frequency and path selection, which is then numerically evaluated in Sec. 3.3.3.

3.3.1 System Model

We consider a TDMA system in which a single IAB donor, featuring a fiber connectivity towards the CN and the Internet, exchanges data with N_U UEs. Without loss of generality, we consider uplink traffic only. To achieve uniform coverage, the donor is aided by N_I IAB nodes, which can be connected either to the former or to neighboring base stations, thus possibly realizing a multi-hop wireless backhaul.

We partition the time resources in T radio subframes of duration $T_{sub} = 1$ ms, and we equip all nodes with buffers. Accordingly, the data that node i transmits to gNB k during subframe t is stored in its buffer $B_k(t)$, and represents either successfully received packets, in the case of the donor, or data to be relayed to the next hop along the path during subframe $t + 1$, in the case of IAB nodes.

We assume that the backhaul links operate *either in the mmWave or in the THz band* and that each IAB node features two Radio Frequency (RF) chains, which are used for the backhaul and the fronthaul communications, respectively. In both cases, gNBs are equipped with directional antennas.

When gNB $k = 0, \dots, N_I$, with index 0 denoting the IAB donor, receives data from node j , packets experience a SINR $\gamma_{s,d}$ which can be expressed as

$$\gamma_{s,d} = \frac{|h_{s,d}^l|^2 \sigma_x^2}{\sigma_n^2 + \sum_{i \in \mathcal{I}} \sigma_i^2}, \quad (3.24)$$

where $h_{s,d}^l$, $l \in \{mW, sT\}$ represents the equivalent channel response between the communication endpoints when using mmWave or sub-THz links, respectively. \mathcal{I} denotes the set of interferers, σ_x^2 , σ_i^2 and σ_n^2 are the powers of the transmitted signal, the i -th received interfering signal, and the thermal noise at the receiver, respectively.

The corresponding access (backhaul) throughput $R_{j,k}^A(t)$ ($R_{j,k}^B(t)$) reads

$$R_{j,k}^A(t) = \frac{1}{T_{sub}} \sum_{l=1}^{B_j^t} \mathbb{1}\left\{\hat{b}_l(\gamma_{j,k}) = b_l\right\}, \quad (3.25)$$

where B_j^t denotes the number of bits transmitted from user (IAB node) j to gNB k during subframe t and $\hat{b}_l(\gamma_{j,k})$ is the l -th decoded bit at the receiver, as a function of $\gamma_{j,k}$.

Our goal is to maximize the average system sum-rate, defined as

$$\bar{R} \doteq \frac{1}{T} \sum_{j=1}^{N_I} \sum_{t=1}^T R_{j,0}^B(t), \quad (3.26)$$

by tuning the carrier frequency (either mmWave or THz) of each backhaul link. We remark that in this metric we take into account only the packets which are received at their final destination, i.e., the IAB donor.

3.3.1.1 Channel Models

mmWave channel model For the mmWave links, we consider the 3GPP 38.901 SCM [23], which models MIMO wireless channels for frequencies between 0.5 and 100 GHz. In particular, [23] outlines the procedures for generating a channel matrix $\mathbf{H}_{s,d}$ whose entries $h_{s,d}^{j,k}$ correspond to the impulse response of the channel between the j -th element of the antenna array of the transmitter (S), and the k -th radiating element of the antenna array of the receiver (D). Then, the channel matrix entries are combined with a frequency-flat path loss term PL .

3.3 High-capacity integrated access and backhaul networks using sub-terahertz links

$$\underset{\mathbf{P}, \{\mathbf{S}(t)\}_t, \mathbf{T}}{\operatorname{argmax}} \bar{R}, \quad (3.28a)$$

$$\text{s.t. C1: } R_{j,k}^B(t)T_{sub} \leq B_j(t) \quad \forall j, \quad \forall t \quad (3.28b)$$

$$\text{C2: } B_j(t+1) = B_j(t) + T_{sub} \left(\sum_{k=1}^{N_U} R_{k,j}^A(t) + \sum_{k=1}^{N_I} R_{k,j}^B(t) - \sum_{k=0}^{N_I} R_{j,k}^B(t) \right) \quad \forall j, \quad \forall t \quad (3.28c)$$

$$\text{C3: } \sum_{k=0}^{N_I} \mathbf{S}[j, k](t) + \sum_{k=1}^{N_I} \mathbf{S}[k, j](t) \leq 1 \quad \forall j, \quad \forall t \quad (3.28d)$$

$$\text{C4: } R_{j,k}^B(t) \mathbf{S}[j, k](t) = R_{j,k}^B(t) \quad \forall j, \quad \forall k, \quad \forall t \quad (3.28e)$$

$$\text{C5: } \sum_{j,k=0}^{N_I} \mathbf{T}[j, k] \leq \rho_{max} \sum_{j,k=0}^{N_I} \mathbf{P}[j, k] \quad (3.28f)$$

When considering analog beamforming at both the transmitter and the receiver, the equivalent channel response $h_{s,d}^{mW}$ can be evaluated as

$$h_{s,d}^{mW} = \sqrt{10^{PL/10}} \cdot \mathbf{w}_d \mathbf{H}_{s,d} \mathbf{w}_s, \quad (3.27)$$

with \mathbf{w}_s and \mathbf{w}_d the beamforming vectors used at S and D, respectively.

THz channel model For sub-THz, we use the physics-based channel modeling approach from [157], which includes molecular absorption and path loss. At THz-band frequencies, molecular absorption, which causes both molecular absorption loss and molecular absorption noise, is the principal factor affecting electromagnetic wave propagation. $h_{s,d}^{tH}$ is the THz-band channel model introduced in [157], with additional transmit and receive antenna gains G_S and G_D , and is given by

$$h_{s,d}^{tH}(f, d) = \frac{c}{4\pi f d} \exp\left(-\frac{k_{abs}(f)d}{2}\right) G_S G_D, \quad (3.29)$$

where c stands for the speed of light and k_{abs} for the medium's molecular absorption coefficient, based on the type and composition of molecules [158].

3.3.2 Sum-rate optimization via THz Link Selection

We define $\mathbf{P} \in \{0, 1\}^{N_l+1 \times N_l+1}$ as the matrix which represents the possible active links among gNBs, i.e., $\mathbf{P}[i, j] = 1$ if and only if the wireless backhaul link between gNBs i and j is a feasible link; index 0 refers to the donor. Similarly, $\mathbf{S}(t) \in \{0, 1\}^{N_l+1 \times N_l+1}$ and $\mathbf{T} \in \{0, 1\}^{N_l+1 \times N_l+1}$ represent the links which are active during subframe t , and whether they use THz spectrum or not, respectively. Our objective is to maximize the average system sum-rate, by choosing whether each link is operating in the THz or the mmWave band and the active links in each subframe. We perform the choice of \mathbf{T} and \mathbf{P} only once, with the goal of reducing the computational complexity of the algorithm.

The optimization problem is thus formulated as (3.28a). Constraint C₁ ensures that nodes do not transmit more data than available in their buffer. C₂ enforces the proper evolution over time of the buffers occupancy, i.e., the buffer occupancy at time t must be equal to the one in subframe $t - 1$, minus (plus) the outgoing (incoming) traffic from other nodes. Constraint C₃ relates to the TDMA mode of operation, and ensures that each backhaul RF chain is used at most for one transmission/reception in any given subframe, while C₄ imposes that only active links can exhibit a positive rate. Finally, with C₅ we set an upper bound ρ_{max} on the maximum percentage of THz links.

3.3.2.1 Backhaul Scheduler

We remark that due to the binary nature of the \mathbf{P} , $\mathbf{S}(t)$ and \mathbf{T} optimization variables, (3.28a) is an Integer Linear Program (ILP), thus NP-hard and not solvable in polynomial time. Therefore, in this section, we present a set of algorithms that solve the path selection and configuration problem heuristically and with low complexity.

Specifically, we first describe the pre-processing steps, referred to as *distance-aware path generation* (Alg. 3) and *THz-link selection* (Alg. 4), which prune the set of possible links established among gNBs and decide which of them are to operate in the THz bands, respectively. Then, we describe the *SINR-based scheduler* (Alg. 5), which differs from the former procedures as it is executed in each subframe to track the dynamic nature of the backhaul network.

The distance-aware path generation algorithm computes the \mathbf{P} matrix, which encodes the potential connections between IAB nodes. \mathbf{P} reduces

Algorithm 3 Distance Aware Path Generation

```

 $d_{max} \leftarrow$  Max distance between IAB nodes of the same tier
 $P = [0]_{N_I+1 \times N_I+1}$ 
for  $n_i = 1, 2, \dots, N_I$  do
     $d_i \leftarrow$  3D distance between  $n_i$  and IAB donor
    if  $d_i < d_{max}$  then
         $P[n_i, 0] = 1$ 
    end if
    for  $n_j = n_1 + 1, \dots, N_I$  do
         $d_{i,j} \leftarrow$  3D distance between  $n_i$  and  $n_j$ 
        if  $d_{i,j} < d_{max}$  then
             $d_j \leftarrow$  3D distance between  $n_j$  and IAB donor
            if  $d_i < d_j$  then
                 $P[n_j, n_i] = 1$ 
            else
                 $P[n_i, n_j] = 1$ 
            end if
        end if
    end for
end for

```

the system complexity by restricting possible paths from each IAB node and by avoiding loops. Specifically, Alg. 3 iterates over each IAB node n_j , establishing a connection towards the donor whenever the distance between them is smaller than d_{max} , i.e., a scenario- and frequency-dependent distance that guarantees a link performance above a certain threshold. In our case, the considered scenario involves a small and dense deployment of IAB nodes, so the path loss distance can be compensated by the antenna gain, and d_{max} for THz and mmWave are assumed to have the same value. Moreover, the proposed pre-processing step performs additional attachments between neighboring nodes, as long as the resulting link exhibits a lower length than d_{max} . The direction of such link is determined in such a way that the destination node is the closer to the donor. Even though this link may be topologically redundant, it can provide an alternative route for load balancing purposes, while still avoiding the creation of cycles.

The THz link selection policy identifies bottleneck links based on two heuristics: 1) links involving IAB nodes which are closer to the donor are more likely to be congested since they are usually used also for relaying traffic of subtending nodes; and 2) the average buffer occupancy provides an estimate of the loads incurred on each link. Accordingly, Alg. 4 partitions the

IAB nodes into disjoint sets, referred to as *tiers*. Nodes are assigned to tiers based on their distance with respect to the donor, with tier 0 indicating the closest level to the donor. Then, the various backhaul links are marked as THz in descending order with respect to the tier of the corresponding transmitting node, until the maximum ratio of non-mmWave links ρ_{max} is reached. Note that the algorithm may eventually reach a tier whose IAB nodes are not all set as THz. In this case, ties within the same tier are broken by sorting its nodes with respect to their average traffic load, which we estimate by measuring the respective buffers. That is to say, nodes with higher buffer occupancy are given priority and thus are set as THz before nodes exhibiting a lower traffic load. Note that this procedure can be based on long-term statistics, thus averaging the load of the nodes over multiple frames.

Finally, the SINR-based scheduler dynamically allocates resources, with the objective of maximizing the average sum rate by choosing a list of paths to be activated in each subframe. The rationale behind the proposed scheme is to schedule links based on their load. Specifically, in Alg. 5 we assign a transmission resource allocation priority which is directly proportional to the buffer occupancy of the transmitting node. Once the first endpoint is chosen, we determine the outgoing link by selecting the one with the highest SINR among those calculated in Alg. 3. Then, we set all links involving the corresponding transmitting and receiving nodes as infeasible (assigning zero to the corresponding transmitting (n) and receiving node (p_n^*) indices in P_{temp}), and repeat the procedure by considering the remaining nodes and links only, thus ensuring that the TDMA constraint is satisfied.

3.3.3 Performance Evaluation

This section introduces a performance evaluation based on a novel simulation setup (Sec. 3.3.3.1) in a dense cellular network (Sec. 3.3.3.2), with a comparison between different results of THz and mmWave networks (Sec. 3.3.3.3).

3.3.3.1 Simulation Setup

We have developed a system-level simulator that runs on top of Sionna [134], an open-source TensorFlow-based GPU-accelerated toolbox, and that includes the IAB networks described in Rel. 17. The proposed simulator, which is written in Python, is a system-level simulator which features 3GPP-compliant

Algorithm 4 THz Link Selection

```

 $\mathbf{N}_T$  = Vector of IAB nodes tier index
 $\mathbf{N}_{sort}$  = Vector of IAB node indices, sorted with respect to their load
 $\mathbf{T} \leftarrow [0]_{N_I+1 \times N_I+1}$ 
 $d_{max} \leftarrow$  Max distance between IAB nodes of the same tier
for  $n = 1, 2, \dots, N_I$  do
     $d \leftarrow$  3D distance between  $n$  and IAB donor
     $\mathbf{N}_T[n] \leftarrow \lfloor d / d_{max} \rfloor$ 
end for
for  $i = 1, \dots, \max(\mathbf{N}_T)$  do
     $\mathbf{N}_T^i \leftarrow \{j \mid \mathbf{N}_T[j] == i\}$ 
     $\mathbf{L}_i \leftarrow$  links in  $\mathbf{P}$  where nodes of  $\mathbf{N}_T^i$  are the transmitting node
    if  $\sum_{j,k} \mathbf{T}[j, k] + \dim(\mathbf{L}_i) < \rho_{max} \sum_{j,k} \mathbf{P}[j, k]$  then
         $\mathbf{T}[j, k] \leftarrow 1 \forall (j, k) \in \mathbf{L}_i$ 
    else
        while  $\sum_{j,k} \mathbf{T}[j, k] < \rho_{max} \sum_{j,k} \mathbf{P}[j, k]$  do
             $n^* \leftarrow \min_n \mid \mathbf{N}_{sort}[n] \cap \mathbf{N}_T^i \neq \emptyset$ 
             $(n^*, k) \leftarrow$  link  $\in \mathbf{L}_i \mid n^*$  is the transmitting node
             $\mathbf{T}[n^*, k] \leftarrow 1; \mathbf{L}_i \leftarrow \mathbf{L}_i \setminus (n^*, k)$ 
        end while
    end if
end for

```

channel modeling and lower layers of the protocol stack. However, it lacks the implementation of 5G NR higher layers. Therefore, we added system-level functions like MAC-level scheduling and RLC-level buffering [159]. In addition, in this research we use the Terasim channel simulator [158] to generate channel responses and integrate them into Sionna. To accomplish this, we generate traces for each IAB node's channel response and load them into Sionna. Terasim channel model integration allows us to generate channels up to 10 THz; in this simulation campaign, the sub-THz carrier frequency is 140 GHz. Several system-level KPIs, including latency, throughput, and packet loss rate, are produced by our simulator.

3.3.3.2 Simulation Scenario

We take a dense cellular base station deployment into account in our models. As shown in Fig. 3.24, we place IAB nodes at a density of 150 gNB/km², thus with an average intersite distance of 40 m. In Table 3.2, the specific simulation settings are displayed. For mmWave, we used the channel model outlined by 3GPP in TR 38.901 [23], a statistical 3GPP channel model for 0.5-100 GHz,

Algorithm 5 SINR-based Scheduler

```

 $N_{sort}$  = Vector of IAB nodes, sorted with respect to their load
 $P_{temp} = P$ 
 $S(t) \leftarrow [0]_{N_I+1 \times N_I+1}$ 
for  $n$  in  $N_{sort}$  do
     $\gamma_{max} \leftarrow -\infty$ 
    for  $i$  in  $0, \dots, N_I$  do
        if  $\gamma_{n,i} > \gamma_{max}$  then
             $\gamma_{max} \leftarrow \gamma_{n,i}$ 
             $p_n^* \leftarrow i$ 
        end if
    end for
     $S(t)[n, p_n^*] \leftarrow 1$ 
     $P_{temp}[:, n], [n, :] \leftarrow [0]; P_{temp}[:, p_n^*], [p_n^*, :] \leftarrow [0]$ 
end for

```

while for sub-THz we used the THz-band channel model introduced in [157] and detailed in Sec. 3.3.1.1. The range of the user rate is 20 Mbps to 500 Mbps. We used a phased array antenna for mmWave and a horn antenna for THz, respectively. In mmWave we do beamforming based on a pre-generated codebook, in order to find the best beam pair for connection. For the purposes of SINR calculation, we assume that each interfering device utilizes the beamforming vector with the greatest SINR towards its intended target. In a similar fashion, both the transmitter and the receiver utilize the beamforming configuration calculated by the hierarchical search technique. We consider a scenario with a single donor to focus on the issues related to the bottleneck in the air interface of the donor itself, while extension to multiple donors is left for future work. We also set $d_{max} = 70$ m, as it has been experimentally shown that sub-THz links can operate in this range also in adverse weather conditions [160].

3.3.3.3 Results

In this section we report the outcomes of our numerical evaluation, focusing on end-to-end metrics measured at the IAB donor. We compare the performance achieved by different backhaul configurations, i.e., different maximum ratios of THz links and bandwidth, in terms of throughput, latency and packet drop ratio. We consider two baselines: *Random Scheduler* (RS) and *Random Links* (RL). The former uses Alg. 4 and chooses at random a feasible set of active links during each subframe. On the contrary, RL randomly picks

3.3 High-capacity integrated access and backhaul networks using sub-terahertz links

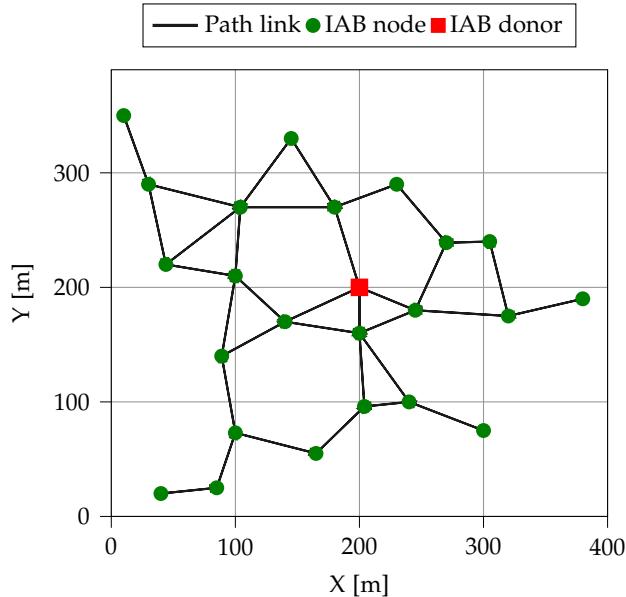


Figure 3.24: Simulation Scenario

Table 3.2: Simulation parameters.

Parameter	Value
Carrier frequency (mmWave)	28 GHz
Bandwidth (mmWave)	400 MHz
Carrier frequency (THz)	140 GHz
Bandwidth (THz)	{10, 32} GHz
IAB RF Chains	2 (1 access + 1 backhaul)
Pathloss model (mmWave)	UMi-Street Canyon [23]
Pathloss model (THz)	Physics-based [157]
Number of IAB nodes N_I	23
Number of users N_U	50
Per-UE source rate	{40, 80, 100, 200} Mbps
ρ_{max}	{0, 0.1, 0.3, 0.5, 0.7, 1}
gNB antenna array	8H × 8V
UE antenna array	4H × 4V
gNB and UE height	15 m and 1.5 m
gNB antenna gain (mmWave)	30 dB
gNB antenna gain (THz)	38 dB
Noise power	10 dBm

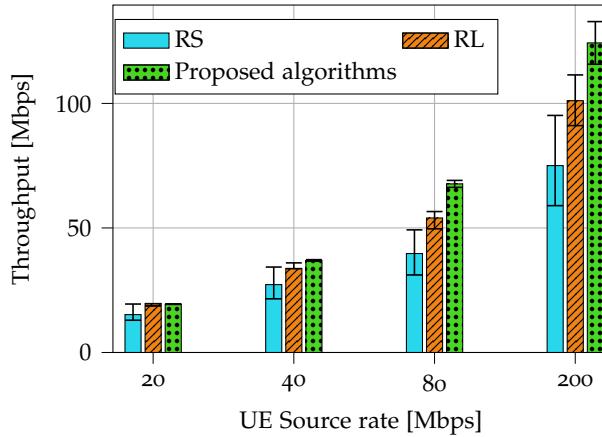


Figure 3.25: Throughput per UE for different schedulers and THz link selection policies, for THz bandwidth 32 GHz and $\rho_{max} = 0.3$.

which links to set as THz, and uses Alg. 5 for scheduling. 10 simulations per configuration are executed, to obtain estimates which are averaged over the realizations of the wireless channels.

Fig. 3.25 reports the UE throughput achieved by the proposed solution, versus that achieved by RS and RL. Focusing on the former, it can be seen that Alg. 5 leads to a throughput increase of up to 40% compared to a random scheduling policy, thanks to the prioritization of the backhaul links incurring a higher load and exhibiting a higher number of subtending IAB nodes. Moreover, Alg. 4 introduces an additional throughput increase of up to 15% compared to RL.

Fig. 3.26 illustrates the UE throughput for various configurations of sub-THz backhauling links and different UE source rates. The performance always improves by adding more bandwidth to the system through sub-THz links, despite the harsher propagation environment at higher frequencies.

The performance gap increases with the user source rate. Indeed, mmWaves successfully sustain a source system rate of 1 Gbps (20 Mbps for 50 UE), but cannot match higher source rates, as the capacity saturates. The configuration with sub-THz links achieves a higher throughput in all scenarios and in particular for $\rho_{max} = 0.3$, 32 GHz achieves the highest throughput for all source rates. It is obvious that increasing the bandwidth improves the performance; nevertheless, increasing the percentage of the THz links from $\rho_{max} = 0.1$ to $\rho_{max} = 0.3$ has a more significant impact on throughput. This

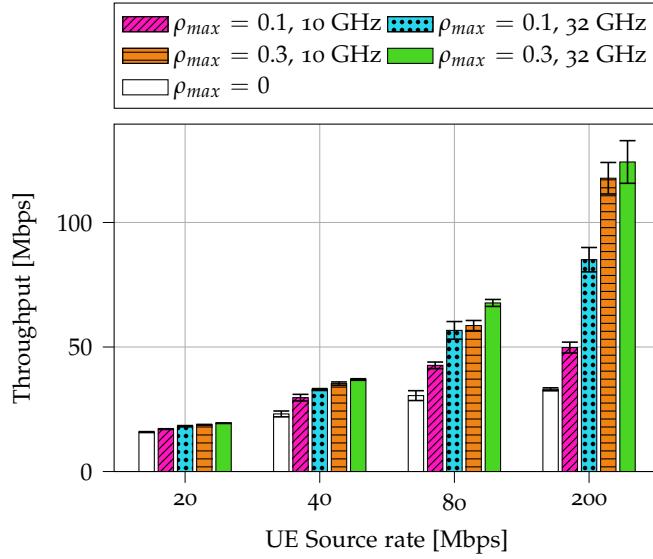


Figure 3.26: Throughput per UE for different configurations.

may be explained by considering the effects of replacing bottleneck backhaul mmWave links with THz links with higher bandwidth.

Similar considerations can be drawn from the results shown in Fig. 3.27, which reports the packet drop percentages for various backhaul configurations. The highest and lowest packet drop percentages across all UE source rates are achieved when using the mmWave and $\rho_{max} = 0.3, 32 \text{ GHz}$ configurations, respectively. Packet drop percentages at 20 Mbps source rates are close to zero for all configurations. The highest packet drop percentages among configurations including THz is $\rho_{max} = 0.1, 10 \text{ GHz}$. It is noteworthy that the system performance is influenced directly by both the THz bandwidth and the link ratio, as seen in Fig 3.26.

Fig. 3.28 depicts the ECDF of the End-to-End (E2E) latency experienced by packets which reach the donor, for different backhaul configurations. Accordingly, both latencies accumulated over the fronthaul and backhaul links are taken into account, from the time packets are generated at the UE until they eventually reach the IAB donor. The plot shows that packet latency decreases as more sub-THz links are added to the network. In accordance with the aforementioned observations (Fig. 3.26 and Fig. 3.27), $\rho_{max} = 0.3, 32 \text{ GHz}$ has the lowest latency, whereas mmWave has the highest latency. The average latency for $\rho_{max} = 0.3, 32 \text{ GHz}$, $\rho_{max} = 0.3, 10 \text{ GHz}$, $\rho_{max} = 0.1, 32 \text{ GHz}$

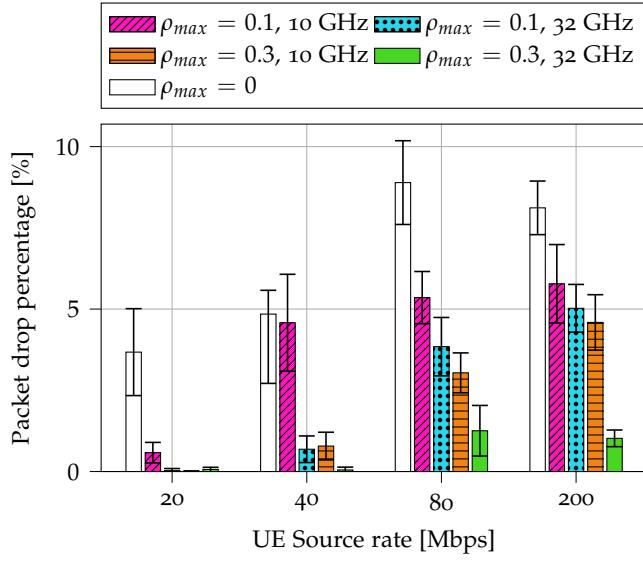


Figure 3.27: Backhaul packet drop percentage for different configurations.

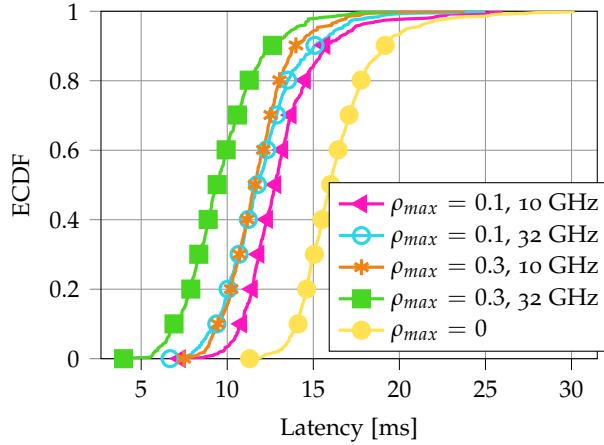


Figure 3.28: E2E latency ECDF for different configurations for 80 Mbps user rate.

GHz, and $\rho_{max} = 0.1, 10 \text{ GHz}$ is approximately 51%, 24%, 24%, and 18% less than in mmWave.

Finally, the average system throughput for different ratios ρ_{max} of THz link is shown in Fig. 3.29. The system throughput increases with the inclusion of additional THz links. The figure also shows that system source rates of 2 Gbps, 4 Gbps, and 10 Gbps can be satisfied by a single donor when ρ_{max} is properly set. However, the larger demand of the 25 Gbps system source rate still cannot be satisfied, as the system becomes saturated. $\rho_{max} = 0.1$ and 1

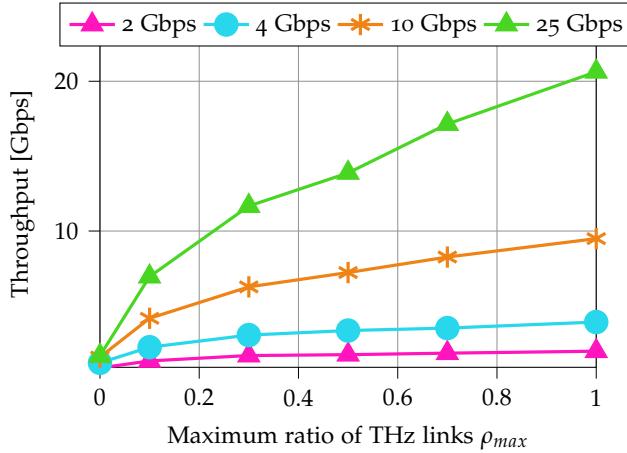


Figure 3.29: System throughput for different source rate and ratio of THz links.

can increase the system throughput by up to four times and twelve times, respectively.

3.4 Conclusions and future work

In this chapter we proposed a semi-centralized resource partitioning scheme for 5G and beyond IAB networks, coupled with a set of allocation policies. We showed that the introduction of this light resource allocation cooperation dramatically improves the end-to-end throughput and delay achieved by the system already, preventing (or at the very least limiting) the insurgence of network congestion in the backhaul links. We provided considerations on the implementation of a semi-centralized resource allocation controller in real world deployments. In particular, we acknowledged that the proposed scheme relies on the assumption of IAB-nodes being capable of exchanging timely feedback information with the IAB-donor. Even though the amount of signaling data which the proposed solution requires is quite low, and its performance is quite robust with respect to an increase of the central allocation period, we argue that this remains a significant constraint. As a consequence, we deem that solutions involving a central controller, which rely on timely exchange of control information with the IAB-donor, are likely to require dedicated control channels, possibly at sub-6 GHz, in order to grant the utmost priority and reliability to the feedback information. Therefore, we conclude that semi-centralized frameworks can bring dramatic per-

formance benefits to IAB networks, although their introduction in 5G and beyond deployments requires additional research efforts. Moreover, we proposed the first reliability-focused scheduling and path selection algorithm for IAB mmWave networks. We illustrated that our RL-based solution can cope with the network dynamics including channel, interference, and load. Furthermore, we demonstrated that Safehaul not only exhibits highly reliable performance in the presence of the above-mentioned network dynamics, but also outperforms the benchmark schemes in terms of throughput, latency and packet-drop rate. The reliability of Safehaul stems from the joint minimization of the average latency, and the expected value of its tail losses, by leveraging CVaR as a risk metric. Finally, we provided the first performance evaluation of the possibilities of sub-terahertz frequencies for 6G IAB using a customized extension of the open-source Sionna simulator. This permits the use of greedy algorithms to evaluate the deployment of mixed mmWave and sub-terahertz links to boost the backhaul network's capacity.

As part of our future work on IAB optimization frameworks, we plan to design efficient machine-learning algorithms which predict the network evolution at the IAB-donor. This improvement will allow us to relax the timely feedback assumption, by increasing the minimum semi-centralized allocation period which leads to performance benefits over distributed strategies. Moreover, we foresee to implement mechanisms which adapt the parameters of the MRBA policy to the system load and configuration, additional resource partitioning strategies, and the generalization of the proposed framework to SDMA systems. Additionally, we identify network reliability as a highly under-explored topic that deserves further investigation. Some interesting research directions in this space are the maximization of reliability under the assumption of statistical system knowledge, or the evaluation of the network's reliability when the functionality of the BAP layer is compromised. Furthermore, our system-level extension to Sionna can be further developed to support an arbitrary number of RF chains and in-band backhauling, allowing more extensive investigation of IAB protocols and architectures. Finally, we will extend the performance analysis of sub-terahertz IAB deployments to cover a broader range of source traffic patterns, scenarios (including multi-donor instances, deployments with lower node density, or more realistic map-based scenarios as in [161, 162]), and protocol stack implementations as future work.

4 *Downlink Clustering-Based Scheduling of IRS-Assisted Communications With Reconfiguration Constraints*

Despite the potential of IAB for improving network coverage, and its lower cost compared to wired-backhaul deployments, the former involves complex signal processing and saturation of the available resources, and may still be too costly and energy-consuming for network operators. In light of this, IRSs are being investigated as solutions to overcome the harsh propagation conditions shown by mmWave and THz bands in a cost- and energy-efficient manner [31]. IRSs are meta-surfaces, whose radiating elements can *passively* tune the phase shift of impinging signals to favorably alter an electromagnetic field towards an intended destination. They can be configured to beamform the reflected signal virtually in any direction, hence acting as a relay to improve the signal quality without an active (power-consuming) amplification [32].

4.1 Prior Work

Despite the substantial research hype, most recent studies on IRSs rely on strong assumptions that do not match real-world deployments. Specifically, a significant body of literature is based on the assumption that IRSs establish an ideal (i.e., fiber-like) control channel with the base station [48, 163–165]. Instead, actual deployments are expected to feature a wireless, i.e., error-prone, IRS control, possibly implemented with low-cost technologies [166, 167]. This introduces constraints on the IRS reconfiguration period, which needs to be synchronized with the base station to beamform the signal towards the UE served during the specific time slot [31], a similar research problem to scheduling in cellular networks.

In this perspective, IRS-assisted downlink scheduling solutions have been widely studied in different domains, each with its own theoretical constraints. For example, in Orthogonal Frequency-Division Multiple Access (OFDMA) user scheduling, all the users scheduled in a given time slot must be served using the same reflection coefficients, due to the lack of frequency selective beamforming capabilities at the IRS. In this context, dynamic optimization schemes, wherein the IRS configurations are adjusted at each time slot, have been studied in [168, 169]. The authors of [170] consider a two-user downlink transmission problem in an IRS-assisted scenario over fading channels, and compare the results of different basic Orthogonal Multiple Access (OMA) and Non-Orthogonal Multiple Access (NOMA) schemes. It is found that, while NOMA is the best solution, by exploiting IRS reconfiguration in each slot of the fading block, TDMA outperforms Frequency Division Multiple Access (FDMA), and its performance is similar to that of NOMA. A hybrid TDMA-NOMA approach, instead, was investigated in an uplink scenario in [171, 172], in the context of a wireless-powered network, where users are grouped based on their channel gains. Then, UEs within the same group transmit in a non-orthogonal fashion, while different groups are assigned to different time slots. Moreover, a user scheduling algorithm based on graph neural networks, able to jointly optimize the IRS configuration and the gNB beamforming in downlink, was recently presented in [173]. Similarly, the authors of [174–176] evaluated the performance of several non-orthogonal downlink scheduling methods, such as Rate-Splitting Multiple Access (RSMA). Finally, IRSs with energy harvesting capabilities are considered in [177]. In this work, the authors propose a trade-off between the system sum capacity and the IRS energetic self-sustainability, with the goal of achieving coverage flexibility and low deployment costs.

Still, most of the literature poses little to no reconfiguration constraints for the IRS. However, early IRS control circuitry prototypes, which have low power consumption (i.e., a few hundreds of mW), have a non-negligible phase-shifts reconfiguration time [178, 179], thus posing additional constraints in the system design. For example, the prototypes in [180] and [181] have a reconfiguration time of a few tens of ms, even though architectures based on Field Programmable Gate Array (FPGA) such as in [182] promise to achieve much lower configuration times, i.e., in the order of tens of microseconds. Still, the overhead (in terms of time) increases as the number of IRS ele-

ments increases, as investigated in [183–185]. In any case, a constraint on the number of reconfigurations (and relative period) is desirable to ensure system synchronization and minimize the IRS downtime during reconfiguration. In this regard, it is of interest to (*i*) investigate the level of performance degradation experienced by IRS-assisted systems when considering practical constraints, including limitations in the number of reconfigurations, and (*ii*) design algorithms that can mitigate these constraints. The limitation on the number of IRS reconfigurations in a given time frame has been initially studied in [178], where the authors evaluate the capacity of both OMA and NOMA schemes of a 2-user IRS-assisted SISO system under Rayleigh fading conditions. Still, additional research efforts are required to fully characterize the impact of IRS reconfigurations constraints on the network.

4.2 Contributions

In this chapter, we propose a TDMA scheduling policy for downlink cellular transmissions based on clustering algorithms, to maximize the sum capacity in IRS-assisted network deployments with practical constraints. Our main contributions are summarized as follows:

- We account for practical IRS limitations by considering a fixed maximum number of reconfigurations of IRS reflecting elements within a time frame, thus setting a simple constraint on the overhead entailed by the control of the IRS.
- We formalize an optimization problem to determine the optimal IRS configurations to maximize the sum capacity while satisfying the reconfiguration per frame constraint. Then, we convert the sum capacity problem into a clustering problem. The latter determines sets of UEs that can be served with the same (possibly suboptimal) IRS configuration while minimizing the related capacity loss.
- We design, as an alternative to typical clustering algorithms based on distance measures, a new class of algorithms which we denote as *capacity-based clustering*. These algorithms adjust the cluster configuration taking into account the sum capacity and the user fairness. Specifically, we propose three clustering algorithms: Capacity-Weighted Cluster-

ing (CWC), which favors users experiencing the best channel conditions, One-Shot Capacity-Based Clustering (OSCBC), which represents a low-complexity alternative to the former, and Inverse Capacity-Weighted Clustering (ICWC), which promotes fairness among the cluster UEs.

- We compare via simulation the performance of distance- and capacity-based clustering in different IRS-assisted scenarios. Extensive numerical results show that scheduling based on clustering can reduce by up to 50% the number of IRS reconfigurations, thus promoting communication efficiency at the expense of a slightly lower sum capacity.

With respect to [186], we introduce new capacity-based clustering strategies to improve fairness and provide more extensive numerical results to demonstrate the scalability of the proposed solutions as a function of the density of UEs and the IRS size. Moreover, we evaluate the performance of the proposed scheduling strategies considering realistic IRS network constraints, including the quantization of phase shifts, and for different channel propagation conditions. In this sense, we provide additional results in terms of the computational complexity of the proposed distance- and capacity-based clustering algorithms, as well as in terms of fairness. Finally, we remark that due to the novelty of the considered scenario, to the best of our knowledge the effectiveness of our solution cannot be directly compared with any works in the literature. Indeed, the most similar works [169, 170, 178, 183] exhibit substantial differences in the considered contexts. More specifically:

- In comparison to [169], which assumes the IRS configurations to be fixed, our work focuses on optimizing user scheduling in conjunction with IRS configurations in a dynamic manner.
- [170, 178] analyze a basic scenario where a single-antenna gNB serves single-antenna UEs, while we consider multiple antennas at both the base station and UEs, and an OFDMA multiuser access scheme.
- While [178] characterizes the capacity region for $K > 1$ UEs, the computational complexity of the proposed scheme limits its applicability to $K = 2$ UEs. In contrast, our solution supports more realistic scenarios, where $K \gg 1$.
- The authors of [183] use the position estimate of a moving UE to minimize the IRS reconfiguration overhead while guaranteeing a minimum

SNR in the single-UE context. In our work, on the other hand, we consider a multi-user scenario and reduce the number of IRS configurations by clustering UEs in the CSI domain.

These fundamental differences in system setup and assumptions prevent a meaningful simulation-based comparison.

4.2.1 Organization and Notation

The rest of the chapter is organized as follows. In Section 4.3, we introduce the system model. In Section 4.4, we present the sum capacity optimization problem. In Section 4.5, we describe the scheduling framework, while in Sections 4.6 and 4.7 we present distance-based and capacity-based clustering algorithms, respectively. In Section 4.8, we show numerical results and compare the different scheduling and clustering solutions. Finally, Section 4.9 draws the main conclusions.

Scalars are denoted by italic letters; vectors and matrices by boldface lowercase and uppercase letters, respectively; sets are denoted by calligraphic uppercase letters. $\text{diag}(\mathbf{a})$ indicates a square diagonal matrix with the elements of \mathbf{a} on the principal diagonal, and $\text{vec}(\mathbf{A})$ denotes the vectorization operator, stacking the columns of matrix \mathbf{A} into a column vector. \mathbf{A}^T and \mathbf{A}^\dagger denote the transpose and the conjugate transpose of matrix \mathbf{A} , respectively. $[\mathbf{A}]_{k\ell}$ denotes the scalar value in the k -th row and ℓ -th column of matrix \mathbf{A} , while $[\mathbf{a}]_k$ denotes the k -th element of vector \mathbf{a} . The imaginary unit is denoted as $j = \sqrt{-1}$, and $\angle a$ denotes the phase of $a \in \mathbb{C}$. The operator \diamond denotes the Khatri-Rao product. Finally, $\mathbb{E}[\cdot]$ denotes statistical expectation.

4.3 System Model

We consider downlink data transmissions for the multi-user MIMO communication system shown in Fig. 4.1, wherein the transmission from the gNB to the K UEs is assisted by an IRS. The gNB and the UEs are equipped with N_g and N_U antennas, respectively. We assume that the direct link between the gNB and the UEs is unavailable due to blockage. As a consequence, the gNB transmits signals to the UEs by exploiting the virtual link offered by the IRS. In this context, the IRS configuration is managed by the gNB through the IRS controller, by exploiting a dedicated link between the gNB

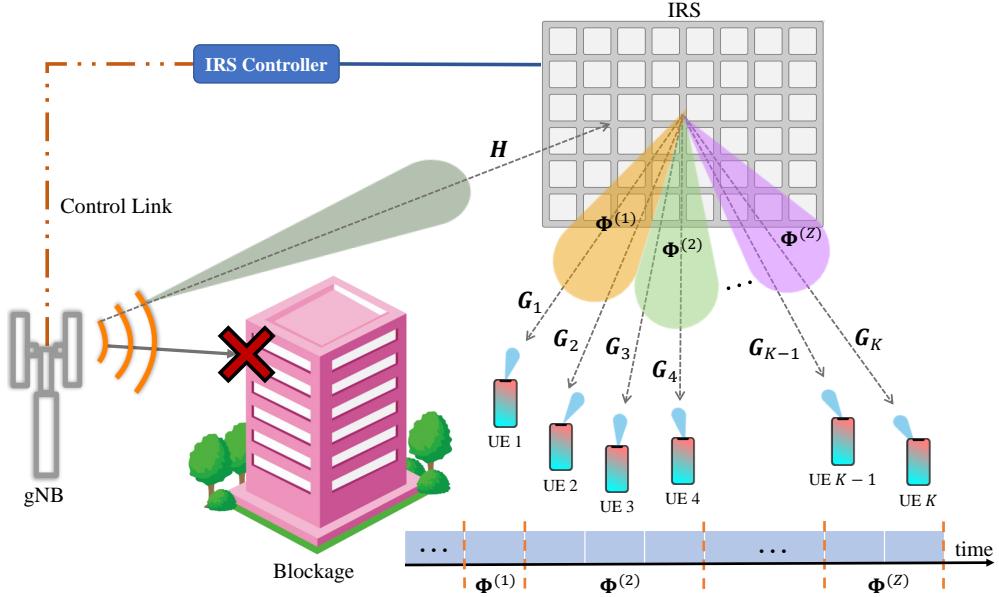


Figure 4.1: Downlink TDMA scheduling for multi-user IRS-aided systems.

and the IRS, thus with no additional communication overhead in the gNB-UE link. Time is divided into frames of K slots, and each UE is served exactly once in a frame in a TDMA fashion, which ensures there is no co-channel interference as UEs are separated in the time domain. In our scenario we expect that the gNB-IRS channel has rank one, i.e., a single dominant path, which effectively prevents multi-stream transmissions and spatial multiplexing. However, when higher-rank channels are available, either multi-stream transmission to each UE or spatial multiplexing can be considered. For the former case, the proposed solution applies straightforwardly. Moreover, our approach can be suitably modified to accommodate for the latter scenario. A detailed investigation of this point is left for future work. We assume that UEs are either static or moving slowly, which is the most typical application scenario for IRS-assisted networks. Under such conditions, the channel coherence time is in the order of 10 ms [187, Fig. 5]. Considering that perfect CSI of all UEs is acquired at the gNB at the beginning of each frame (a realistic assumption that does not affect the proposed scheduling framework for IRS communication), it is reasonable to conclude that the channel remains

constant throughout the whole time frame. Here, we assume that the CSI is available for any IRS configuration.

4.3.1 IRS Model

Each of the N_I elements of the IRS acts independently as an omnidirectional antenna unit that reflects the impinging electromagnetic field by introducing a tunable phase shift on the baseband-equivalent signal. We denote as $\phi_n = e^{j\theta_n}$ the reflection coefficient of the n -th IRS element, where $\theta_n \in \mathcal{P}_\theta$ is the induced phase shift, and \mathcal{P}_θ is the set of possible phase shifts. Recent works argue that continuous phase shifts are hardly implementable in practice [188]. Therefore, we consider both continuous and quantized phase shifts. While in the former case the set of phase shifts is $\mathcal{P}_\theta = [-\pi, \pi)$, in the latter we have $\mathcal{P}_\theta = \left\{0, \frac{2\pi}{2^b}, \dots, \frac{2\pi(2^b-1)}{2^b}\right\}$ where $b > 0$ is the number of bits employed to quantize the phase shifts.

We denote with $\mathbf{H} \in \mathbb{C}^{N_I \times N_g}$ the channel matrix between the IRS and the gNB, and with $\mathbf{G}_k \in \mathbb{C}^{N_U \times N_I}$ the channel matrix of the link between the IRS and UE k , respectively. We consider single-stream transmissions,¹ with $\mathbf{w}_k \in \mathbb{C}^{N_g \times 1}$ and $\mathbf{v}_k \in \mathbb{C}^{N_U \times 1}$ defined as the beamforming vectors at the gNB and UE k , respectively. Let x_k be the single-stream signal transmitted by the gNB to UE k ; the received post-processing signal can be expressed as

$$z_k = \mathbf{v}_k^\top \mathbf{G}_k \boldsymbol{\Phi} \mathbf{H} \mathbf{w}_k x_k + \mathbf{v}_k^\top \mathbf{n}_k, \quad (4.1)$$

where $\mathbf{n}_k \in \mathbb{C}^{N_U \times 1}$ represents the circularly symmetric complex Gaussian noise vector with entries having zero mean and variance σ_n^2 , while $\boldsymbol{\Phi} \in \mathbb{C}^{N_I \times N_I}$ is the IRS configuration, i.e., a diagonal matrix defined as $\boldsymbol{\Phi} = \text{diag}(\phi_1, \dots, \phi_{N_I})$. Note that different, and specific, IRS configurations can be adopted for different UEs. Accordingly, in the rest of the chapter we let $\boldsymbol{\Phi}_k$ be the IRS configuration adopted when UE k is served.

The SNR at UE k under IRS configuration $\boldsymbol{\Phi}_k$ is

$$\Gamma_k(\boldsymbol{\Phi}_k) = \frac{|\mathbf{v}_k^\top \mathbf{G}_k \boldsymbol{\Phi}_k \mathbf{H} \mathbf{w}_k|^2 \sigma_x^2}{|\mathbf{v}_k|^2 \sigma_n^2}, \quad (4.2)$$

¹The assumption of single-stream transmissions is justified by the rank of the cascade channel matrix, which is likely equal to one. This conclusion comes from the considerations reported in [189–191], and has been verified numerically for the considered setup.

where σ_x^2 is the power of the transmitted signal. To maximize the SNR of a given UE, a specific IRS configuration should be adopted, tailored to the UE position in the cell and the channel conditions. However, the goal of this work is to limit the number of IRS reconfigurations to comply with realistic overhead constraints, as well as to improve the communication efficiency, and algorithms seeking to comply with these requirements will be presented in Section 4.5.

4.4 Sum Capacity Optimization Problem

We impose a constraint on the number of IRS reconfigurations per time frame, with the goal of either limiting the reconfiguration,² or accounting for practical limitations that might arise in realistic deployments. On the downside, achieving this objective usually leads to SNR degradation as sub-optimal IRS configurations might be adopted for some UEs. To mitigate this effect, we formulate a constrained optimization problem on the average cell sum capacity. Specifically, we assume the following conditions:

1. at most Z IRS reconfigurations can occur per time frame;
2. the gNB serves K UEs by partitioning them into Z disjoint subsets $\mathcal{U}_1, \dots, \mathcal{U}_Z, Z \leq K$;
3. for each UE in \mathcal{U}_z , the same IRS configuration $\Phi^{(z)}$ is used, i.e., $\Phi_k = \Phi^{(z)}, \forall k \in \mathcal{U}_z, \forall 1 \leq z \leq Z$.

Then, the achievable rate of UE $k \in \mathcal{U}_z$ is

$$R_k(\Phi^{(z)}) = \log_2 \left(1 + \Gamma_k(\Phi^{(z)}) \right), \quad (4.3)$$

where $\Gamma_k(\Phi^{(z)})$ is the SNR experienced by the k -th UE while configuration $\Phi^{(z)}$ is adopted at the IRS, i.e., the configuration shared by all UEs belonging to subset \mathcal{U}_z .

²We remark that the gNB typically communicates a (possibly new) IRS configuration in each TTI. The reconfiguration constraint introduced in the proposed IRS scheduling framework is able to reduce this overhead by a factor $Z/K \leq 1$.

Let $\mathcal{I} = \{\Phi^{(1)}, \Phi^{(2)}, \dots, \Phi^{(Z)}\}$ be the set of IRS configurations corresponding to subsets $\mathcal{U}_1, \dots, \mathcal{U}_Z$. The system sum capacity within a time frame is defined as

$$C(\mathcal{U}_1, \dots, \mathcal{U}_Z, \mathcal{I}) = B \sum_{z=1}^Z \sum_{k \in \mathcal{U}_z} R_k(\Phi^{(z)}), \quad (4.4)$$

where B is the transmission bandwidth. The optimization problem is then formulated as

$$\max_{\mathcal{U}_1, \dots, \mathcal{U}_Z, \mathcal{I}} C(\mathcal{U}_1, \dots, \mathcal{U}_Z, \mathcal{I}), \quad (4.5a)$$

$$\text{s.t. } \angle[\Phi^{(z)}]_{n,n} \in \mathcal{P}_\theta, \quad \forall n, z. \quad (4.5b)$$

Problem (4.5) determines the optimal grouping strategy for the UEs subsets $\mathcal{U}_1, \dots, \mathcal{U}_Z$, and assigns the best IRS configuration accordingly. Therefore, (4.5) is both continuous (i.e., the optimization of the IRS configuration) and combinatorial (i.e., the grouping of the UEs), and can be thus classified as a Mixed Integer Nonlinear Programming (MINLP) problem. Moreover, the following theorem holds.

Theorem 3. *The sum capacity maximization problem (4.5) is NP-complete.*

Proof: First, we observe that the problem falls within the general NP class. This is because if (4.5) is solved to find $\mathcal{U}_1, \dots, \mathcal{U}_Z, \mathcal{I}$, both the sum capacity and the phase-shift constraints (4.5b) could be verified in polynomial time. To prove that the problem is NP-complete, we set \mathcal{I} , and consider the simplified problem

$$\max_{\mathcal{U}_1, \dots, \mathcal{U}_Z} C(\mathcal{U}_1, \dots, \mathcal{U}_Z, \mathcal{I}). \quad (4.6)$$

This problem can be viewed as a multi-knapsack problem with different clusters $\mathcal{U}_1, \dots, \mathcal{U}_Z$ as knapsacks, and the goal is to maximize the total system capacity. This is known to be NP-hard, as it is a generalization of the classic knapsack problem. The original sum capacity maximization problem (4.5), where we consider the additional degrees of freedom of the IRS configurations, remains NP-hard, thus making the problem NP-complete. ■

Given the inherent problem complexity, we adopt heuristic clustering algorithms to obtain approximated, though close-to-optimal, solutions, as described in Section 4.5.

4.5 Heuristic Sum Capacity Maximization

In this section, we provide heuristic solutions to (4.5). First, we present two clustering-based approaches to identify and group UEs with a similar optimal IRS configuration. Then, we solve the scheduling problem on the identified clusters with a TDMA approach [192]. We compute the UEs clusters by first estimating the optimal individual IRS configurations, denoted as Φ_k^* , $1 \leq k \leq K$, i.e., the IRS configurations leading to the maximum capacity for each UE k , as described in Section 4.5.1. These configurations would solve (4.5) for $Z = K$, as in this case all UEs are served in a TDMA fashion and with their optimal IRS configuration. The phase coefficients of the optimal IRS configuration matrices are then chosen as the initial points of a procedure leveraging clustering algorithms in the N_I -dimensional space, as explained in Section 4.5.2.

4.5.1 Optimal Individual IRS Configurations

In MIMO systems, both the gNB and the UEs adopt properly tuned beamformers to match the signal transmissions and receptions to the spatial direction providing the highest channel gain [31]. For the optimization of the IRS configuration of each individual UE, we adopt a procedure similar to that presented in [185], focusing on single-stream transmissions and, without loss of generality, on UE k .

For a given IRS configuration, the optimal beamforming vectors v_k and w_k coincide with the singular vectors corresponding to the highest singular value of the wireless channel matrix. In particular, we calculate the SVD of the overall cascade channel matrix

$$\mathbf{G}_k \boldsymbol{\Phi}_k \mathbf{H} = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^\dagger, \quad (4.7)$$

where the right and left singular vectors of $\mathbf{G}_k \boldsymbol{\Phi}_k \mathbf{H}$ are the columns of \mathbf{V} and \mathbf{U} , and the corresponding singular values are the diagonal entries of $\boldsymbol{\Sigma}$.

In our formulation, the IRS configuration $\boldsymbol{\Phi}_k$ is one of the optimization variables. Indeed, given v_k and w_k , we can solve

$$\boldsymbol{\Phi}_k^* = \underset{\boldsymbol{\Phi}_k}{\operatorname{argmax}} \ R_k(\boldsymbol{\Phi}_k), \quad (4.8a)$$

$$\text{s.t. } \angle[\boldsymbol{\Phi}_k]_{n,n} \in \mathcal{P}_\theta, \quad \forall n, \quad (4.8b)$$

where R_k is the achievable rate of user k , $1 \leq k \leq K$, according to (4.3).

The derivation of the optimal IRS configuration of user k requires the alignment of the channel phase coefficients. According to [193], the cascade channel can be expressed as

$$\begin{aligned} \mathbf{v}_k^T \mathbf{G}_k \boldsymbol{\Phi}_k \mathbf{H} \mathbf{w}_k &= \text{vec}\left(\mathbf{v}_k^T \mathbf{G}_k \boldsymbol{\Phi}_k \mathbf{H} \mathbf{w}_k\right) \\ &= \left(\mathbf{w}_k^T \mathbf{H}^T \diamond \mathbf{v}_k^T \mathbf{G}_k\right) \text{diag}(\boldsymbol{\Phi}_k) \end{aligned} \quad (4.9)$$

where \diamond denotes the Khatri-Rao product operator and $\text{diag}(\boldsymbol{\Phi}_k)$ return the column vector with all the elements in the diagonal of $\boldsymbol{\Phi}_k$. Then, it is sufficient to observe that the SNR is maximized when the phase shifts introduced by the IRS are aligned with the phase shifts accumulated along the various paths, i.e.,

$$\theta_{k,n} = -(\angle \left[\left(\mathbf{w}_k^T \mathbf{H}^T \diamond \mathbf{v}_k^T \mathbf{G}_k \right) \right]_n), \quad \forall n. \quad (4.10)$$

Note that, in general, we need to know the estimated phase shift of each component resulting from the Khatri-Rao product in (4.10), rather than the exact phase coefficients of \mathbf{H} and \mathbf{G}_k . Moreover, as pointed out in [193], for structured channel models adopted at mmWaves, where multipath scattering is sparse and propagation is often dominated by strong specular components, the estimation of the separated channel matrices \mathbf{H} and \mathbf{G}_k can be simply accommodated by optimizing a limited number of parameters.

Taking into account the possible quantization, the optimal phase shifts are given by

$$\angle[\boldsymbol{\Phi}_k^*]_{n,n} \leftarrow \underset{\psi \in \mathcal{P}_\theta}{\text{argmin}} \left(\angle e^{j(\theta_{k,n} - \psi)} \right), \quad \forall n. \quad (4.11)$$

To overcome the interdependence between optimal IRS configurations and beamforming vectors, we propose an iterative alternate optimization approach. We first estimate the optimal beamforming vectors for a given IRS configuration using (4.7). Then, we plug the derived beamformers into (4.8a), and obtain the corresponding optimal IRS configuration. We repeat this two-step procedure until convergence, which, for practical purposes, is assumed to be reached when the difference between the achievable rates R_k , $\forall k$, in two consecutive iterations is lower than a tolerance $\nu > 0$. This procedure is summarized in Algorithm 6, where t is the iteration index. The number of iterations grows with the numbers of antennas and IRS phase shifters. How-

Algorithm 6 Iterative Alternate IRS Optimization

Input: \mathbf{G}_k, \mathbf{H}
Output: Φ_k^*

- 1: $t \leftarrow 0$
- 2: $\mathbf{v}_k, \mathbf{w}_k \leftarrow \mathbf{1}$
- 3: **repeat**
- 4: $\theta_{k,n} \leftarrow -(\angle[(\mathbf{w}_k^\top \mathbf{H}^\top \diamond \mathbf{v}_k^\top \mathbf{G}_k)]_n), \quad \forall n$
- 5: $\angle[\Phi_{k,t}]_{n,n} \leftarrow \operatorname{argmin}_{\psi \in \mathcal{P}_\theta} (\angle e^{j(\theta_{k,n} - \psi)})$
- 6: $\mathbf{U}, \mathbf{\Sigma}, \mathbf{V}^\dagger \leftarrow \text{SVD of } \mathbf{v}_k^\top \mathbf{G}_k \Phi_k \mathbf{H} \mathbf{w}_k$
- 7: $\mathbf{v}_k \leftarrow \text{column of } \mathbf{V} \text{ corresponding to}$
 the largest singular value
- 8: $\mathbf{w}_k \leftarrow \text{column of } \mathbf{U} \text{ corresponding to}$
 the largest singular value
- 9: $t \leftarrow t + 1$
- 10: **until** $|R_k(\Phi_{k,t}) - R_k(\Phi_{k,t-1})| < \nu$
- 11: $\Phi_k^* \leftarrow \Phi_{k,t}$

ever, from preliminary simulations, and based on the set of parameters we considered (see Section 4.8), convergence is typically reached in less than 10 iterations.

4.5.2 Clustering-based TDMA Scheduling

For an approximated but close-to-optimal solution to (4.5), we resort to a clustering-based approach. Our proposed clustering algorithms estimate both the subsets of UEs $\mathcal{U}_1, \dots, \mathcal{U}_Z$, and the relative set of IRS configurations \mathcal{I} . We operate on the *phase vector space*, i.e., the points to be clustered are identified by the IRS phase shifts vector $[\angle\phi_1, \dots, \angle\phi_{N_l}]^\top = [\theta_1, \dots, \theta_{N_l}]^\top$, which maps each IRS configuration Φ to a point in $[-\pi, \pi]^{N_l}$. In case of quantized phase shifts, the phase vector space is a lattice in the continuous space $[-\pi, \pi]^{N_l}$.

The general clustering-based procedure works as follows:

- *Step 1:* find $\Phi_k^*, \forall k$, i.e., the optimal individual IRS configurations for each UE as in Section 4.5.1;
- *Step 2:* build UE subsets $\mathcal{U}_z, z = 1, \dots, Z$, by using a clustering algorithm, according to Sections 4.6 and 4.7;
- *Step 3:* assign $\Phi^{(z)}$ to all UEs $\in \mathcal{U}_z$.

4.6 Distance-Based Clustering Algorithms

The core idea of this procedure is to use clustering algorithms to group UEs, and assign the respective IRS configurations, which are mapped to the *centroid* of the cluster. In the case of quantized phase shifts, once the clustering procedure is performed, clusters may share the same centroid and be merged. Therefore, Z represents the *maximum number of clusters*, not the effective number. Moreover, we remark that the procedure above does not rely on the assumption of perfect CSI, as the grouping strategy (Step 2) and the individual optimization step (Step 1) are performed independently. Nevertheless, in the case of imperfect CSI, the estimated individual optimal configurations may differ from the actual optimal configurations, leading to a suboptimal grouping.

In the following, we propose different techniques to build the clusters based on either a distance metric (Section 4.6) or the achievable rate (Section 4.7).

4.6 Distance-Based Clustering Algorithms

The class of distance-based clustering contains methods that group data points based on their similarity or dissimilarity according to a distance metric. This approach has several advantages, including the efficiency in handling large datasets, and the flexibility to adapt to many different scenarios of interest. However, distance-based clustering can be sensitive to the choice of the distance metric (which depends on the nature of the data and the clustering problem), and the initialization values. Moreover, in our specific case, it does not take into account the achievable rate, which is not directly related to the distance among the points in the phase vector space.

Since the scalar field is the range $[-\pi, \pi]$, the adopted distance has to take into account the circularity of data. However, the convergence to a local minimum for most of the clustering algorithms is guaranteed only if the points to be clustered belong to a Euclidean space. For distance-based algorithms, we thus define the bijective mapping function $f : \mathcal{P}_\theta^{N_1} \rightarrow \mathbb{R}^{2N_1}$ as

$$\begin{aligned} f(\boldsymbol{\theta}) &= f([\theta_1, \dots, \theta_{N_1}]) \\ &= [\cos(\theta_1), \sin(\theta_1), \dots, \cos(\theta_{N_1}), \sin(\theta_{N_1})], \end{aligned} \quad (4.12)$$

and then define the pairwise distance between two generic IRS configurations α and β as

$$\delta(\alpha, \beta) = ||f(\alpha) - f(\beta)||, \quad (4.13)$$

i.e., the Euclidean distance between the mapping on the unit N_I -sphere of their respective phase vectors. In the following, with a slight abuse of notation, $f(\Phi)$ maps the phases of the complex entries in the diagonal of Φ as in (4.12), and $\delta(\Phi_1, \Phi_2)$ denotes the pairwise distance between the phases of the elements in the diagonal of matrices Φ_1 and Φ_2 . The sum of squared distances is defined as

$$J(\mathcal{U}_1, \dots, \mathcal{U}_Z, \mathcal{I}) = \sum_{z=1}^Z \sum_{k \in \mathcal{U}_z} \delta(\Phi_k^*, \Phi^{(z)})^2, \quad (4.14)$$

and the distance-based clustering schemes are used to solve the following problem:

$$\min_{\mathcal{U}_1, \dots, \mathcal{U}_Z, \mathcal{I}} J(\mathcal{U}_1, \dots, \mathcal{U}_Z, \mathcal{I}), \quad \text{s.t. } (4.5b). \quad (4.15)$$

We consider and compare some of the most popular distance-based clustering algorithms, namely, K-means, agglomerative hierarchical clustering, and K-medoids.

K-Means (KM). KM clustering [194] aims at finding Z disjoint clusters minimizing the within-cluster squared Euclidean distances. Here, we consider the generalized Lloyd algorithm [195], which randomly selects Z points in the space of phase vectors as the initial centroids. In our setup, to ensure optimal performance when $Z = K$, we force the algorithm initialization to a random selection among the phase vectors of the optimal individual IRS configurations derived in Section 4.5.1. Then, in the *assignment step* KM assigns each data point to the closest centroid, according to the specified distance metric. In the subsequent *update step*, the set of centroids is re-computed as the average of the data points that belong to each cluster. These steps are repeated until either convergence or a maximum number of iterations I_{\max}^{KM} is reached.

Agglomerative Hierarchical Clustering (HC). The agglomerative HC [196] partitions a set of data points into disjoint clusters by iteratively merging points into clusters until a target number of partitions is met. In our setup, clusters are initialized as the optimal phase vectors, which thus act as the respective

centroids. Then, the average distance between all pairs of data points in any pair of clusters is evaluated. The closest pair of clusters are merged into a new single cluster, whose centroid is computed as the mean of its data points. The procedure is repeated until the number of clusters is Z .

K-Medoids (KMed). KMed [197] is a clustering technique similar to KM, but instead of the mean of the data points within each cluster, it uses the medoid, i.e., the data point that is closest to the center of the cluster. In our setup, we consider the Partition Around Medoids (PAM) method [198], which starts by randomly selecting Z medoids among the optimal phase vectors and assigns each point to the cluster with the closest medoid. In each iteration, the algorithm evaluates potential *swaps* of medoids with non-medoids. A swap is accepted only if it results in a lower value of the sum of the squared distances to all other data points within the same cluster. The algorithm continues until the medoids no longer change.

Theorem 4. *The proposed distance-based clustering techniques converge to a local minimum of (4.14).*

Proof: The proof directly derives from the well-known results of clustering with Euclidean distance. The exact proofs for each of the considered algorithms under distance metric (4.13) are reported in the Appendix. ■

4.7 Capacity-Based Clustering Algorithms

The distance-based clustering techniques presented in Section 4.6 do not directly take into account the actual capacity achievable by the UEs, which is a crucial factor for the sum capacity maximization (4.5). Thus, in the following we propose original capacity-based clustering algorithms that go beyond the state of the art, namely CWC (Section 4.7.1), OSCBC (Section 4.7.2), and ICWC (Section 4.7.3).

4.7.1 Capacity-Weighted Clustering (CWC)

Similarly to distance-based clustering, also CWC proceeds iteratively. However, the stopping condition is based on the variation of the sum capacity of each cluster, rather than on the distance between the centroids. In this approach, the clustering algorithm itself weighs the UEs based on their achievable capacity, so that the parameters of the resulting clusters are closer to

those preferred by the UEs with higher rates, thus promoting the maximization of the sum capacity.

Let $\Phi_i^{(z)}$ be the IRS configuration of cluster $\mathcal{U}_{z,i}$ at iteration i . UEs are initially sorted in decreasing order of achievable rate. The algorithm then selects the Z UEs providing the highest achievable capacity based on the expression in (4.3) with their optimal IRS configurations. Without loss of generality, we let $z = 1, \dots, Z$ be the index of those UEs, and set $\Phi_1^{(z)} = \Phi_z^*$, $\forall 1 \leq z \leq Z$, as the centroids of the initial clusters $\mathcal{U}_{1,0}, \dots, \mathcal{U}_{Z,0}$. In the following, for simplicity, we denote with $z_{k,i}$ the cluster such that $k \in \mathcal{U}_{z,i}$. Each UE $k > Z$ is assigned to the cluster whose centroid provides the lowest rate difference with respect to its ideal configuration. Let $R_k(\Phi_k^*)$ be the maximum achievable rate of UE k , obtained from the solution of problem (6). UE k is assigned to cluster

$$z_{k,i} = \operatorname{argmin}_z [R_k(\Phi_k^*) - R_k(\Phi_i^{(z)})], \quad (4.16)$$

where $R_k(\Phi_i^{(z)})$ is the rate achieved by UE k adopting the IRS configuration of cluster z at iteration i . Note that, despite being always non-negative, the rate difference in (4.16) cannot be considered a distance metric as, in general, it does not satisfy the triangle inequality. However, we prove that, as the distance from the optimal configuration increases, the corresponding rate decreases, thus supporting the use of the rate difference as a clustering criterion.

Theorem 5. *Given the optimal IRS configuration Φ_k^* , the rate $R_k(\Phi)$ is monotonically decreasing with respect to the magnitude of any phase shifts error ϵ .*

Proof: Let $\epsilon \in [-\pi, \pi]$ be an arbitrary error phase shift, and consider the configuration $\Phi_k^\epsilon = \Phi_k^* \mathbf{E}$, where $\mathbf{E} = \operatorname{diag}(e^{j\epsilon}, 1, \dots, 1)$, i.e., the suboptimal configuration where only the first IRS element is affected by the error ϵ . Assuming, without loss of generality, that $N_g = N_k = 1$ and $\sigma_x = \sigma_n = 1$, the rate $R_k(\Phi_k^\epsilon)$ is proportional to $\Gamma_k(\Phi_k^\epsilon)$ when using configuration Φ_k^ϵ . The SNR $\Gamma_k(\Phi_k^\epsilon)$ can be written as

$$\Gamma_k(\Phi_k^\epsilon) = |\mathbf{g}_k \Phi_k^* \mathbf{E} \mathbf{h}|^2 \quad (4.17)$$

$$= |[\mathbf{g}_k]_1 [\Phi_k^* \mathbf{E}]_{1,1} [\mathbf{h}]_1 + A|^2, \quad (4.18)$$

where $A = \sum_{n=2}^{N_I} [\mathbf{g}_k]_n [\Phi_k^*]_{n,n} [\mathbf{h}]_n$. Since Φ_k^* is the optimal configuration, it satisfies (4.10). It follows that $A \in \mathbb{R}^+$, so (4.17) can be further manipulated into

$$\Gamma_k(\Phi_k^\epsilon) = \|[\mathbf{g}_k]_1\| \|[\mathbf{h}]_1\| e^{j\epsilon} + A|^2 \quad (4.19)$$

$$= A^2 + (\|[\mathbf{g}_k]_1\| \|[\mathbf{h}]_1\|)^2 + 2A \|[\mathbf{g}_k]_1\| \|[\mathbf{h}]_1\| \cos(\epsilon). \quad (4.20)$$

Finally, we evaluate the sign of the derivative of $\Gamma_k(\Phi_k^\epsilon)$ with respect to the error ϵ as

$$\frac{\partial \Gamma_k(\Phi_k^\epsilon)}{\partial \epsilon} = -2A \|[\mathbf{g}_k]_1\| \|[\mathbf{h}]_1\| \sin(\epsilon), \quad (4.21)$$

and observe that $\Gamma_k(\Phi_k^\epsilon)$, and therefore $R_k(\Phi_k^\epsilon)$, is strictly decreasing for $0 < |\epsilon| \leq \pi$. ■

After all the remaining UEs have been assigned to the corresponding clusters, the coordinates of the centroids are updated. At iteration $i + 1$, the new IRS configuration (centroid) of cluster $\mathcal{U}_{z,i+1}$ is computed as the average of the data points in the cluster, weighted by their achievable rate, i.e.,

$$\Phi_{i+1}^{(z)} = f^{-1} \left(\frac{\sum_{k \in \mathcal{U}_z} f(\Phi_k^*) R_k(\Phi_k^*)}{\sum_{k \in \mathcal{U}_z} R_k(\Phi_k^*)} \right). \quad (4.22)$$

Also, in the case of phase shift quantization, an additional approximation step must be performed as

$$\angle[\Phi_{i+1}^{(z)}]_{n,n} \leftarrow \operatorname{argmin}_{\psi \in \mathcal{P}_\theta} \left(\angle e^{j(\angle[\Phi_{i+1}^{(z)}]_{n,n} - \psi)} \right), \quad \forall n. \quad (4.23)$$

This two-step procedure is repeated until convergence, which is reached when the rate difference between two consecutive iterations is lower than the sum capacity tolerance $\mu > 0$.

The rationale behind the algorithm is that, based on the initial centroid assignment, the UEs experiencing the best channel conditions, i.e., those dominating the system sum capacity, are initially served with their optimal (individual) IRS configurations. Even after the adjustment of the clusters, these UEs will always get the largest weight coefficient within the cluster. The remaining UEs, instead, will be penalized by the configuration constraints, but their impact on the sum capacity will be limited. The whole workflow of the CWC procedure is summarized in Algorithm 7.

Algorithm 7 CWC Algorithm

Input: $Z, \mathbf{H}, \mathbf{G}_k, \forall k$
Output: $\mathcal{U}_1, \dots, \mathcal{U}_Z, \mathcal{I}$

- 1: Compute Φ_k^* , $\forall k$ with the procedure of Algorithm 6
- 2: Sort the UEs in decreasing order of $R_k(\Phi_k^*)$
- 3: Select the Z UEs providing the highest $R_k(\Phi_k^*)$,
- 4: Set $\Phi_1^{(z)} = \Phi_k^*, z = 1, \dots, Z$ as the initial centroids.
- 5: **repeat**
- 6: **for** each UE k **do**
- 7: $z_{k,i} \leftarrow \operatorname{argmin}_z R_k(\Phi_k^*) - R_k(\Phi_i^{(z)})$
- 8: **end for**
- 9: **for** each cluster z **do**
- 10: Compute $\Phi_{i+1}^{(z)}$ as per (4.22), (4.23)
- 11: **end for**
- 12: $i \leftarrow i + 1$
- 13: **until** $\left| \sum_{k \in \mathcal{U}_z} R_k(\Phi_i^{(z)}) - \sum_{k \in \mathcal{U}_z} R_k(\Phi_{i-1}^{(z)}) \right| < \mu$
- 14: Assign $\Phi^{(z)}$ to all $k \in \mathcal{U}_z$.

4.7.2 One-Shot Capacity-Based Clustering (OSCBC)

The main drawback of CWC presented in Section 4.7.1 is that it requires solving problem (4.16) at each iteration, relative to all the UEs in each cluster. Considering massive MIMO systems, the CWC procedure could become exceedingly complex, as it requires the SVD computation of extremely large matrices. Therefore, we propose another lower-complexity clustering algorithm, denoted as OSCBC.

As in CWC, also in OSCBC: (i) the UEs are sorted in decreasing order of achievable rate; (ii) the Z IRS configurations of the Z UEs experiencing the highest rates are chosen as initial centroids for the clusters; and (iii) the remaining UEs are assigned to the closest centroid in terms of circular distance, as per (4.13). Then, compared to CWC, instead of recomputing the coordinates of the centroids at each iteration until convergence, the algorithm stops right after the initial association. Therefore, with OSCBC the computed centroids are the optimal configurations relative to the Z UEs achieving the highest individual rate, which provides suboptimal (non-optimized) performance for the rest of the UEs in the clusters.

Table 4.1: Computational complexity of distance-based vs. capacity-based clustering.

Clustering algorithm	Computational complexity
KM (Lloyd)	$O(IZKN_I)$
KMed (PAM)	$O(Z^3K^2N_I)$
HC	$O(K^3N_I)$
CWC/ICWC	$O(IZKN_gN_I^2)$
OSCBC	$O(Z(K - Z)N_I)$

4.7.3 Inverse Capacity-Weighted Clustering (ICWC)

The CWC algorithm is designed to optimize the capacity of the UEs experiencing the best channel conditions and is unfair to the other UEs in the system, which may use suboptimal IRS configurations. Therefore, we propose an additional variation of CWC, named ICWC, with the goal of achieving higher fairness among the UEs in the system. In ICWC, while the cluster association principle of (4.16) is preserved, the initial condition is reversed. Specifically: (i) UEs are sorted in increasing order of achievable rate; (ii) the initial configurations of the clusters $\Phi_1^{(z)} = \Phi_z^*, z = 1, \dots, Z$, are based on the optimal configurations of the UEs with the worst channel conditions. The remaining $k > Z$ UEs are associated as in (4.16). Then, at iteration i , the IRS configuration is updated as

$$\Phi_{i+1}^{(z)} = f^{-1} \left(\frac{\sum_{k \in \mathcal{U}_z} \Phi_k^* R_k^{-1}(\Phi_k^*)}{\sum_{k \in \mathcal{U}_z} R_k^{-1}(\Phi_k^*)} \right), \quad (4.24)$$

and the discretization step (4.23) is performed (if needed). As in CWC, convergence is achieved if the rate difference between two consecutive iterations is lower than the tolerance μ . While ICWC obtains lower sum capacity than CWC, it can provide significant improvements in terms of fairness, especially from the perspective of the UEs with the worst channel conditions.

4.7.4 Computational Complexity

The computational complexity is evaluated as the number of iterations required for the clustering algorithms to: (i) obtain the optimal IRS configura-

tion of each UE; (ii) partition the UEs into disjoint subsets, or clusters, based on distance or capacity metrics; and (iii) for each cluster, find the best IRS configuration to serve the corresponding UEs.

Specifically, at each iteration, the main source of complexity is the computation of the overall cascade channel matrix $\mathbf{G}_k \Phi_k \mathbf{H}$, which has complexity $O(N_g N_I^2 + N_g N_I N_U)$. Additionally, in the case of quantized IRS phase shifts, after obtaining the optimal beamformers, the optimal phase shifts for the IRS are obtained through an exhaustive search over the set of possible phase shifts \mathcal{P}_θ , yielding a complexity $O(2^b N_I)$.

Notice that different clustering algorithms, in general, require a different number of iterations I to reach convergence, thus possibly introducing practical limitations. Moreover, the complexity introduced in each iteration depends on the clustering algorithm itself. In Table 4.1 and in the following text we characterize the computational complexity of each of the clustering algorithms presented in Sections 4.6 and 4.7.

Distance-based clustering. These algorithms do not require specific initialization. For KM, based on the Lloyd implementation in [195], each iteration involves calculating the distances between data points and centroids. As a result, the computational complexity is influenced by the number of iterations required for convergence, the number of data points, the number of clusters, and the dimensionality of data, resulting in an overall complexity $O(I Z K N_I)$. KMed can be solved with the PAM algorithm [198], so the computational complexity is $O(Z^3 K^2 N_I)$ due to the pairwise distance computations between data points and medoids. Finally, the computational complexity of the agglomerate HC is primarily determined by the computation of pairwise distances among all data points, resulting in a total complexity $O(K^3 N_I)$ [199].

Capacity-based clustering. The complexity of the OSCBC algorithm is dominated by the centroid assignment upon initialization, which has complexity $O(Z(K - Z) N_I)$. Instead, for the CWC and ICWC algorithms, the complexity is $O(I Z K N_g N_I^2)$, as demonstrated in the following theorem.

Theorem 6. *The time complexity of CWC and ICWC scales quadratically with N_I as $O(I Z K N_g N_I^2)$.*

Proof: Capacity-based clustering requires an initialization stage where the algorithm selects the Z UEs providing the highest (or lowest) $R_k(\Phi_k^*)$,

resulting in a complexity of $O(K \log K)$ due to the sorting of K scalars. In the subsequent iterations:

1. Both CWC and ICWC compute the rate difference between each UE and the Z centroids. The complexity of computing $R_k(\Phi_i^{(z)})$ can be dominated either by the matrix multiplication in (4.2), or by the SVD for the single stream beamforming which require, respectively, $O(N_g N_I^2 + N_g N_I N_U)$ and $O(N_g N_U \min(N_g, N_U))$ operations for each UE and each centroid.
2. The computation of the centroids as per (4.22)-(4.24) requires $N_I + 1$ scalar operations per UE, which has negligible complexity with respect to the rate computation.

In typical IRS-assisted systems, $N_I \gg N_g > N_U$. Therefore, the complexity at each iteration is dominated by the channel matrix product, and the overall algorithm complexity is $O(IZKN_g N_I^2)$. ■

4.8 Numerical Results

After presenting our various simulation scenarios and evaluation metrics in Sections 4.8.1 and 4.8.2, respectively, we assess in Section 4.8.3 the scheduling performance of an IRS-assisted network with practical constraints.

4.8.1 Simulation Parameters

Our simulation parameters are reported in Table 4.2.

Scenario. All devices are assumed to lie on a 2D plane, and we consider an Urban Microcell (UMi) scenario, according to the 3GPP nomenclature [23], with the gNB placed at the center. According to the 3GPP specifications, the coverage area of the gNB is characterized by an average radius of 167 m and is assumed to lie in the positive x -axis region.

We assume that $K = 100$ UEs are randomly deployed according to a uniform distribution within the cell area, to be served in downlink by the gNB, assisted by an IRS at coordinates (75, 100) m. The gNB is equipped with a UPA with $8H \times 8V$ antennas (i.e., $N_g = 64$), and the UEs with Uniform Linear Arrays (ULAs) of $2H \times 1V$ antennas (i.e., $N_U = 2$). For the IRS, if not otherwise specified, we adopt a $40H \times 80V$ reflective panel ($N_I = 3200$).

Table 4.2: Simulation parameters.

Parameter	Value
Carrier frequency	28 GHz
Total bandwidth (B)	100 MHz
Noise power spectral density	-174 dBm/Hz
Number of UEs (K)	100
gNB antenna array (N_g)	8H×8V
gNB transmit power	33 dBm
UE antenna array (N_U)	2H×1V
IRS elements (N_I)	{10H×20V, 20H×40V, 40H×80V, 60H×120V}
Phase shift quant. bits (b)	{unquantized, 1-bit, 2-bits}
LoS probability (p_{LoS})	Eq. (4.26)
Individual rate opt. tolerance (ν)	10^{-6} [bit/s/Hz]
KM max. iterations (I_{max}^{KM})	50
CWC/ICWC rate tolerance (μ)	10^{-3} [bit/s]

Channel and Frame Structure. The system operates at a carrier frequency of 28 GHz (that is in the lower part of the mmWave bands), the transmission power at the gNB is set to 33 dBm, the noise power spectral density at the receivers is -174 dBm/Hz, and the total system bandwidth is 100 MHz. We consider the fourth numerology of the NR frame structure [200], wherein each 10 ms frame is split into 160 slots. With this assumption, as already pointed out in Section ??, channels can be considered constant over the entire frame duration. We consider the 3GPP TR 38.901 spatial channel model [23], which supports a wide range of frequencies, from 0.5 to 100 GHz (and including therefore our carrier frequency of 28 GHz), and can be integrated with realistic beamforming models. As such, channel matrices, and multipath fading, are computed based on the superposition of N different clusters, each of which consists of M rays that arrive (depart) to (from) the antenna arrays with specific angles and powers. Based on [23], and using the simplifications

proposed in [58], the generic entry $[A]_{pq}$ of the channel matrix can then be computed as:

$$\begin{aligned} [A]_{pq} = \gamma \sum_{n=1}^N \sqrt{\frac{P_n}{M}} \sum_{m=1}^M & \bar{\mathbf{F}}_{rx}\left(\theta_{n,m}^A, \phi_{n,m}^A\right) \\ & \times \begin{bmatrix} e^{j\Phi_{n,m}^{\theta,\theta}} & \sqrt{K_{n,m}^{-1}} e^{j\Phi_{n,m}^{\theta,\phi}} \\ \sqrt{K_{n,m}^{-1}} e^{j\Phi_{n,m}^{\phi,\theta}} & e^{j\Phi_{n,m}^{\phi,\phi}} \end{bmatrix} \\ & \times \bar{\mathbf{F}}_{tx}\left(\theta_{n,m}^D, \phi_{n,m}^D\right) \\ & \times e^{j\bar{\mathbf{k}}_{rx,n,m}^T \bar{\mathbf{d}}_{rx,p} e^{j\bar{\mathbf{k}}_{tx,n,m}^T \bar{\mathbf{d}}_{tx,q}}}, \end{aligned} \quad (4.25)$$

where γ is the Large-Scale Fading Coefficient (LSFC) of the considered link, which incorporates the path loss and shadowing terms. For a complete description of the remaining terms appearing in (4.25) we refer the interested reader to [58]. Specifically, while the gNB and the IRS can be assumed to operate in LoS, the path loss between a generic UE k and the IRS is modeled based on the following channel conditions:

- *NLoS*: UE k is in NLoS with the IRS;
- *deterministic LoS* (LoS): UE k is in LoS with the IRS;
- *probabilistic LoS* (LoS): the IRS-UE k link is in LoS ∨ NLoS with respective probabilities $p_k^{\text{LoS}}(d_k) \vee 1 - p_k^{\text{LoS}}(d_k)$, with

$$p_k^{\text{LoS}}(d_k) = \begin{cases} 1 & \text{if } d_k \leq 18, \\ \frac{18}{d_k} + \left(1 - \frac{18}{d_k}\right) e^{-\frac{d_k}{36}} & \text{if } d_k > 18, \end{cases} \quad (4.26)$$

where d_k is the distance (in m) between the IRS and UE k . In the considered UMi scenario, and based on 3GPP specifications [23], the average LoS probability in (4.26) is 0.35.

For each wireless link, based on the presence of the LoS component, the path loss is then derived according to [23, Table 7.4.1-1], with shadowing standard deviation set to $\sigma_{SF} = 0$. For the optimal individual IRS configuration (Section 4.5.1), we set $\nu = 10^{-6}$ [bit/s/Hz].

Clustering algorithms. In the following subsections, we present extensive simulation results to compare the performance of distance-based (KM, HC, KMed) vs. capacity-based (CWC, OSCBC, ICWC) clustering algorithms to perform scheduling in an IRS system with reconfiguration constraints. The

KM clustering has been implemented with the Lloyd algorithm [195] with a maximum number of $I_{\max}^{KM} = 50$ iterations. Instead, for both CWC and ICWC, we set $\mu = 10^{-3}$ [bit/s].

As an upper bound to the system performance, we also consider an “unclustered” scheduling, wherein we assume that all UEs are served with their optimal IRS configuration. This scheduling clearly violates the constraint on the maximum numbers of reconfiguration per frame, but can be regarded as the limit case when $Z = K$, i.e., all UEs belong to a cluster with cardinality one. As such, it is a suitable approach for benchmarking the performance of more practical schemes.

4.8.2 Performance Metrics

The performance of the proposed clustering-based scheduling techniques is evaluated in terms of average sum capacity and fairness, as a function of the numbers of both clusters and UEs, under different channel conditions, IRS dimensions, and degrees of quantization for the phase shifts.

Average sum capacity. It is derived from (4.4) as

$$\bar{C} = \frac{1}{K} \mathbb{E}[C(\mathcal{U}_1, \dots, \mathcal{U}_Z, \mathcal{I})], \quad (4.27)$$

where the expectation is computed across the different channel realizations. Moreover, as each UE is served in its specific slot, we average over the TDMA frame length, dividing the empirical expectation by the number of UEs (slots) K .

Fairness. We consider the 95% percentile of the achieved individual user capacity, computed as

$$C_{95\%} = \frac{B}{K} \inf\{x : CDF(x) \geq 0.95\}, \quad (4.28)$$

where $CDF(\cdot)$ is the empirical cumulative distribution function of $R_k(\Phi^{(z)})$, $\forall k, z$. Notice that the 95% percentile of the user capacity is a practical and meaningful way to evaluate fairness, as it measures the performance of the majority of the UEs, excluding only the top 5%.

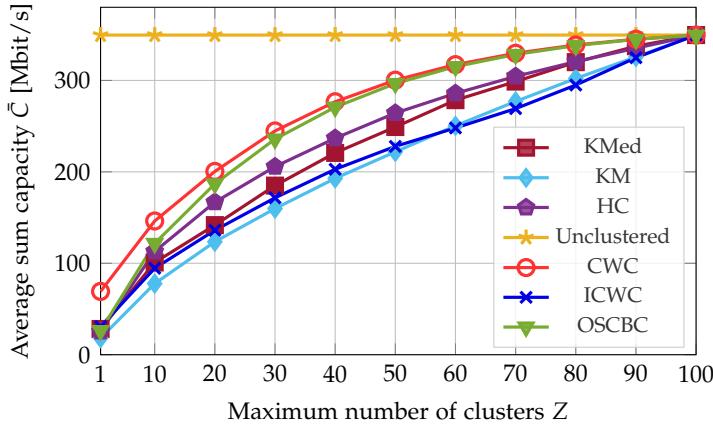


Figure 4.2: Average sum capacity as a function of the maximum number of clusters Z , for an unquantized $40H \times 80V$ IRS, and considering a pLoS channel for the IRS-UEs links.

4.8.3 Scheduling Performance

In this section, we compare the IRS scheduling performance considering distance-based vs. capacity-based clustering, and as a function of different channel conditions, reconfiguration constraints, and degrees of quantization of the phase shifts.

Impact of the clustering algorithm. First, Fig. 4.2 displays the average sum capacity \bar{C} per slot as a function of the number of clusters Z , for unquantized IRS phase shifts, and considering a pLoS channel for the IRS-UEs links. It is evident that all the scheduling policies perform better whenever Z increases, and converge to the “unclustered” policy when $Z = K$. In fact, increasing the number of clusters corresponds to a smaller intra-cluster average distance, which eventually becomes zero when $Z = K$. Among the considered clustering policies, CWC and OSCBC provide the highest sum capacity, as they are designed to maximize \bar{C} , and choose the IRS configurations of the UEs that achieve the highest rate. Instead, distance-based clustering achieves worse performance as it does not exploit the knowledge of the rate achievable with different IRS configurations when building the clusters. As expected, ICWC is designed to promote fairness, thus performs worse than both CWC and OSCBC in terms of sum capacity; still, it achieves similar performance as distance-based clustering. Finally, the gap between CWC and OSCBC is almost negligible: this implies that a single iteration in the clustering process

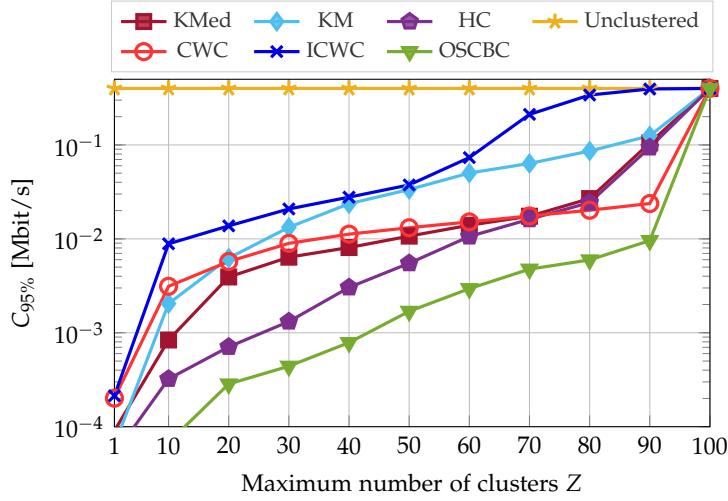


Figure 4.3: 95% percentile of the user capacity as a function of the maximum number of clusters Z , for an unquantized $40H \times 80V$ IRS, and considering a pLoS channel for the IRS-UEs links.

is enough to achieve good sum capacity, while also promoting lower computational complexity as reported in Table 4.1, which demonstrates the good scalability of the proposed techniques.

Fig. 4.3 compares the fairness performance of the different clustering algorithms, measured as the 95% percentile of the average sum capacity $C_{95\%}$, as a function of the maximum number of clusters Z in pLoS conditions. Our results identify ICWC as the best clustering approach in terms of fairness, which comes at the cost of a lower sum capacity, as shown in Fig. 4.2. Therefore, there exists a trade-off between the achievable sum capacity and fairness. We also observe that OSCBC achieves very low fairness, as the UEs with worst channel conditions are forced to aggregate to the strongest UEs, thus via a suboptimal IRS configuration. On the other hand, we see that CWC is more than acceptable in terms of fairness, and achieves comparable performance than most of the distance-based clustering algorithms. Furthermore, $C_{95\%}$ increases as Z increases, and eventually approaches the “unclustered” baseline for $Z = K$. This is due to the fact that the LoS probability in the pLoS scenario increases with the number of clusters, i.e., as the inter-cluster distance becomes smaller, which permits to experience better channel conditions, thus a higher capacity, even for the worst UEs. Finally, despite performing worse than their capacity-based counterparts, the distance-based

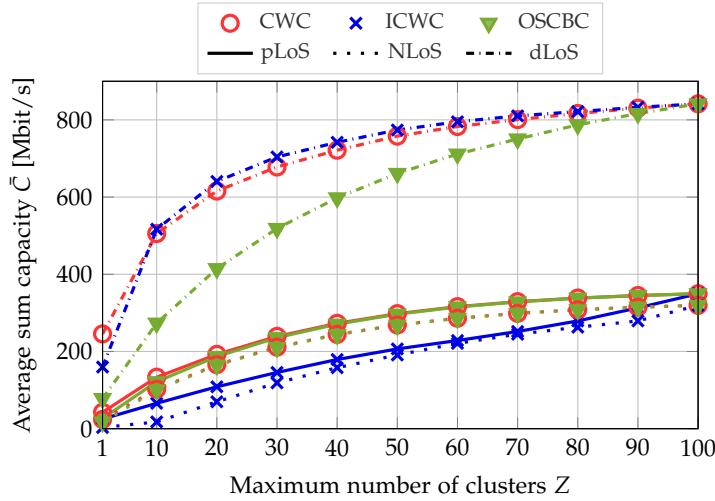


Figure 4.4: Average sum capacity as a function of the maximum number of clusters Z , for $N_I = 3200$, unquantized phase shifts, and for different channel conditions.

methods are a viable alternative for constrained IRS control nodes thanks to their lower computational complexity. In such cases, HC is to be preferred for sum-capacity maximization, while KM is the best alternative to capacity-based algorithms when the 95% percentile of the average sum capacity represents the metric of interest.

Impact of the channel. From the above results, we concluded that distance-based clustering provides lower sum capacity and fairness compared to capacity-based scheduling, so the rest of our simulation campaign has been focused on the latter. Figs. 4.4 and 4.5 display the average sum capacity and the 95% percentile, respectively, for CWC, ICWC, and OSCBC in different channel conditions. First, we observe that in the dLoS scenario, where UEs are in LoS with the IRS, the sum capacity is up to 2.6 (2.4) times higher than in the NLoS (pLoS) scenario for $Z = K$. This is mainly due to the fact that NLoS links experience (i) a higher path loss, and (ii) the lack of a dominant multipath component, thus of a clear steering direction for the IRS beam, which deteriorates the link quality. In particular, in the pLoS scenario the LoS probability decreases exponentially with the distance, therefore, the UEs that are far from the IRS typically operate in NLoS. For similar reasons, both CWC and ICWC in the dLoS scenario start to reach stability in terms of ca-

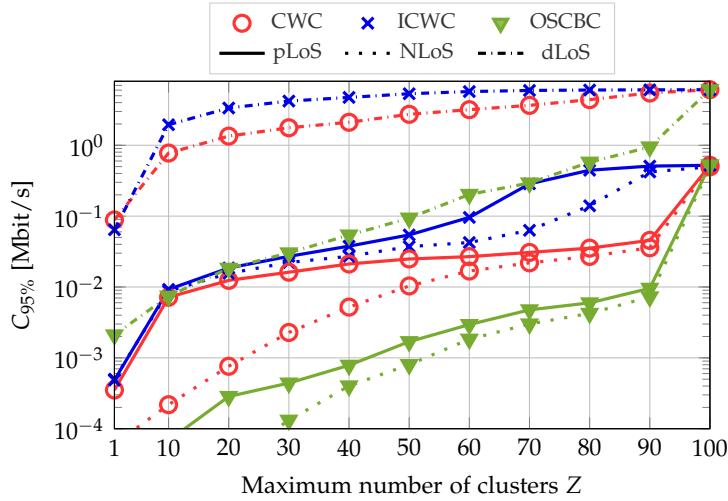


Figure 4.5: 95% percentile of the user capacity as a function of the maximum number of clusters Z , for $N_I = 3200$, unquantized phase shifts, and for different channel conditions.

capacity with a relatively lower number of clusters than in the pLoS and NLoS scenarios.

As expected, OSCBC performs worse than its competitors, and the gap is even more significant in the dLoS scenario (around -30% in terms of sum capacity). The bad performance of OSCBC compared to CWC and ICWC is confirmed also in terms of fairness, as illustrated in Fig. 4.5 (see, in particular, the zoom for $50 \leq Z \leq 90$).

Finally, even though ICWC is not explicitly designed to maximize the sum capacity, it shows similar performance (if not even slightly better) as CWC in the dLoS scenario. The rationale behind this behavior is not clear and deserves more investigation. Most likely, it is related to the fact that, in the dLoS scenario, all UEs have similar channel conditions, which permits ICWC to choose, on average, a good IRS configuration even among the worst UEs in the clusters.

Impact of the IRS configuration. Figs. 4.6 and 4.7 show the impact of the number of IRS radiating elements on the system performance when considering the CWC and ICWC clustering algorithms. As expected, both fairness (measured in terms of the 95% percentile of the average sum capacity) and sum capacity increase as the IRS is larger and operates with more reflecting elements, regardless of the number of clusters. For example, we observe that

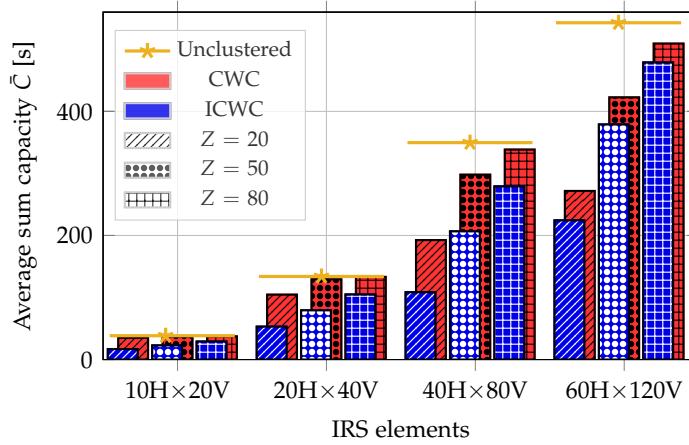


Figure 4.6: Average sum capacity for CWC and ICWC as a function of the number of reflecting elements at the IRS, for unquantized phase shifts, and considering a pLoS channel for the IRS-UEs links.

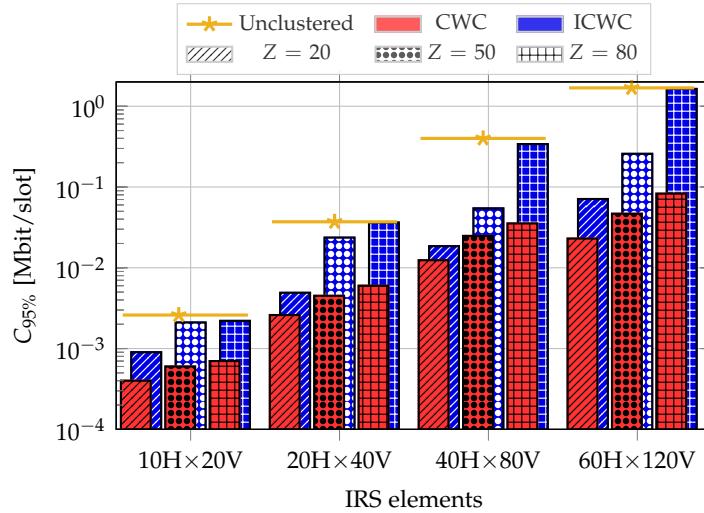


Figure 4.7: 95% percentile of the user capacity for CWC and ICWC as a function of the number of reflecting elements at the IRS, for unquantized phase shifts, and considering a pLoS channel for the IRS-UEs links.

CWC is able to approach the optimal sum capacity with as few as 20 clusters for small-sized IRS, i.e., with $10H \times 20V$ or $20H \times 40V$ arrays. The same trends are shown also in Fig. 4.7 in terms of fairness. Still, notice that \bar{C} is below 100 Mbps, which is not compatible with the requirement of most 5G applications when the IRS is made of fewer than 200 elements, which justifies the use of larger IRS panels [165].

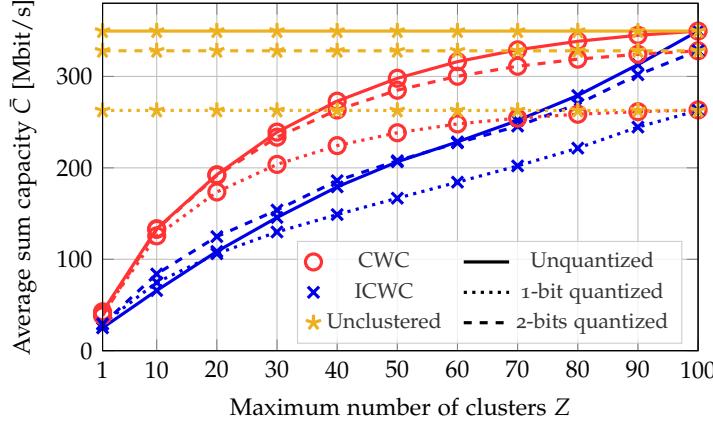


Figure 4.8: Average sum capacity as a function of the maximum number of clusters Z , for $N_I = 3200$ and for different degrees of quantization of the phase shifts, and considering a pLoS channel for the IRS-UEs links.

Nevertheless, we still observe that the number of reflecting elements has an impact on the number of clusters that are needed to provide maximum performance. Indeed, the number of possible IRS configurations increases as we consider larger IRS antennas. In turn, this decreases the likelihood of UEs having the same (or similar) ideal configurations, and therefore, it increases the probability of being associated with increasingly suboptimal centroids if the number of clusters is small. However, if the number of phase shifters is large, the suboptimality is mitigated by the increasing number of reconfigurations. Typically, the IRS reconfiguration cost is proportional to the number of radiating elements, and therefore the specific reconfiguration cost may be different for different IRSs in general. A detailed quantitative analysis of this issue would need to go into the specifics of the various IRS architectures, which goes beyond the scope of the present work, and will be considered in our future work.

Impact of quantization. Figs. 4.8 and 4.9 display the average sum capacity and the 95% percentile, respectively, as a function of the maximum number of clusters Z for CWC and ICWC, and of the number of quantization bits b of the phase shifts. Notice that energy and hardware constraints pose a limit to b [201], which implies restricting the infinite set of possible IRS configurations to a finite set of cardinality 2^{bN_I} . Moreover, the quantization constraint affects the beamforming capabilities of the IRS [164], with negative implications for the resulting achievable sum capacity. In [186], results were obtained consid-

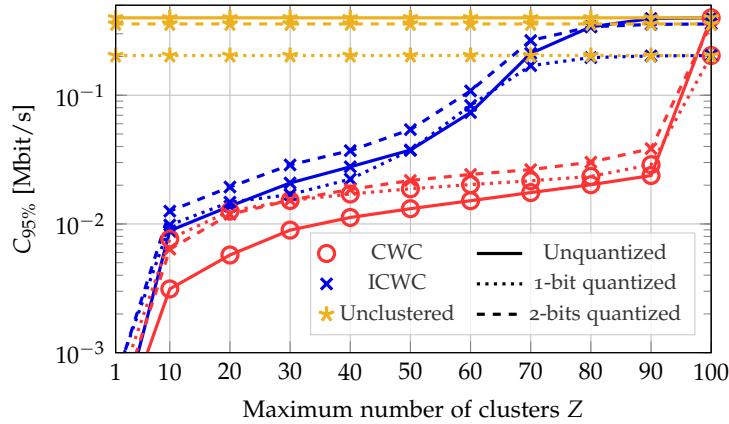


Figure 4.9: 95% percentile of the user capacity as a function of the maximum number of clusters Z , for $N_I = 3200$ and for different degrees of quantization of the phase shifts, and considering a pLoS channel for the IRS-UEs links.

ering that the quantization was performed only at the end of the clustering procedure. Here, instead, we assume that the quantization of the phase shifts is taken into account from the initial optimization stage. The results reveal that the use of non-ideal phase shifters leads to a 30% degradation in the sum capacity when using $b = 1$ at the IRS, while the performance is close to the unquantized baseline if more quantization bits are used. Furthermore, it is shown that the gap between quantized and the unquantized performance increases with Z . As a result, 1-bit quantization is sufficient to guarantee a performance comparable to the unquantized case with a small number of clusters, while more quantization bits are needed to achieve higher capacity. In any case, we can conclude that our proposed capacity-based clustering algorithms are robust to phase-shift quantization.

Scalability. Finally, we prove the scalability performance of the proposed clustering algorithms. To do so, we first show the performance of capacity-based clustering as a function of the number of UEs in Figs 4.10 and 4.11. In particular, we compare CWC and ICWC with HC as a function of the ratio K/N , for an unquantized $40H \times 80V$ IRS, $K = \{50, 100, 150\}$ UEs, and considering a pLoS channel for the IRS-UEs links. The results are in line with the plots in Figs 4.2 and 4.3, which demonstrates the scalability of the proposed clustering techniques for different numbers of UEs. Finally, Fig. 4.12 depicts the average minimum number of IRS configurations Z_{min} needed to achieve 80% of the maximum achievable sum capacity (“unclustered” base-

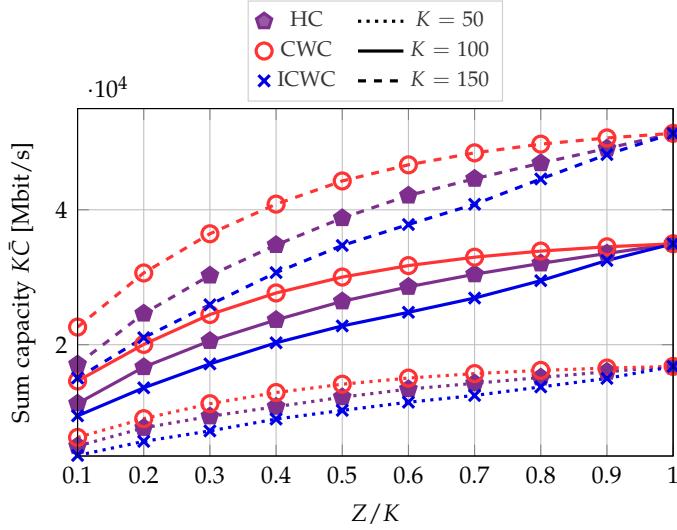


Figure 4.10: Sum capacity as a function of the maximum number of clusters over the number of UEs Z/K , for different values of K , for an unquantized $40H \times 80V$ IRS, $K = \{50, 100, 150\}$, and considering a pLoS channel for the IRS-UEs links. For readability, the results are shown without averaging over the TDMA frame length.

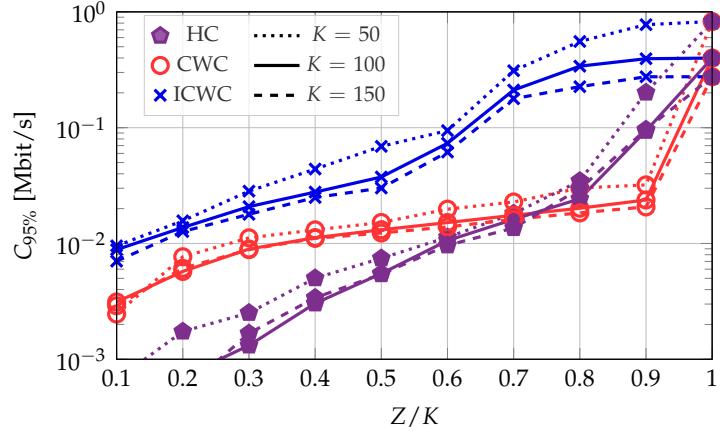


Figure 4.11: 95% of the user capacity as a function of the maximum number of clusters over the number of UEs Z/K , for different values of K , for an unquantized $40H \times 80V$ IRS, $K = \{50, 100, 150\}$ UEs, and considering a pLoS channel for the IRS-UEs links.

line) as a function of the number of UEs K in the system. Notably, we observe that CWC and OSCBC are confirmed to be the best algorithms to optimize the sum capacity, even for a limited number of IRS configurations. For example, both solutions achieve 80% of the maximum sum capacity with less

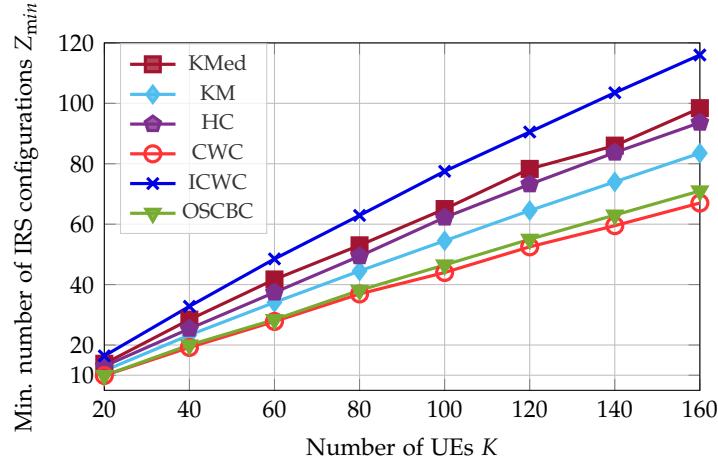


Figure 4.12: Minimum number of IRS configurations (clusters) Z_{min} to achieve 80% of the maximum achievable sum capacity, for an unquantized $40H \times 80V$ IRS, and considering a pLoS channel for the IRS-UEs links.

than half the number of configurations than in the “unclustered” deployment. Moreover, we recognize the same trends as in the previous results. Specifically, capacity-based clustering outperforms distance-based clustering and requires fewer IRS reconfigurations to maximize the sum capacity (up to –37% considering CWC vs. KMed). Furthermore, the gap increases as the number of UEs increases.

4.9 Conclusions

We considered a MIMO cellular network, in which a gNB serving multiple UEs is assisted by an IRS acting as a relay. Notably, we considered practical constraints on the IRS reconfiguration period. We studied a TDMA scheduling for downlink transmissions, and formulated an optimization problem to maximize the average sum capacity, subject to a fixed number of IRS reconfigurations per time frame. We first discussed an iterative algorithm to obtain the optimal IRS configuration of each UE. Then, we proposed clustering-based scheduling algorithms, which group UEs with similar (ideal) IRS configurations based on either a distance metric or the achievable capacity, to mitigate the performance degradation due to the constraint in the number of possible reconfigurations. Different clustering algorithms were numerically evaluated in terms of computational complexity, sum capacity, and

fairness under different channel conditions, as a function of the size of the IRS size and the number of users, and with or without quantization of phase shifts. The results showed that capacity-based clustering outperforms distance-based clustering, and can achieve up to 85% of the sum capacity obtained in an ideal deployment (with no reconfiguration constraints), reducing by 50% the number of IRS reconfigurations.

Appendix

Proof of Theorem 4

Proof for KM (Lloyd): In the assignment step, each UE k is assigned to the cluster z that minimizes the squared distance $\delta(\Phi_k^*, \Phi^{(z)})^2$. This guarantees that the total sum $J(\mathcal{U}_1, \dots, \mathcal{U}_Z, \mathcal{I})$ does not increase. Then, in the update step, $\Phi^{(z)}$ is recalculated as the average Φ_k^* within each cluster, so as to minimize the intra-cluster sum of squared distances $\sum_{k \in \mathcal{U}_z} \delta(\Phi_k^*, \Phi^{(z)})^2$, for all z . Therefore, the conditions of [202, Lemma 5] are satisfied, which ensures the convergence to a local minimum. Notice that [202] does not specify the number of iterations needed to reach convergence, which could be large in the case of highly dimensional spaces. Therefore, in practice, we limit the maximum number of iterations to I_{\max}^{KM} . ■

Proof for Agglomerative HC: At each step, clusters are merged to minimize the increase of the total intra-cluster sum of squared distances. This is equivalent to choosing the merged cluster that results in the smallest increase of $J(\mathcal{U}_1, \dots, \mathcal{U}_Z, \mathcal{I})$. Then, as in the update step of KM, the average of the data points minimizes $\sum_{k \in \mathcal{U}_z} \delta(\Phi_k^*, \Phi^{(z)})^2$, for all z . Once the number of clusters Z reaches the desired value, convergence to a local minimum is reached. ■

Proof for KMed (PAM): Since, at each iteration, a swap is performed only when it leads to a lower value of the intra-cluster sum of squares, $J(\mathcal{U}_1, \dots, \mathcal{U}_Z, \mathcal{I})$ does not increase over different iterations. Given the finite number of data points and possible configurations, the algorithm is guaranteed to converge to a configuration where no swap can further decrease the objective function, thus reaching a local minimum. ■

5 Conclusions

One of the key innovations of 5G is its ability to harness mmWave frequencies, thereby unlocking substantial unused radio resources that were previously inaccessible. The next iteration of mobile networks, i.e., 6G is predicted to extend the bandwidth capabilities even further by expanding the supported spectrum bands to encompass THz frequencies. While these advancements enable unprecedented data rates and ultra-low latency, mmWave and THz frequencies are hindered by challenging propagation conditions that impede the provision of ubiquitous high-speed wireless connectivity. This thesis tackles these limitations by studying innovative coverage enhancement solutions which hold promise for the widespread adoption of mmWave and THz deployments in 6G cellular networks, overcoming their unfavorable propagation characteristics.

Specifically, in Chapter 2 we proposed a semi-centralized resource partitioning scheme for 5G and beyond IAB networks, coupled with a set of allocation policies. Our analysis demonstrated that the integration of this lightweight resource allocation cooperation significantly enhances the system end-to-end throughput and delay, while preventing or mitigating network congestion in the backhaul links. We also provided considerations for implementing a semi-centralized resource allocation controller in real-world deployments. Moreover, we introduced the first reliability-focused scheduling and path selection algorithm specifically designed for IAB networks. Our RL-based solution was shown to effectively account for the complexities of the network, including channel variations, interference, and load dynamics. Results demonstrated that our proposed algorithm not only achieves highly reliable performance in the presence of these challenges, but also outperforms benchmark schemes in terms of throughput, latency, and packet-drop rate. The reliability of our approach is rooted in its ability to jointly minimize the average latency and the expected value of tail losses through the use of CVaR as a risk metric. Additionally, we conducted the first comprehensive

5 Conclusions

evaluation of the potential benefits of sub-terahertz frequencies for 6G IAB networks, utilizing a customized extension of the open-source Sionna simulator. This simulation framework enabled us to assess the feasibility of employing mixed mmWave and sub-terahertz links in conjunction with greedy algorithms to optimize the deployment of backhaul networks.

In Chapter 3, we proposed a signal model for 5G NR deployments featuring IRSs and AF relays based on the 3GPP TR 38.901 channel model. Our simulation framework provides numerical guidelines for dimensioning IRS/AF-assisted networks. We also developed an ns-3 implementation of the 3GPP channel model for NTN scenarios, which we open-sourced and validated against 3GPP calibration results. Additionally, we presented optimizations for simulating MIMO wireless channels in ns-3, including improved linear algebra routines and a performance-oriented statistical channel model that significantly reduces simulation time.

In Chapter 4, we analyzed the performance of IRS-aided cellular deployments with practical constraints on the IRS reconfiguration period. In this context, we maximized the average sum capacity, subject to a fixed number of IRS reconfigurations per time frame. We proposed a clustering-based scheduling algorithms, which groups users with similar (optimal) IRS configurations based on either a distance metric or the achievable capacity, thereby mitigating the impact of the reconfiguration limit. The efficacy of the proposed algorithms was assessed in terms of computational complexity, sum capacity, and fairness across diverse channel conditions, as a function of the IRS configuration and the number of users. Our results indicate that capacity-based clustering outperforms distance-based approaches, yielding up to 85% of the sum capacity achievable in an ideal scenario (i.e., with no reconfiguration constraints), while reducing the number of IRS reconfigurations by up to 50%.

Acronyms

3GPP	3rd Generation Partnership Project
5G	5th generation
5GC	5G Core
AF	Amplify-and-Forward
AI	Artificial Intelligence
AM	Acknowledged Mode
BA	Backlog Avoidance
BAP	Backhaul Adaptation Protocol
BBU	Baseband Unit
BS	Base Station
BSR	Buffer Status Report
BWP	Bandwidth Part
C-RAN	Cloud Radio Access Network
CDF	Cumulative Distribution Function
CE	Control Element
CIR	Channel Impulse Response
CMOS	Complementary Metal-Oxide Semiconductor
CN	Core Network
CNR	Carrier-to-Noise Ratio
CQI	Channel Quality Information
CSI	Channel State Information

Acronyms

CVaR	Conditional Value at Risk
CWC	Capacity-Weighted Clustering
DAG	Directed Acyclic Graph
DCI	Downlink Control Information
DL	Downlink
dLoS	deterministic LoS
DRL	Deep Reinforcement Learning
DU	Distributed Unit
E2E	End-to-End
ECDF	Empirical Cumulative Distribution Function
ECEF	Earth-Centered Earth-Fixed
EESM	Exponential Effective SNR Mapping
EM	Electromagnetic
eMBB	enhanced Mobile Broadband
FDD	Frequency Division Duplexing
FDMA	Frequency Division Multiple Access
FPGA	Field Programmable Gate Array
FSPL	Free Space Path Loss
FTR	Fluctuating Two-Ray
GEO	Geostationary Equatorial Orbit
gNB	Next Generation Node Base
HAP	High Altitude Platform
HC	Hierarchical Clustering
HetNet	Heterogeneous Network
HPBW	Half Power Beamwidth
IA	Initial Access
IAB	Integrated Access and Backhaul

Acronyms

ICWC	Inverse Capacity-Weighted Clustering
ILP	Integer Linear Program
IRS	Intelligent Reflective Surface
ITU	International Telecommunications Union
KM	K-Means
KMed	K-Medoids
KPI	Key Performance Indicator
L ₂ SM	Link-to-System Mapping
LCG	Logical Channel Group
LEO	Low Earth Orbit
LoS	Line-of-Sight
LP	Linear Programming
LSFC	Large-Scale Fading Coefficient
LTE	Long Term Evolution
m-MIMO	massive MIMO
MAC	Medium Access Control
MBS	Macro Base Station
MC	Markov Chain
MDP	Markov Decision Process
MEC	Mobile Edge Cloud
MIMO	Multiple Input, Multiple Output
MINLP	Mixed Integer Nonlinear Programming
MLR	Maximum-local-rate
mmWave	millimeter wave
MNO	Mobile Network Operator
MRBA	Max-Rate Backlog Avoidance
MSR	Max Sum-Rate

Acronyms

MU	Multi-User
MWM	Maximum Weighted Matching
NLoS	Non-Line-of-Sight
NOMA	Non-Orthogonal Multiple Access
NR	New Radio
NTN	Non-Terrestrial Networks
NUM	Network Utility Maximization
OFDM	Orthogonal Frequency Division Multiplexing
OFDMA	Orthogonal Frequency-Division Multiple Access
OMA	Orthogonal Multiple Access
OSCBC	One-Shot Capacity-Based Clustering
PAM	Partition Around Medoids
PDCP	Packet-Data Convergence Protocol
PDF	Probability Density Function
PER	Packet Error Rate
PHY	Physical
PL	Path Loss
pLoS	probabilistic LoS
PSD	Power Spectrum Density
QD	Quasi Deterministic
QoS	Quality of Service
RAN	Radio Access Network
RAT	Radio Access Technology
RF	Radio Frequency
RIC	RAN Intelligent Controller
RL	Reinforcement Learning
RLC	Radio Link Control

RLF	Radio Link Failure
RNTI	Radio Network Temporary Identifier
RRH	Remote Radio Head
RSMA	Rate-Splitting Multiple Access
RT	Ray Tracer
SBS	Small Base Station
SCAROS	Scalable and Robust Self-backhauling Solution
SCM	Spatial Channel Model
SDMA	Spatial Division Multiple Access
SF	Shadow Fading
SI	Study Item
SIMD	Single Instruction, Multiple Data
SINR	Signal-to-Interference-plus-Noise Ratio
SISO	Single Input Single Output
SNR	Signal-to-Noise Ratio
ST	Spanning Tree
STL	Standard Template Library
SVD	Singular Value Decomposition
TB	Transport Block
TCP	Transmission Control Protocol
TDD	Time Division Duplexing
TDMA	Time Division Multiple Access
THz	terahertz
TTI	Transmission Time Interval
UAS	Unmanned Aerial System
UAV	Unmanned Aerial Vehicles
UDP	User Datagram Protocol

Acronyms

UE	User Equipment
UL	Uplink
ULA	Uniform Linear Array
UMa	Urban Macro
UMi	Urban Microcell
UPA	Uniform Planar Array
URLLC	Ultra-Reliable Low-Latency Communication
VR	Virtual Reality
VSAT	Very Small Aperture Terminal
WB	Wideband
WLAN	Wireless Local Area Network
WPANs	Wireless Personal Area Networks
XR	Extended Reality
ZF	Zero-Forcing

Publications

International journals

- [J1] **M. Pagin**, T. Zugno, M. Polese and M. Zorzi, "Resource Management for 5G NR Integrated Access and Backhaul: A Semi-Centralized Approach", *IEEE Transactions on Wireless Communications*, vol. 21, no. 2, pp. 753-767, Feb. 2022.
- [J2] *Early access, to appear in IEEE/ACM Transactions on Networking*. A. A. Gargari, A. Ortiz, **M. Pagin**, W. de Sombre, M. Zorzi and A. Asadi, "Risk-Averse Learning for Reliable mmWave Self-Backhauling", 2024.
- [J3] *To appear in IEEE Transactions on Communications*. A. Rech, L. Badia, S. Tomasin, **M. Pagin**, M. Giordani, J. Gambini and M. Zorzi, "Downlink Clustering-Based Scheduling of IRS-Assisted Communications With Re-configuration Constraints", 2024.
- [J4] *Under review in IEEE Transactions on Wireless Communications*. M. Rawat, **M. Pagin**, M. Giordani, L.-A. Dufrene, Q. Lampin and M. Zorzi, "Optimizing Energy Efficiency of 5G RedCap Beam Management for Smart Agriculture Applications", 2024.
- [J5] *To be submitted to IEEE Transactions on Communications*. **M. Pagin**, L. Badia and M. Zorzi, "Interaction of Sources With Full Online Information as a Markov Game of Age of Information".
- [J6] *To be submitted to IEEE/ACM Transactions on Networking*. **M. Pagin**, A. Traspadini, M. Giordani and M. Zorzi, "End-to-end Simulation of 5G NR Integrated Access and Backhaul Networks for Remote Connectivity".

Acronyms

Conference proceedings

- [C1] **M. Pagin**, S. Lagén, B. Bojovic, M. Polese and M. Zorzi, "Improving the Efficiency of MIMO Simulations in ns-3", Proceedings of the 2023 Workshop on ns-3 (WNS3 23), Association for Computing Machinery, New York, NY, USA, 2023.
- [C2] **M. Pagin** et al. "End-to-End Simulation of 5G Networks Assisted by IRS and AF Relays", 20th Mediterranean Communication and Computer Networking Conference (MedComNet), 2022.
- [C3] A. A. Gargari, A. Ortiz, **M. Pagin**, A. Klein, M. Hollick, M. Zorzi and A. Asadi, "Safehaul: Risk-Averse Learning for Reliable mmWave Self-Backhauling in 6G Networks", IEEE INFOCOM, 2023.
- [C4] H. Poddar, T. Yoshimura, **M. Pagin**, T. Rappaport, A. Ishi and M. Zorzi, "ns-3 Implementation of Sub-Terahertz and Millimeter Wave Drop-based NYU Channel Model (NYUSIM)", Proceedings of the 2023 Workshop on ns-3 (WNS3 23), Association for Computing Machinery, New York, NY, USA, 2023.
- [C5] A. A. Gargari, **M. Pagin**, A. Ortiz, N. M. Gholian, M. Polese and M. Zorzi, "Demo:SeBaSi system-level Integrated Access and Backhaul simulator for self-backhauling", 2023 IEEE 24th International Symposium on a World of Wireless, Mobile and Multimedia Networks (WoWMoM), Boston, MA, USA, 2023.
- [C6] M. Sandri, **M. Pagin**, M. Giordani and M. Zorzi, "Implementation of a Channel Model for Non-Terrestrial Networks in ns-3", Proceedings of the 2023 Workshop on ns-3, 2023.
- [C7] **M. Pagin**, T. Zugno, M. Giordani, L.-A. Dufrene, Q. Lampin and M. Zorzi, "5G NR-Light at Millimeter Waves: Design Guidelines for Mid-Market IoT Use Cases," 2023 International Conference on Computing, Networking and Communications (ICNC), Honolulu, HI, USA, 2023.
- [C8] A. Gargari, **M. Pagin**, M. Polese and M. Zorzi, "6G Integrated Access and Backhaul Networks with Sub-Terahertz Links," 2023 18th Wireless On-Demand Network Systems and Services Conference (WONS), Madonna di Campiglio, Italy, 2023.

- [C9] A. Rech, **M. Pagin**, S. Tomasin, F. Moretto, L. Badia, M. Giordani, J. Gambini and M. Zorzi, "Downlink TDMA Scheduling for IRS-aided Communications with Block-Static Constraints", 2023 IEEE Wireless Communications and Networking Conference (WCNC), 2023.
- [C10] **M. Pagin**, L. Badia and M. Zorzi, "A Markov Game of Age of Information From Strategic Sources With Full Online Information", 2023 IEEE International Conference on Communications (ICC), 2023.
- [C11] H. Poddar, T. Yoshimura, **M. Pagin**, T.T. Rappaport, A. Ishii and M. Zorzi, "Full Stack End-To-End mmWave Simulations Using 3GPP and NYUSIM Channel Model in ns-3", IEEE International Conference on Communications (ICC), 2023.

Acronyms

Bibliography

- [1] ITU-R. *Future Technology Trends of Terrestrial International Mobile Telecommunications Systems Towards 2030 and Beyond*. Report M.2516-o. Nov. 2022.
- [2] ITU-R. *IMT Vision - Framework and Overall Objectives of the Future Development of IMT for 2020 and Beyond*. Recommendation ITU-R M.2083. Sept. 2015.
- [3] A. Ghosh, A. Maeder, M. Baker, and D. Chandramouli. “5G Evolution: A View on 5G Cellular Technology Beyond 3GPP Release 15”. In: *IEEE Access* 7 (Sept. 2019), pp. 127639–127651.
- [4] M. Polese, M. Giordani, M. Mezzavilla, S. Rangan, and M. Zorzi. “6G Enabling Technologies”. In: *6G mobile wireless networks*. Springer, 2021, pp. 25–41.
- [5] J. Hoydis, F. A. Aoudia, A. Valcarce, and H. Viswanathan. “Toward a 6G AI-Native Air Interface”. In: *IEEE Communications Magazine* 59.5 (May 2021), pp. 76–81.
- [6] C. Han and Y. Chen. “Propagation Modeling for Wireless Communications in the Terahertz Band”. In: *IEEE Communications Magazine* 56.6 (June 2018), pp. 96–101. ISSN: 1558-1896. DOI: [10.1109/MCOM.2018.1700898](https://doi.org/10.1109/MCOM.2018.1700898).
- [7] J. M. Jornet and I. F. Akyildiz. “Channel modeling and capacity analysis for electromagnetic wireless nanonetworks in the terahertz band”. In: *IEEE Trans. Wireless Commun.* 10.10 (Aug. 2011), pp. 3211–3221.
- [8] M. Polese, J. M. Jornet, T. Melodia, and M. Zorzi. “Toward End-to-End, Full-Stack 6G Terahertz Networks”. In: *IEEE Commun. Mag.* 58.11 (Nov. 2020), pp. 48–54. ISSN: 1558-1896. DOI: [10.1109/MCOM.001.2000224](https://doi.org/10.1109/MCOM.001.2000224).

Bibliography

- [9] M. Shafi, J. Zhang, H. Tataria, A. F. Molisch, S. Sun, T. S. Rappaport, F. Tufvesson, S. Wu, and K. Kitao. "Microwave vs. Millimeter-Wave Propagation Channels: Key Differences and Impact on 5G Cellular Systems". In: *IEEE Communications Magazine* 56.12 (Dec. 2018), pp. 14–20.
- [10] Y. Heng, J. G. Andrews, J. Mo, V. Va, A. Ali, B. L. Ng, and J. C. Zhang. "Six Key Challenges for Beam Management in 5.5G and 6G Systems". In: *IEEE Communications Magazine* 59.7 (July 2021), pp. 74–79.
- [11] M. Zhang, M. Polese, M. Mezzavilla, J. Zhu, S. Rangan, S. Panwar, and M. Zorzi. "Will TCP work in mmWave 5G cellular networks?" In: *IEEE Communications Magazine* 57.1 (2019), pp. 65–71.
- [12] M. Polese, L. Bonati, S. D’Oro, S. Basagni, and T. Melodia. "Understanding O-RAN: Architecture, Interfaces, Algorithms, Security, and Research Challenges". In: *IEEE Communications Surveys & Tutorials* (Jan. 2023), Early Access.
- [13] M. Polese, L. Bonati, S. D’Oro, S. Basagni, and T. Melodia. "ColO-RAN: Developing Machine Learning-Based xApps for Open RAN Closed-Loop Control on Programmable Experimental Platforms". In: *IEEE Transactions on Mobile Computing* (July 2022), pp. 1–14.
- [14] A. Lacava, M. Polese, R. Sivaraj, R. Soundrarajan, B. S. Bhati, T. Singh, T. Zugno, F. Cuomo, and T. Melodia. "Programmable and Customized Intelligence for Traffic Steering in 5G Networks using Open RAN Architectures". In: *arXiv preprint arXiv:2209.14171* (2022).
- [15] A. Ashtari Gargari, A. Ortiz, M. Pagin, A. Klein, M. Hollick, M. Zorzi, and A. Asadi. "Safehaul: Risk-Averse Learning for Reliable mmWave Self-Backhauling in 6G Networks". In: *IEEE Conference on Computer Communications (INFOCOM)*. New York, USA, 2023.
- [16] J. Hoydis, F. A. Aoudia, S. Cammerer, M. Nimier-David, N. Binder, G. Marcus, and A. Keller. "Sionna RT: Differentiable Ray Tracing for Radio Propagation Modeling". In: *2023 IEEE Globecom Workshops (GC Wkshps)*. 2023, pp. 317–321. DOI: [10.1109/GCWkshps58843.2023.10465179](https://doi.org/10.1109/GCWkshps58843.2023.10465179).

Bibliography

- [17] M. Mezzavilla, M. Zhang, M. Polese, R. Ford, S. Dutta, S. Rangan, and M. Zorzi. "End-to-End Simulation of 5G mmWave Networks". In: *IEEE Communications Surveys & Tutorials* 20.3 (Third quarter 2018), pp. 2237–2263. ISSN: 1553-877X. DOI: [10.1109/COMST.2018.2828880](https://doi.org/10.1109/COMST.2018.2828880).
- [18] N. Patriciello, S. Lagen, B. Bojovic, and L. Giupponi. "An E2E Simulator for 5G NR Networks". In: *Simulation Modelling Practice and Theory* 96 (Nov. 2019), p. 101933.
- [19] D. Magrin, S. Avallone, S. Roy, and M. Zorzi. "Validation of the ns-3 802.11ax OFDMA implementation". In: *Proceedings of the 2021 Workshop on Ns-3*. WNS3 '21. Virtual Event, USA: Association for Computing Machinery, 2021, pp. 1–8. ISBN: 9781450390347. DOI: [10.1145/3460797.3460798](https://doi.org/10.1145/3460797.3460798). URL: <https://doi.org/10.1145/3460797.3460798>.
- [20] H. Assasa, N. Grosheva, T. Ropitault, S. Blandino, N. Golmie, and J. Widmer. "Implementation and evaluation of a WLAN IEEE 802.11ay model in network simulator ns-3". In: *Proceedings of the 2021 Workshop on Ns-3*. WNS3 '21. Virtual Event, USA: Association for Computing Machinery, 2021, pp. 9–16. ISBN: 9781450390347. DOI: [10.1145/3460797.3460799](https://doi.org/10.1145/3460797.3460799). URL: <https://doi.org/10.1145/3460797.3460799>.
- [21] A. Jönsson, D. Åkerman, E. Fitzgerald, C. Nyberg, B. E. Priyanto, and K. Agardh. "Modeling, implementation and evaluation of IEEE 802.11ac in NS-3 for enterprise networks". In: *2016 Wireless Days (WD)*. 2016, pp. 1–6. DOI: [10.1109/WD.2016.7461452](https://doi.org/10.1109/WD.2016.7461452).
- [22] T. Zugno, M. Polese, N. Patriciello, B. Bojović, S. Lagen, and M. Zorzi. "Implementation of a Spatial Channel Model for Ns-3". In: *Proceedings of the 2020 Workshop on ns-3*. Gaithersburg, MD, USA, 2020.
- [23] 3GPP. *Study on channel model for frequencies from 0.5 to 100 GHz*. Technical Report (TR) 38.901. Version 16.1.0. 3rd Generation Partnership Project (3GPP), Jan. 2020.
- [24] P. Testolina, M. Lecci, M. Polese, M. Giordani, and M. Zorzi. "Scalable and Accurate Modeling of the Millimeter Wave Channel". In: *International Conference on Computing, Networking and Communications (ICNC)*. IEEE. Big Island, Hawaii, USA, 2020.

Bibliography

- [25] S. Lagen, K. Wanuga, H. Elkotby, S. Goyal, N. Patriciello, and L. Giupponi. "New Radio Physical Layer Abstraction for System-Level Simulations of 5G Networks". In: *International Conference on Communications (ICC)*. Virtual Event: IEEE, 2020.
- [26] 3GPP. "Study on New Radio (NR) to support non-terrestrial networks". In: *TR 38.811* (2020).
- [27] M. Lecci, P. Testolina, M. Polese, M. Giordani, and M. Zorzi. "Accuracy vs. Complexity for mmWave Ray-Tracing: A Full Stack Perspective". In: *IEEE Trans. Wireless Commun.* (June 2021).
- [28] T. Zugno, M. Polese, N. Patriciello, B. Bojović, S. Lagén, and M. Zorzi. "Implementation of a Spatial Channel Model for ns-3". In: *Proceedings of the 2020 Workshop on ns-3*. Gaithersburg, MD, USA: ACM, June 2020.
- [29] 3GPP. *Study on Channel Model for Frequencies from 0.5 to 100 GHz*. TR 38.901 (Rel. 15), V15.0.0. 2019.
- [30] G. Guennebaud, B. Jacob, et al. *Eigen v3*. <http://eigen.tuxfamily.org>. 2010.
- [31] R. Flamini, D. De Donno, J. Gambini, F. Giuppi, C. Mazzucco, A. Milani, and L. Resteghini. "Towards a Heterogeneous Smart Electromagnetic Environment for Millimeter-Wave Communications: An Industrial Viewpoint". In: *IEEE Trans. Antennas Propag.* 70.10 (Oct. 2022), pp. 8898–8910.
- [32] E. Björnson, Ö. Özdogan, and E. G. Larsson. "Intelligent reflecting surface versus decode-and-forward: How large surfaces are needed to beat relaying?" In: *IEEE Commun. Lett.* 9.2 (Feb. 2020), pp. 244–248.
- [33] C. Huang, A. Zappone, G. C. Alexandropoulos, M. Debbah, and C. Yuen. "Reconfigurable intelligent surfaces for energy efficiency in wireless communication". In: *IEEE Trans. Wireless Commun.* 18.8 (June 2019), pp. 4157–4170.
- [34] Q. Wu and R. Zhang. "Intelligent reflecting surface enhanced wireless network: Joint active and passive beamforming design". In: *IEEE Global Communications Conference (GLOBECOM)*. IEEE. 2018.

- [35] E. Basar and I. Yildirim. "SimRIS Channel Simulator for Reconfigurable Intelligent Surface-Empowered Communication Systems". In: *IEEE Latin-American Conference on Communications (LATINCOM)*. 2020.
DOI: [10.1109/LATINCOM50620.2020.9282349](https://doi.org/10.1109/LATINCOM50620.2020.9282349).
- [36] P. K. Gkonis, P. T. Trakadas, and D. I. Kaklamani. "A Comprehensive Study on Simulation Techniques for 5G Networks: State of the Art Results, Analysis, and Future Challenges". In: *Electronics* 9.3 (Mar. 2020), p. 468.
- [37] L. Bonati, M. Polese, S. D'Oro, S. Basagni, and T. Melodia. "Open, Programmable, and Virtualized 5G Networks: State-of-the-Art and the Road Ahead". In: *Computer Networks (COMNET)* 182 (Aug. 2020).
- [38] F. Wilhelmi, M. Carrascosa, C. Cano, A. Jonsson, V. Ram, and B. Belalata. "Usage of network simulators in machine-learning-assisted 5G/6G networks". In: *IEEE Wireless Communications* 28.1 (Feb. 2021), pp. 160–166.
- [39] 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: *2019 16th IEEE Annual Consumer Communications & Networking Conference (CCNC)*. IEEE. 2019, pp. 1–6.
- [40] G. Nardini, D. Sabella, G. Stea, P. Thakkar, and A. Virdis. "Simu5G—An OMNeT++ Library for End-to-End Performance Evaluation of 5G Networks". In: *IEEE Access* 8 (2020), pp. 181176–181191.
- [41] S. Pratschner, B. Tahir, L. Marijanovic, M. Mussbah, K. Kirev, R. Nissel, S. Schwarz, and M. Rupp. "Versatile mobile communications simulation: The Vienna 5G link level simulator". In: *EURASIP Journal on Wireless Communications and Networking* 2018.1 (Sept. 2018), pp. 1–17.
- [42] M. K. Müller, F. Ademaj, T. Dittrich, A. Fastenbauer, B. R. Elbal, A. Nabavi, L. Nagel, S. Schwarz, and M. Rupp. "Flexible multi-node simulation of cellular mobile communications: the Vienna 5G System Level Simulator". In: *EURASIP Journal on Wireless Communications and Networking* 2018.1 (Sept. 2018).
- [43] C.-K. Jao, C.-Y. Wang, T.-Y. Yeh, C.-C. Tsai, L.-C. Lo, J.-H. Chen, W.-C. Pao, and W.-H. Sheen. "WiSE: a system-level simulator for 5G mobile networks". In: *IEEE Wireless Communications* 25.2 (Apr. 2018), pp. 4–7.

Bibliography

- [44] M. Drago, T. Zugno, M. Polese, M. Giordani, and M. Zorzi. "MilliCar: An ns-3 Module for MmWave NR V2X Networks". In: *Proceedings of the 2020 Workshop on ns-3*. Gaithersburg, MD, USA, June 2020.
- [45] 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)*. Barcelona, Spain, 2018.
- [46] K. Heimann, B. Sliwa, M. Patchou, and C. Wietfeld. "Modeling and Simulation of Reconfigurable Intelligent Surfaces for Hybrid Aerial and Ground-Based Vehicular Communications". In: *24th International ACM Conference on Modeling, Analysis and Simulation of Wireless and Mobile Systems* (2021).
- [47] 3GPP. *NR - User Equipment (UE) radio transmission and reception; Part 1: Range 1 Standalone - Rel. 16*. TS 38.101-1. 2020.
- [48] Q. Wu and R. Zhang. "Intelligent Reflecting Surface Enhanced Wireless Network via Joint Active and Passive Beamforming". In: *IEEE Trans. Wireless Commun.* 18.11 (Nov. 2019), pp. 5394–5409. doi: [10.1109/TWC.2019.2936025](#).
- [49] M. Giordani, M. Polese, A. Roy, D. Castor, and M. Zorzi. "A Tutorial on Beam Management for 3GPP NR at mmWave Frequencies". In: *IEEE Commun. Surveys Tuts.* 21.1 (Firstquarter 2019), pp. 173–196. doi: [10.1109/COMST.2018.2869411](#).
- [50] 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*. 2012, pp. 55–60.
- [51] E. Björnson, Ö. Özdogan, and E. G. Larsson. "Reconfigurable intelligent surfaces: Three myths and two critical questions". In: *IEEE Communications Magazine* 58.12 (2020), pp. 90–96.
- [52] T. G. Roberts and S. Kaplan. "Space launch to low earth orbit: How much does it cost". In: *Civil and Commercial Space Space Security* (2020).

Bibliography

- [53] A. Traspadini, M. Giordani, G. Giambene, and M. Zorzi. "Real-Time HAP-Assisted Vehicular Edge Computing for Rural Areas". In: *IEEE Wireless Communications Letters* (2023).
- [54] D. Wang, A. Traspadini, M. Giordani, M.-S. Alouini, and M. Zorzi. "On the Performance of Non-Terrestrial Networks to Support the Internet of Things". In: *Asilomar Conference on Signals, Systems, and Computers* (2022).
- [55] M. Giordani and M. Zorzi. "Satellite Communication at Millimeter Waves: a Key Enabler of the 6G Era". In: *IEEE International Conference on Computing, Networking and Communications (ICNC)* (2020).
- [56] ITU. "Attenuation by atmospheric gases and related effects". In: *Recommendation P.676* (2013).
- [57] ITU. "Ionospheric propagation data and prediction methods required for the design of satellite services and systems". In: *Recommendation P.531* (2012).
- [58] T. Zugno, M. Polese, N. Patriciello, B. Bojović, S. Lagen, and M. Zorzi. "Implementation of a Spatial Channel Model for ns-3". In: *Proc. ACM WNS3*. 2020.
- [59] I. A. of Oil and G. Producers. "Geomatics Guidance Notes 7, part 2: Coordinate Conversions and Transformations including Formulas". In: (2021).
- [60] 3GPP. "Simulation assumptions and parameters for FDD HeNB RF requirements". In: *ETSI TR 136 921* (2009).
- [61] G. Calcev and M. Dillon. "Antenna tilt control in CDMA networks". In: *Proceedings of the 2nd annual international workshop on Wireless internet* (2006).
- [62] 3GPP. "Technical Specification Group Radio Access Network. Solutions for NR to support non-terrestrial networks (NTN)". In: *TR 38.821* (2019).
- [63] J. G. Andrews, T. Bai, M. N. Kulkarni, A. Alkhateeb, A. K. Gupta, and R. W. Heath. "Modeling and Analyzing Millimeter Wave Cellular Systems". In: *IEEE Transactions on Communications* 65.1 (Jan. 2017), pp. 403–430.

Bibliography

- [64] M. Lecci, P. Testolina, M. Giordani, M. Polese, T. Ropitault, C. Gentile, N. Varshney, A. Bodi, and M. Zorzi. "Simplified Ray Tracing for the Millimeter Wave Channel: A Performance Evaluation". In: *Information Theory and Applications Workshop (ITA)*. San Diego, CA USA: IEEE, Feb. 2020.
- [65] T. Zugno, M. Drago, S. Lagén, Z. Ali, and M. Zorzi. "Extending the ns-3 Spatial Channel Model for Vehicular Scenarios". In: *Proceedings of the 2021 Workshop on ns-3*. Virtual Event, USA: Association for Computing Machinery, 2021.
- [66] A. Ramos, Y. Estrada, M. Cantero, J. Romero, D. Martin-Sacristán, S. Inca, M. Fuentes, and J. Monserrat. "Implementation and Calibration of the 3GPP Industrial Channel Model for ns-3". In: *Proceedings of the 2022 Workshop on ns-3*. Virtual Event, USA: Association for Computing Machinery, 2022.
- [67] S. Jin, S. Roy, W. Jiang, and T. R. Henderson. "Efficient Abstractions for Implementing TGn Channel and OFDM-MIMO Links in ns-3". In: *Proceedings of the 2020 Workshop on ns-3*. Gaithersburg, MD, USA: Association for Computing Machinery, 2020.
- [68] S. Jin, S. Roy, and T. R. Henderson. "Efficient PHY Layer Abstraction for Fast Simulations in Complex System Environments". In: *IEEE Transactions on Communications* 69.8 (May 2021), pp. 5649–5660.
- [69] S. Jin, S. Roy, and T. R. Henderson. "EESM-log-AR: an Efficient Error Model for OFDM MIMO Systems Over Time-Varying Channels". In: *Proceedings of the 2021 Workshop on ns-3*. Virtual Event, USA: Association for Computing Machinery, 2021.
- [70] Y. Liu, S. K. Crisp, and D. M. Blough. "Performance Study of Statistical and Deterministic Channel Models for mmWave Wi-Fi Networks in ns-3". In: *Proceedings of the 2021 Workshop on ns-3*. Virtual Event, USA: Association for Computing Machinery, 2021.
- [71] Anuraag Bodi, Steve Blandino, Neeraj Varshney, Jiayi Zhang, Tanguy Ropitault, Mattia Lecci, Paolo Testolina, Jian Wang, Chiehping Lai, and Camillo Gentile. *The NIST Q-D Channel Realization Software*. 2021.

Bibliography

- [72] M. Pagan, M. Giordani, A. A. Gargari, A. Rech, F. Moretto, S. Tomasin, J. Gambini, and M. Zorzi. "End-to-End Simulation of 5G Networks Assisted by IRS and AF Relays". In: *Proc. IEEE MedComNet*. Paphos, Cyprus, 2022.
- [73] 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)*. Kalamata, Greece, 2018.
- [74] M. Rebato, L. Resteghini, C. Mazzucco, and M. Zorzi. "Study of Realistic Antenna Patterns in 5G mmWave Cellular Scenarios". In: *Proc. IEEE ICC*. 2018, pp. 1–6.
- [75] H. Asplund, D. Astely, P. von Butovitsch, T. Chapman, M. Frenne, F. Ghasemzadeh, M. Hagström, B. Hogan, G. Jöngren, J. Karlsson, F. Kronestedt, and E. Larsson. "Chapter 4 - Antenna Arrays and Classical Beamforming". In: *Advanced Antenna Systems for 5G Network Deployments*. Academic Press, 2020, pp. 89–132.
- [76] D. Chizhik, G. J. Foschini, and R. A. Valenzuela. "Capacities of Multi-element Transmit and Receive Antennas: Correlations and Keyholes". In: *Electronics Letters* 36.13 (June 2000), p. 1.
- [77] M. N. Kulkarni, E. Visotsky, and J. G. Andrews. "Correction Factor for Analysis of MIMO Wireless Networks with Highly Directional Beamforming". In: *IEEE Wireless Communications Letters* 7.5 (Mar. 2018), pp. 756–759.
- [78] M. D. Yacoub. "The κ - μ Distribution and the η - μ Distribution". In: *IEEE Antennas and Propagation Magazine* 49.1 (Feb. 2007), pp. 68–81.
- [79] S. L. Cotton. "Human Body Shadowing in Cellular Device-to-Device Communications: Channel Modeling Using the Shadowed κ - μ Fading Model". In: *IEEE J. Sel. Areas Commun.* 33.1 (Nov. 2014), pp. 111–119.
- [80] T. Mavridis, L. Petrillo, J. Sarrazin, A. Benlarbi-Delai, and P. De Doncker. "Near-Body Shadowing Analysis at 60 GHz". In: *IEEE Transactions on Antennas and Propagation* 63.10 (July 2015), pp. 4505–4511.

Bibliography

- [81] J. M. Romero-Jerez, F. J. Lopez-Martinez, J. F. Paris, and A. J. Goldsmith. "The Fluctuating Two-Ray Fading Model: Statistical Characterization and Performance Analysis". In: *IEEE Trans. Wireless Commun.* 16.7 (May 2017), pp. 4420–4432.
- [82] T. W. Anderson and D. A. Darling. "A Test of Goodness of Fit". In: *Journal of the American statistical association* 49.268 (1954), pp. 765–769.
- [83] D. Magrin, D. Zhou, and M. Zorzi. "A Simulation Execution Manager for ns-3: Encouraging Reproducibility and Simplifying Statistical Analysis of ns-3 Simulations". In: *Proceedings of the 22nd International ACM Conference on Modeling, Analysis and Simulation of Wireless and Mobile Systems*. Miami Beach, FL USA, 2019.
- [84] M. Pagin, T. Zugno, M. Polese, and M. Zorzi. "Resource Management for 5G NR Integrated Access and Backhaul: A Semi-Centralized Approach". In: *IEEE Trans. Wireless Commun.* 21.2 (July 2022), pp. 753–767.
- [85] M. Sandri and M. Pagin. *ns-3-NTN TR 38.811 channel model implementation*. <https://gitlab.com/mattiasandri/ns-3-ntn/-/tree/ntn-dev>. 2024.
- [86] X. Xu, Y. Pan, P. P. M. Y. Lwin, and X. Liang. "3D holographic display and its data transmission requirement". In: *International Conference on Information Photonics and Optical Communications*. 2011, pp. 1–4. DOI: [10.1109/IPOC.2011.6122872](https://doi.org/10.1109/IPOC.2011.6122872).
- [87] 3GPP. *NR; NR and NG-RAN Overall description; Stage-2*. Technical Specification (TS) 38.300. v16.2.0. 3rd Generation Partnership Project (3GPP), July 2020.
- [88] F. Z. Yousaf, M. Bredel, S. Schaller, and F. Schneider. "NFV and SDN—Key technology enablers for 5G networks". In: *IEEE J. Sel. Areas Commun.* 35.11 (Nov. 2017), pp. 2468–2478.
- [89] 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". In: *IEEE J. Sel. Areas Commun.* 35.6 (June 2017), pp. 1201–1221. ISSN: 0733-8716.

Bibliography

- [90] W. Saad, M. Bennis, and M. Chen. "A Vision of 6G Wireless Systems: Applications, Trends, Technologies, and Open Research Problems". In: *IEEE Network* 34.3 (May 2020), pp. 134–142.
- [91] F. Khan and Z. Pi. "mmWave mobile broadband (MMB): Unleashing the 3–300 GHz spectrum". In: *34th IEEE Sarnoff Symposium*. Princeton, NJ, USA, 2011.
- [92] S. Rangan, T. S. Rappaport, and E. Erkip. "Millimeter-Wave Cellular Wireless Networks: Potentials and Challenges". In: *Proceedings of the IEEE* 102.3 (Mar. 2014), pp. 366–385.
- [93] I. A. Hemadeh, K. Satyanarayana, M. El-Hajjar, and L. Hanzo. "Millimeter-Wave Communications: Physical Channel Models, Design Considerations, Antenna Constructions, and Link-Budget". In: *IEEE Communications Surveys & Tutorials* 20.2 (Dec. 2017), pp. 870–913.
- [94] F. Gómez-Cuba, E. Erkip, S. Rangan, and F. J. González-Castaño. "Capacity Scaling of Cellular Networks: Impact of Bandwidth, Infrastructure Density and Number of Antennas". In: *IEEE Trans. Wireless Commun.* 17.1 (Nov. 2017), pp. 652–666.
- [95] M. Polese, M. Giordani, T. Zugno, A. Roy, S. Goyal, D. Castor, and M. Zorzi. "Integrated Access and Backhaul in 5G mmWave Networks: Potential and Challenges". In: *IEEE Communications Magazine* 58.3 (Mar. 2020), pp. 62–68.
- [96] C. Madapatha, B. Makki, C. Fang, O. Teyeb, E. Dahlman, M.-S. Alouini, and T. Svensson. "On Integrated Access and Backhaul Networks: Current Status and Potentials". In: *IEEE Open Journal of the Communications Society* 1 (Sept. 2020), pp. 1374–1389. doi: [10.1109/OJCOMS.2020.3022529](https://doi.org/10.1109/OJCOMS.2020.3022529).
- [97] C. Saha and H. S. Dhillon. "Millimeter Wave Integrated Access and Backhaul in 5G: Performance Analysis and Design Insights". In: *IEEE J. Sel. Areas Commun.* 37.12 (Dec. 2019), pp. 2669–2684.
- [98] 3GPP. NR; *Integrated Access and Backhaul (IAB) radio transmission and reception*. Technical Specification (TS) 38.174. v0.1.0. June 2020.

Bibliography

- [99] M. N. Islam, S. Subramanian, and A. Sampath. "Integrated Access Backhaul in Millimeter Wave Networks". In: *IEEE Wireless Communications and Networking Conference (WCNC)*. San Francisco, CA, USA, 2017, pp. 1–6.
- [100] M. N. Islam, N. Abedini, G. Hampel, S. Subramanian, and J. Li. "Investigation of performance in integrated access and backhaul networks". In: *IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS)*. Honolulu, HI, USA, 2018.
- [101] M. N. Kulkarni, A. Ghosh, and J. G. Andrews. "Max-min rates in self-backhauled millimeter wave cellular networks". In: *arXiv preprint arXiv:1805.01040* (2018).
- [102] W. Lei, Y. Ye, and M. Xiao. "Deep Reinforcement Learning-Based Spectrum Allocation in Integrated Access and Backhaul Networks". In: *IEEE Transactions on Cognitive Communications and Networking* 6.3 (Sept. 2020), pp. 970–979. doi: [10.1109/TCCN.2020.2992628](https://doi.org/10.1109/TCCN.2020.2992628).
- [103] M. E. Rasekh, D. Guo, and U. Madhow. "Interference-aware routing and spectrum allocation for millimeter wave backhaul in urban picocells". In: *53rd Annual Allerton Conference on Communication, Control, and Computing*. Monticello, IL, USA, 2015.
- [104] M. Bilal, M. Kang, S. C. Shah, and S.-G. Kang. "Time-Slotted Scheduling Schemes for Multi-hop Concurrent Transmission in WPANs with Directional Antenna". In: *ETRI Journal* 36.3 (June 2014), pp. 374–384.
- [105] R. L. Cruz and A. V. Santhanam. "Optimal routing, link scheduling and power control in multihop wireless networks". In: *22nd Annual Joint Conference of the IEEE Computer and Communications Societies (INFOCOM 2003)*. San Francisco, CA, USA, 2003.
- [106] C. Saha, M. Afshang, and H. S. Dhillon. "Bandwidth partitioning and downlink analysis in millimeter wave integrated access and backhaul for 5G". In: *IEEE Trans. Wireless Commun.* 17.12 (Dec. 2018), pp. 8195–8210.
- [107] R. Singh and P. Kumar. "Throughput optimal decentralized scheduling of multihop networks with end-to-end deadline constraints: Unreliable links". In: *IEEE Transactions on Automatic Control* 64.1 (Jan. 2019), pp. 127–142.

- [108] B. Ji, C. Joo, and N. Shroff. "Throughput-optimal scheduling in multi-hop wireless networks without per-flow information". In: *IEEE/ACM Transactions on Networking* 21.2 (Apr. 2013), pp. 634–647.
- [109] T. K. Vu, C.-F. Liu, M. Bennis, M. Debbah, and M. Latva-Aho. "Path selection and rate allocation in self-backhauled mmWave networks". In: *IEEE Wireless Communications and Networking Conference (WCNC)*. Barcelona, Spain, 2018.
- [110] J. Garcia-Rois, F. Gomez-Cuba, M. R. Akdeniz, F. J. Gonzalez-Castano, J. C. Burguillo, S. Rangan, and B. Lorenzo. "On the analysis of scheduling in dynamic duplex multihop mmWave cellular systems". In: *IEEE Trans. Wireless Commun.* 14.11 (Nov. 2015), pp. 6028–6042.
- [111] F. Gomez-Cuba and M. Zorzi. "Optimal link scheduling in millimeter wave multi-hop networks with space division multiple access". In: *2016 Information Theory and Applications Workshop (ITA)*. La Jolla, CA, USA, 2016.
- [112] F. Gómez-Cuba and M. Zorzi. "Optimal Link Scheduling in Millimeter Wave Multi-hop Networks with MU-MIMO radios." In: *IEEE Trans. Wireless Commun.* 19.3 (Mar. 2020), pp. 1839–1854.
- [113] L. Tassiulas and A. Ephremides. "Stability properties of constrained queueing systems and scheduling policies for maximum throughput in multihop radio networks". In: *29th IEEE Conference on Decision and Control*. Honolulu, HI, USA, 1990.
- [114] F. Kelly. "Charging and rate control for elastic traffic". In: *European Transactions on Telecommunications* 8.1 (Jan. 1997), pp. 33–37.
- [115] F. P. Kelly, A. K. Maulloo, and D. K. Tan. "Rate control for communication networks: shadow prices, proportional fairness and stability". In: *Journal of the Operational Research society* 49.3 (Apr. 1998), pp. 237–252.
- [116] 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)*. Abu Dhabi, United Arab Emirates, 2018.
- [117] 3GPP. *NR; Study on integrated access and backhaul*. Technical Specification (TS) 38.874. v.16.0.0. Jan. 2019.

Bibliography

- [118] S. Dutta, M. Mezzavilla, R. Ford, M. Zhang, S. Rangan, and M. Zorzi. “Frame structure design and analysis for millimeter wave cellular systems”. In: *IEEE Trans. Wireless Commun.* 16.3 (Mar. 2017), pp. 1508–1522.
- [119] B. Korte and J. Vygen. *Combinatorial Optimization*. Springer Berlin Heidelberg, 2002.
- [120] H. N. Gabow. “Data Structures for Weighted Matching and Nearest Common Ancestors with Linking”. In: *Proceedings of the First Annual ACM-SIAM Symposium on Discrete Algorithms*. SODA '90. San Francisco, California, USA: Society for Industrial and Applied Mathematics, 1990, pp. 434–443. ISBN: 0898712513.
- [121] 3GPP. *NR; Medium Access Control (MAC) protocol specification*. Technical Specification (TS) 38.321. v16.1.0. July 2020.
- [122] 3GPP. *CSI feedback for Type I codebook*. Technical Document (TDoc) R1-1713763. Huawei, HiSilicon, Aug. 2017.
- [123] 3GPP. *NR; Backhaul Adaptation Protocol (BAP) specification*. Technical Specification (TS) 38.340. v16.4.0. Mar. 2021.
- [124] 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*. Miami, Florida, USA, 2011.
- [125] C. Pan, H. Zhu, N. J. Gomes, and J. Wang. “Joint precoding and RRH selection for user-centric green MIMO C-RAN”. In: *IEEE Trans. Wireless Commun.* 16.5 (Mar. 2017), pp. 2891–2906.
- [126] A. Alizadeh and M. Vu. “Load balancing user association in millimeter wave MIMO networks”. In: *IEEE Trans. Wireless Commun.* 18.6 (Mar. 2019), pp. 2932–2945.
- [127] X. Huang, G. Xue, R. Yu, and S. Leng. “Joint scheduling and beamforming coordination in cloud radio access networks with QoS guarantees”. In: *IEEE Trans. Veh. Technol.* 65.7 (Aug. 2015), pp. 5449–5460.

- [128] H. T. Nguyen, H. D. Tuan, T. Q. Duong, H. V. Poor, and W.-J. Hwang. “Nonsmooth optimization algorithms for multicast beamforming in content-centric fog radio access networks”. In: *IEEE Trans. Signal Process.* 68 (Jan. 2020), pp. 1455–1469.
- [129] G. Kwon and H. Park. “Joint user association and beamforming design for millimeter wave UDN with wireless backhaul”. In: *IEEE J. Sel. Areas Commun.* 37.12 (Oct. 2019), pp. 2653–2668.
- [130] A. Pizzo and L. Sanguinetti. “Optimal design of energy-efficient millimeter wave hybrid transceivers for wireless backhaul”. In: *Proc. IEEE WiOpt*. June 2017.
- [131] A. Ortiz, A. Asadi, G. H. Sim, D. Steinmetzer, and M. Hollick. “SCAROS: A Scalable and Robust Self-Backhauling Solution for Highly Dynamic Millimeter-Wave Networks”. In: *IEEE J. Sel. Areas Commun.* 37.12 (Oct. 2019), pp. 2685–2698.
- [132] R. T. Rockafellar and S. Uryasev. “Optimization of Conditional Value-at-Risk”. In: *Journal of Risk* 2 (2000), pp. 21–41.
- [133] H. Levy. *Stochastic Dominance*. Springer Science+Business Media New York, 1998.
- [134] J. Hoydis, S. Cammerer, F. A. Aoudia, A. Vem, N. Binder, G. Marcus, and A. Keller. “Sionna: An Open-Source Library for Next-Generation Physical Layer Research”. In: *arXiv preprint arXiv:2203.11854* (Mar. 2022).
- [135] M. Polese, M. Giordani, A. Roy, D. Castor, and M. Zorzi. “Distributed path selection strategies for integrated access and backhaul at mmWaves”. In: *Proc. IEEE Globecom*. Dec. 2018.
- [136] R. Rockafellar and S. Uryasev. “Conditional value-at-risk for general loss distributions”. In: *Journal of Banking & Finance* 26.7 (2002), pp. 1443–1471.
- [137] G. C. Pflug. “Some Remarks on the Value-at-Risk and the Conditional Value-at-Risk”. In: *Probabilistic Constrained Optimization: Methodology and Applications*. Ed. by S. P. Uryasev. Boston, MA: Springer US, 2000, pp. 272–281. ISBN: 978-1-4757-3150-7. DOI: [10.1007/978-1-4757-3150-7_15](https://doi.org/10.1007/978-1-4757-3150-7_15). URL: https://doi.org/10.1007/978-1-4757-3150-7_15.
- [138] 3GPP. *NR; Overall description; Stage-2*. Technical Specification (TS) 38.300. Version 15.6.0. 3rd Generation Partnership Project (3GPP), June 2019.

Bibliography

- [139] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. 2nd. Edition. MIT Press, 2018.
- [140] Z. Luo and S. Ou. “The almost sure convergence rate of the estimator of optimized certainty equivalent risk measure under α -mixing sequences”. In: *Communications in Statistics - Theory and Methods* 46.16 (Apr. 2017), pp. 8166–8177.
- [141] P. Auer, N. Cesa-Bianchi, and P. Fischer. “Finite-Time Analysis of the Multiarmed Bandit Problem”. In: *Machine Learning* 47.2–3 (May 2002), pp. 235–256.
- [142] 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: *Proc. IEEE CAMAD*. Oct. 2018.
- [143] M. Mezzavilla, M. Zhang, M. Polese, R. Ford, S. Dutta, S. Rangan, and M. Zorzi. “End-to-End Simulation of 5G mmWave Networks”. In: *IEEE Commun. Surveys Tuts.* 20.3 (Apr. 2018), pp. 2237–2263.
- [144] 3GPP. NR; *Study on integrated access and backhaul*. Technical Specification (TS) 38.874. v.0.6.0. 2018.
- [145] C. Madapatha, B. Makki, C. Fang, O. Teyeb, E. Dahlman, M.-S. Alouini, and T. Svensson. “On integrated access and backhaul networks: Current status and potentials”. In: *IEEE Open J. Commun. Soc.* 1 (Sept. 2020), pp. 1374–1389.
- [146] D. Yuan, H.-Y. Lin, J. Widmer, and M. Hollick. “Optimal joint routing and scheduling in millimeter-wave cellular networks”. In: *Proc. IEEE INFOCOM*. Oct. 2018, pp. 1205–1213.
- [147] S. Hur, T. Kim, D. J. Love, J. V. Krogmeier, T. A. Thomas, and A. Ghosh. “Millimeter Wave Beamforming for Wireless Backhaul and Access in Small Cell Networks”. In: *IEEE Transactions on Communications* 61.10 (Jan. 2013), pp. 4391–4403.
- [148] Y. Zhu, Y. Niu, J. Li, D. O. Wu, Y. Li, and D. Jin. “QoS-aware scheduling for small cell millimeter wave mesh backhaul”. In: *Proc. IEEE ICC*. July 2016.
- [149] I. F. Akyildiz, J. M. Jornet, and C. Han. “Terahertz band: Next frontier for wireless communications”. In: *Physical Communication* 12 (Sept. 2014), pp. 16–32.

- [150] A. Singh, M. Andrello, N. Thawdar, and J. M. Jornet. "Design and Operation of a Graphene-Based Plasmonic Nano-Antenna Array for Communication in the Terahertz Band". In: *IEEE J. Sel. Areas Commun.* 38.9 (Sept. 2020), pp. 2104–2117. ISSN: 1558-0008. DOI: [10.1109/JSAC.2020.3000881](https://doi.org/10.1109/JSAC.2020.3000881).
- [151] S. Ghafoor, N. Boujnah, M. H. Rehmani, and A. Davy. "MAC Protocols for Terahertz Communication: A Comprehensive Survey". In: *IEEE Commun. Surveys Tuts.* 22.4 (Fourth Quarter 2020), pp. 2236–2282. ISSN: 1553-877X. DOI: [10.1109/COMST.2020.3017393](https://doi.org/10.1109/COMST.2020.3017393).
- [152] M. Polese, V. Ariyarathna, P. Sen, J. V. Siles, F. Restuccia, T. Melodia, and J. M. Jornet. "Dynamic spectrum sharing between active and passive users above 100 GHz". In: *Communications Engineering* 1.1 (2022), pp. 1–9.
- [153] G. Gougeon, Y. Corre, M. Z. Aslam, S. Bicaïs, and J.-B. Doré. "Assessment of sub-THz Mesh Backhaul Capabilities from Realistic Modelling at the PHY Layer". In: *14th European Conference on Antennas and Propagation (EuCAP)*. 2020, pp. 1–5. DOI: [10.23919/EuCAP48036.2020.9135258](https://doi.org/10.23919/EuCAP48036.2020.9135258).
- [154] A. A. Raja, M. A. Jamshed, H. Pervaiz, and S. A. Hassan. "Performance Analysis of UAV-assisted Backhaul Solutions in THz enabled Hybrid Heterogeneous Network". In: *Conference on Computer Communications Workshops (INFOCOM WKSHPS)*. 2020, pp. 628–633. DOI: [10.1109/INFOCOMWKSHPS50562.2020.9163026](https://doi.org/10.1109/INFOCOMWKSHPS50562.2020.9163026).
- [155] H. Jiang, Y. Niu, B. Ai, Z. Zhong, and S. Mao. "QoS-Aware Bandwidth Allocation and Concurrent Scheduling for Terahertz Wireless Backhaul Networks". In: *IEEE Access* 8 (July 2020), pp. 125814–125825. DOI: [10.1109/ACCESS.2020.3007865](https://doi.org/10.1109/ACCESS.2020.3007865).
- [156] 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, pp. 1–6. DOI: [10.1109/ICC.2018.8422149](https://doi.org/10.1109/ICC.2018.8422149).
- [157] J. M. Jornet and I. F. Akyildiz. "Channel Modeling and Capacity Analysis for Electromagnetic Wireless Nanonetworks in the Terahertz Band".

Bibliography

- In: *IEEE Trans. Wireless Commun.* 10.10 (Aug. 2011), pp. 3211–3221. DOI: [10.1109/TWC.2011.081011.100545](https://doi.org/10.1109/TWC.2011.081011.100545).
- [158] Z. Hossain, Q. Xia, and J. M. Jornet. “TeraSim: An ns-3 extension to simulate Terahertz-band communication networks”. In: *Nano Communication Networks* 17 (Sept. 2018), pp. 36–44.
- [159] A. Ashtari Gargari, A. Ortiz, M. Pagin, A. Klein, M. Hollick, M. Zorzi, and A. Asadi. “Safehaul: Risk-Averse Learning for Reliable mmWave Self-Backhauling in 6G Networks”. In: *IEEE Conference on Computer Communications (INFOCOM)*. New York, USA, 2023.
- [160] P. Sen, J. Hall, M. Polese, V. Petrov, D. Bodet, F. Restuccia, T. Melodia, and J. M. Jornet. “Terahertz Communications Can Work in Rain and Snow: Impact of Adverse Weather Conditions on Channels at 140 GHz”. In: *Proceedings of the 6th ACM Workshop on Millimeter-Wave and Terahertz Networks and Sensing Systems. mmNets ’22*. Sydney, NSW, Australia: Association for Computing Machinery, 2022, pp. 13–18. ISBN: 9781450395090.
- [161] G. Gemmi, R. L. Cigno, and L. Maccari. “On the Properties of Next Generation Wireless Backhaul”. In: *IEEE Transactions on Network Science and Engineering* 10.1 (Jan. 2023), pp. 166–177.
- [162] G. Gemmi, R. Lo Cigno, and L. Maccari. “On Cost-Effective, Reliable Coverage for LoS Communications in Urban Areas”. In: *IEEE Transactions on Network and Service Management* 19.3 (Sept. 2022), pp. 2767–2779. ISSN: 1932-4537.
- [163] Q. Wu and R. Zhang. “Towards Smart and Reconfigurable Environment: Intelligent Reflecting Surface Aided Wireless Network”. In: *IEEE Commun. Mag.* 58.1 (Jan. 2020), pp. 106–112. DOI: [10.1109/MCOM.001.1900107](https://doi.org/10.1109/MCOM.001.1900107).
- [164] S. Abeywickrama, R. Zhang, Q. Wu, and C. Yuen. “Intelligent Reflecting Surface: Practical Phase Shift Model and Beamforming Optimization”. In: *IEEE Trans. Commun.* 68.9 (Sept. 2020), pp. 5849–5863. DOI: [10.1109/TCOMM.2020.3001125](https://doi.org/10.1109/TCOMM.2020.3001125).
- [165] M. Pagin, M. Giordani, A. A. Gargari, A. Rech, F. Moretto, S. Tomasin, J. Gambini, and M. Zorzi. “End-to-End Simulation of 5G Networks Assisted by IRS and AF Relays”. In: *Proc. IEEE MedComNet*. 2022.

- [166] R. Liu, Q. Wu, M. Di Renzo, and Y. Yuan. "A path to smart radio environments: An industrial viewpoint on reconfigurable intelligent surfaces". In: *IEEE Wireless Commun.* 29.1 (Jan. 2022), pp. 202–208.
- [167] C. Liaskos, S. Nie, A. Tsoliariadou, A. Pitsillides, S. Ioannidis, and I. Akyildiz. "Realizing wireless communication through software-defined hypersurface environments". In: *Proc. IEEE WoWMoM*. 2018.
- [168] Y. Yang, S. Zhang, and R. Zhang. "IRS-Enhanced OFDMA: Joint Resource Allocation and Passive Beamforming Optimization". In: *IEEE Wireless Commun. Lett.* 9.6 (June 2020), pp. 760–764. DOI: [10.1109/LWC.2020.2968303](https://doi.org/10.1109/LWC.2020.2968303).
- [169] J. Lee, J. Choi, and J. Kang. "Harmony Search-Based Optimization for Multi-RISs MU-MISO OFDMA Systems". In: *IEEE Wireless Commun. Lett.* 12.2 (Feb. 2023), pp. 257–261. DOI: [10.1109/LWC.2022.3222455](https://doi.org/10.1109/LWC.2022.3222455).
- [170] Y. Guo, Z. Qin, Y. Liu, and N. Al-Dhahir. "Intelligent Reflecting Surface Aided Multiple Access Over Fading Channels". In: *IEEE Trans. Commun.* 69.3 (Mar. 2021), pp. 2015–2027. DOI: [10.1109/TCOMM.2020.3042277](https://doi.org/10.1109/TCOMM.2020.3042277).
- [171] D. Zhang, Q. Wu, M. Cui, G. Zhang, and D. Niyato. "Throughput Maximization for IRS-Assisted Wireless Powered Hybrid NOMA and TDMA". In: *IEEE Wireless Commun. Lett.* 10.9 (Sept. 2021), pp. 1944–1948. DOI: [10.1109/LWC.2021.3087495](https://doi.org/10.1109/LWC.2021.3087495).
- [172] H. Al-Obiedollah, H. A. B. Salameh, K. Cumanan, Z. Ding, and O. A. Dobre. "Self-Sustainable Multi-IRS-Aided Wireless Powered Hybrid TDMA-NOMA System". In: *IEEE Access* 11 (June 2023), pp. 57428–57436.
- [173] Z. Zhang, T. Jiang, and W. Yu. "Learning Based User Scheduling in Reconfigurable Intelligent Surface Assisted Multiuser Downlink". In: *IEEE J. Sel. Topics Signal Process.* 16.5 (May 2022), pp. 1026–1039. DOI: [10.1109/JSTSP.2022.3178213](https://doi.org/10.1109/JSTSP.2022.3178213).
- [174] A. Bansal, K. Singh, B. Clerckx, C.-P. Li, and M.-S. Alouini. "Rate-Splitting Multiple Access for Intelligent Reflecting Surface Aided Multi-User Communications". In: *IEEE Trans. Veh. Technol.* 70.9 (Sept. 2021), pp. 9217–9229. DOI: [10.1109/TVT.2021.3102212](https://doi.org/10.1109/TVT.2021.3102212).

Bibliography

- [175] H. Fu, S. Feng, and D. W. Kwan Ng. "Resource Allocation Design for IRS-Aided Downlink MU-MISO RSMA Systems". In: *Proc. IEEE ICC Wkshps.* 2021.
- [176] B. Zhuo, J. Gu, W. Duan, X. Gu, G. Zhang, M. Wen, and P.-H. Ho. "Partial Non-Orthogonal Multiple Access: A New Perspective for RIS-Aided Downlink". In: *IEEE Wireless Commun. Lett.* 11.11 (Nov. 2022), pp. 2395–2399. DOI: [10.1109/LWC.2022.3204666](https://doi.org/10.1109/LWC.2022.3204666).
- [177] S. Hu, Z. Wei, Y. Cai, C. Liu, D. W. K. Ng, and J. Yuan. "Robust and Secure Sum-Rate Maximization for Multiuser MISO Downlink Systems With Self-Sustainable IRS". In: *IEEE Trans. Commun.* 69.10 (Oct. 2021), pp. 7032–7049. DOI: [10.1109/TCOMM.2021.3097140](https://doi.org/10.1109/TCOMM.2021.3097140).
- [178] X. Mu, Y. Liu, L. Guo, J. Lin, and N. Al-Dhahir. "Capacity and Optimal Resource Allocation for IRS-Assisted Multi-User Communication Systems". In: *IEEE Trans. Commun.* 69.6 (June 2021), pp. 3771–3786. DOI: [10.1109/TCOMM.2021.3062651](https://doi.org/10.1109/TCOMM.2021.3062651).
- [179] ETSI. *Reconfigurable Intelligent Surfaces (RIS); Communication Models, Channel Models, Channel Estimation and Evaluation Methodology*. ETSI GR RIS 003 V1.1.1. 2023.
- [180] M. Rossanese, P. Mursia, A. Garcia-Saavedra, V. Sciancalepore, A. Asadi, and X. Costa-Perez. "Designing, Building, and Characterizing RF Switch-Based Reconfigurable Intelligent Surfaces". In: *Proc. ACM WiNTECH.* 2022.
- [181] G. C. Alexandropoulos, D.-T. Phan-Huy, K. D. Katsanos, M. Crozzoli, H. Wyneersch, P. Popovski, P. Ratajczak, Y. Bénédic, M.-H. Hamon, S. H. Gonzalez, et al. "RIS-enabled smart wireless environments: Deployment scenarios, network architecture, bandwidth and area of influence". In: *EURASIP J. on Wirel. Commun. and Netw.* 2023.1 (Oct. 2023), p. 103.
- [182] L. Yezhen, R. Yongli, Y. Fan, X. Shenheng, and Z. Jiannian. "A novel 28 GHz phased array antenna for 5G mobile communications". In: *ZTE Communications* 18.3 (2020), pp. 20–25.
- [183] V. Jamali, G. C. Alexandropoulos, R. Schober, and H. V. Poor. "Low-to-Zero-Overhead IRS Reconfiguration: Decoupling Illumination and

- Channel Estimation". In: *IEEE Commun. Lett.* 26.4 (Apr. 2022), pp. 932–936. DOI: [10.1109/LCOMM.2022.3141206](https://doi.org/10.1109/LCOMM.2022.3141206).
- [184] Q.-U.-A. Nadeem, A. Kammoun, A. Chaaban, M. Debbah, and M.-S. Alouini. "Asymptotic Max-Min SINR Analysis of Reconfigurable Intelligent Surface Assisted MISO Systems". In: *IEEE Trans. Wireless Commun.* 19.12 (Dec. 2020), pp. 7748–7764. DOI: [10.1109/TWC.2020.2986438](https://doi.org/10.1109/TWC.2020.2986438).
- [185] X. Qian, M. Di Renzo, V. Sciancalepore, and X. Costa-Pérez. "Joint Optimization of Reconfigurable Intelligent Surfaces and Dynamic Metasurface Antennas for Massive MIMO Communications". In: *Proc. IEEE SAM workshop*. 2022.
- [186] A. Rech, M. Pagin, S. Tomasin, F. Moretto, L. Badia, M. Giordani, J. Gambini, and M. Zorzi. "Downlink TDMA Scheduling for IRS-aided Communications with Block-Static Constraints". In: *Proc. IEEE WCNC Wkshps*. 2023.
- [187] Y. Lu, M. Koivisto, J. Talvitie, M. Valkama, and E. S. Lohan. "Positioning-aided 3D beamforming for enhanced communications in mmWave mobile networks". In: *IEEE Access* 8 (2020), pp. 55513–55525.
- [188] X. Tan, Z. Sun, D. Koutsonikolas, and J. M. Jornet. "Enabling Indoor Mobile Millimeter-wave Networks Based on Smart Reflect-arrays". In: *Proc. IEEE INFOCOM*. 2018. DOI: [10.1109/infocom.2018.8485924](https://doi.org/10.1109/infocom.2018.8485924).
- [189] Z.-Q. He and X. Yuan. "Cascaded Channel Estimation for Large Intelligent Metasurface Assisted Massive MIMO". In: *IEEE Commun. Lett.* 9.2 (Feb. 2020), pp. 210–214. DOI: [10.1109/LWC.2019.2948632](https://doi.org/10.1109/LWC.2019.2948632).
- [190] J. Rains, A. Tukmanov, Q. Abbasi, and M. Imran. "RIS-enhanced MIMO channels in urban environments: Experimental insights". In: *Proc. EuCAP*. 2024.
- [191] T. S. Rappaport, S. Sun, and M. Shafi. "Investigation and Comparison of 3GPP and NYUSIM Channel Models for 5G Wireless Communications". In: *Proc. IEEE VTC-Fall*. 2017.
- [192] L. Anchora, L. Badia, E. Karipidis, and M. Zorzi. "Capacity gains due to orthogonal spectrum sharing in multi-operator LTE cellular networks". In: *Proc. IEEE ISWCS*. 2012.

Bibliography

- [193] A. L. Swindlehurst, G. Zhou, R. Liu, C. Pan, and M. Li. "Channel Estimation With Reconfigurable Intelligent Surfaces—A General Framework". In: *Proc. IEEE* 110.9 (Sept. 2022), pp. 1312–1338. doi: [10.1109/JPROC.2022.3170358](https://doi.org/10.1109/JPROC.2022.3170358).
- [194] L. Rokach and O. Maimon. "Clustering methods". In: *Data mining and knowledge discovery handbook*. Springer, 2005, pp. 321–352.
- [195] S. Lloyd. "Least squares quantization in PCM". In: *IEEE Trans. Inf. Theory* 28.2 (Mar. 1982), pp. 129–137. doi: [10.1109/TIT.1982.1056489](https://doi.org/10.1109/TIT.1982.1056489).
- [196] F. Murtagh and P. Contreras. "Algorithms for hierarchical clustering: an overview". In: *Wiley Interdisciplinary Rev.: Data Mining Knowl. Discovery* 2.1 (2012), pp. 86–97.
- [197] L. Kaufman and P. J. Rousseeuw. "Clustering by means of medoids". In: *Rep. Fac. Math. Inf.* Vol. 87. 3. 1987.
- [198] M. Van der Laan, K. Pollard, and J. Bryan. "A new partitioning around medoids algorithm". In: *J. Stat. Computat. Simul.* 73.8 (2003), pp. 575–584.
- [199] D. Xu and Y. Tian. "A comprehensive survey of clustering algorithms". In: *Annals of Data Science* 2 (June 2015), pp. 165–193.
- [200] 3GPP. 5G; NR; *Physical channels and modulation*. TS 38.211 (Rel. 16). 2020.
- [201] M. Rivera, M. Chegini, W. Jaafar, S. Alfattani, and H. Yanikomeroglu. "Optimization of Quantized Phase Shifts for Reconfigurable Smart Surfaces Assisted Communications". In: *Proc. IEEE CCNC*. 2022.
- [202] M. Sabin and R. Gray. "Global convergence and empirical consistency of the generalized Lloyd algorithm". In: *IEEE Trans. Inf. Theory* 32.2 (Feb. 1986), pp. 148–155. doi: [10.1109/TIT.1986.1057168](https://doi.org/10.1109/TIT.1986.1057168).