

Dissertation

Towards Massive Connectivity via Uplink Code-Domain NOMA

Author

Dipl.-Ing. Bashar Tahir

Matriculation Number: 01476041

Advisor

Associate Prof. Dipl.-Ing. Dr.techn. Stefan Schwarz

Assisting Advisor

Univ.Prof. Dipl.-Ing. Dr.techn. Markus Rupp

Reviewers

Prof. Daniel Benevides da Costa, Ph.D.

Senior Lecturer (Associate Prof.) Yuanwei Liu, Ph.D.

Institute of Telecommunications
Technische Universität Wien
Vienna, Austria

February, 2022

Abstract

Future mobile networks are envisioned to provide wireless access to a massive number of devices. The substantial increase in connectivity comes mainly from machine-type communication (MTC), for which a large number of low-rate transmissions take place. Accommodating access for such a large number of user equipments (UEs) can be inefficient if applied to current network architectures, which are mainly based on orthogonal multiple access (OMA) and scheduling-based transmissions. This is due to the resulting control overhead and increased access-delay. The framework of non-orthogonal multiple access (NOMA) has attracted attention recently as a promising solution to tackle these issues. It allows multiple UEs to access the network simultaneously over the same resources, and provides naturally, the support for grant-free access, in which no explicit scheduling of the UEs is required.

Motivated by the potential benefits of NOMA in enabling massive connectivity, this dissertation focuses on studying uplink code-domain NOMA, where multiple UEs access the network via short non-orthogonal spreading signatures. The dissertation consists of three major parts corresponding to the three building-blocks of the NOMA communication chain: the transmitter, the receiver, and the channel. In the first part, we consider the codebook design problem, that is, the design of the spreading signatures across the different UEs. The formulation we consider leads us to constructing codebooks of a Grassmannian nature. We propose an iterative algorithm for constructing such codebooks, with close-to-optimal correlation properties, leading to enhanced performance under low-complexity suboptimal detection. We then extend the codebook design problem to the case where the UEs are available as groups, such as in cells, or spatial clusters. We propose to jointly design the codebooks across the different groups via an alternating projection algorithm, and show that such a joint design can improve the performance of UEs suffering from strong inter-group interference.

The second part of the dissertation is then concerned with the receiver side, aiming to reduce the complexity of the detection procedure. We consider the user activity detection in the context of grant-free access under a practical frame-structure. We formulate the activity detector based on subspace methods, and address the influence of the channel. Namely, we show that strong time-frequency correlation of the channel can prevent successful detection of the active set of UEs. To address that, we propose overlaying the pilot sequences with user-specific masking sequences, which results effectively in a decorrelation of the channel. We also consider, on the other hand, the influence of strong time-frequency selectivity, for which we investigate different pilots' allocation strategies. We then focus on reducing the detection

complexity of the data part of the transmission. We utilize the time-frequency correlation of the channel to reduce the number of calculated filters for the spreading blocks over the time-frequency grid. Also, by assuming the base station (BS) is equipped with a sufficient number of receive antennas, we show how that combination of the code- and spatial-domains allows us to replace exact minimum mean square error (MMSE) filtering with a low-complexity approximation requiring no inverse calculation, while resulting only in a small performance loss.

The last part investigates the controllability of the channel via reconfigurable intelligent surfaces (RISs). We consider first a RIS-assisted two-UE NOMA uplink, where part of the surface elements are configured to boost the signal of the first UE, while the other part is used to boost the second one. By approximating the receive powers as gamma random variables, tractable expressions for the outage probability under interference cancellation (IC) are derived. We show how the optimization of the RIS impacts the NOMA detection performance, and identify robust operation points that guarantee reliable link quality for both UEs. Finally, we consider the combination of a K -UE code-domain NOMA uplink with RISs, in the context of a cluster-based massive multiple-input multiple-output (MIMO) deployment. We investigate the optimization of the RIS under such a setup, and show a solution based on a semi-definite relaxation of the problem. The results show that our proposed approach can substantially increase the number of UEs supported by the system.

Acknowledgments

First, I would like to thank my advisor, Prof. Stefan Schwarz, for his support and assistance throughout this work, for welcoming me into his lab, and for giving me the freedom in pursuing my own ideas. I am also thankful to my second advisor, Prof. Markus Rupp, for his support and encouragement over the years, and I am grateful to him for giving me the opportunity to join the group and start my academic career. I am also grateful to the Christian Doppler Laboratory and the involved industrial partners for their financial support.

A special thanks to Prof. Daniel Benevides da Costa and Prof. Yuanwei Liu for agreeing to review my dissertation and be part of my defense examination board.

During my stay the last couple of years at the Institute of Telecommunications, I have had the pleasure to meet many awesome people, with whom I have made friends and shared lots of constructive discussions and endless hours of fun. I am thankful to them, and I am glad that I have got the chance to know them.

Finally, I want to thank my family. They were always there for me, and have always encouraged me to go after my goals.

Contents

1	Introduction	1
1.1	Motivation and Scope of Dissertation	1
1.2	Structure and Contribution	4
2	System Model and Methodology	7
2.1	K -UE NOMA Uplink	7
2.2	Received Signal Processing	10
2.3	Simulation Setup	12
3	Transmit-Side Optimization: NOMA Codebook Design	15
3.1	The Codebook Design Problem	15
3.2	Construction of Grassmannian Codebooks	17
3.2.1	Preliminaries	18
3.2.2	Minimizing Coherence by Euclidean Distance Maximization	19
3.2.3	Collision-Based Packing	21
3.2.4	Algorithm Performance	24
3.3	Multi-Group Joint Codebook Design	27
3.3.1	Cross-Codebook Optimization	28
3.3.2	Alternating Projection	31
3.3.3	Case Study: Multi-Cell Deployment	33
3.4	Final Remarks	35
4	Receive-Side Processing: NOMA Activity and Data Detection	37
4.1	Activity Detection in Grant-Free NOMA	38
4.2	Optimizing Activity Detection via MUSIC	40
4.2.1	Considered Model for Subspace Detection	40
4.2.2	Impact of Strong Time-Frequency Correlation	44
4.2.3	Impact of Strong Time-Frequency Selectivity	48
4.3	Reducing Data Detection Complexity	51
4.3.1	Exploiting Time-Frequency Correlation	52
4.3.2	Exploiting the Spatial Domain	57
4.4	Final Remarks	60

5	Controlling the Channel: RIS-Assisted Uplink NOMA	61
5.1	Reconfigurable Intelligent Surfaces	62
5.2	RIS-Assisted Two-User NOMA Uplink	64
5.2.1	Two-UE System Model	64
5.2.2	Outage Analysis	66
5.2.3	Analysis of an Example Scenario	71
5.3	Combination with Code-Domain NOMA	75
5.3.1	System Model Combined with Code-Domain NOMA	76
5.3.2	Sum-Rate Optimized Phase-Shifts	78
5.3.3	Proposed Optimization Approach	79
5.3.4	Investigation of an Example Scenario	81
5.4	Final Remarks	84
6	Conclusion and Outlook	85
6.1	Summary of Contribution	85
6.2	Possible Future Work Directions	86
	List of Abbreviations	89
	Notation	91
	Bibliography	93

1

Introduction

Mobile networks have been evolving rapidly over the past couple of decades. Starting with analog networks providing mainly voice-only services in the early 1980s, to supporting today's digitally-connected world allowing for, literally, an infinite amount of content types. The growth of these networks and the adoption of new technologies have been on the rise since then, and this trend is expected to continue over the next decade. Consider, for example, the recent mobility report by Ericsson of November 2021 [1]. It shows almost a 10-fold increase of the global mobile traffic from 6.7 exabytes/month in 2016 to 65 exabytes/month by the end of 2021. This is forecast to increase to 288 exabytes/month by the end of 2027. The massive growth in traffic comes generally from two aspects: first, the increase of services requiring high data-rate transmission, such as various smartphone applications requiring continuous Internet access, high definition video streaming, virtual reality, etc; second, the increase in the number of connected users or devices, themselves. The increase in data-rate comes naturally as more throughput-demanding applications are developed. The increase in the number of connections, however, not only comes from the direct increase of the human mobile subscribers, but also from the envisioned massive connectivity of machine-type communication (MTC). Machine-type traffic is typically characterized by its low data-rate transmission coming from a large number of devices that are trying to access the network, such as home devices, sensors on the streets, controllers in factories, falling under the general umbrella of the Internet-of-things (IoT). Therefore, it has been important not only to improve the data-rate per device, but also to evolve mobile networks such that they are capable of supporting a large number of accessing devices.

1.1 Motivation and Scope of Dissertation

Multiple access (MA) techniques have been key in enabling the evolution into what mobile networks are today. So far, commercial networks have been mainly based on orthogonal multiple access (OMA). This began with the single-carrier era of frequency-division multiple access (FDMA) in 1st generation (1G) networks, time-

division multiple access (TDMA) in 2G, code-division multiple-access (CDMA) in 3G, to the modern era of multi-carrier modulation with orthogonal frequency-division multiple access (OFDMA) in 4G/long-term evolution (LTE), maintained over to the recently deployed 5G standard. The main feature of OMA is the operational simplicity: the base station (BS) provides access to the user equipments (UEs) by allocating them disjoint sets of time, frequency, and/or code resources. This allows for low-complexity detection at the receiver side, since the transmissions of the UEs do not interfere with each other, i.e., orthogonal, and therefore simple per-UE detection is optimal.

Non-Orthogonal Multiple Access

Recently, non-orthogonal multiple access (NOMA) has attracted major attention within the research community. With NOMA, the UEs contest the same time-frequency resources and therefore are intentionally allowed to interfere with each other. This can be done via pure power-domain superposition, or in combination with code-domain techniques, such as spreading with short non-orthogonal sequences, providing further interference suppression via code-domain processing.

NOMA can bring many benefits to the table. In the context of downlink transmission, NOMA has been shown to have a larger achievable rate region compared to OMA [2–4]; moreover, it can achieve optimal points on that region with improved fairness between the UEs compared to OMA. These gains can also be observed in the context of cluster-based multiple-input multiple-output (MIMO) deployments; UEs lying within the same spatial cluster, i.e., served by the same beam, can be served via NOMA, which further generalizes the gains to MIMO systems [5–7]. From a connectivity point-of-view, the number of UEs that can simultaneously access the network with OMA depends on the orthogonal resources' granularity of the system, and how the BS schedules the UEs on these resources. NOMA paves the way into supporting massive connectivity by allowing multiple UEs to access the network over the same resources. This, on the one hand, increases the number of UEs accessing the network simultaneously, and on the other hand, reduces the access latency to the network [8–10]. This is especially important on the uplink side, since that is where most of the massive MTC overhead will take place. The disadvantage of NOMA transmissions is the increased detection complexity, as multiple UEs now interfere with each other. This requires, in general, receivers performing joint multi-user detection, typically involving the application of interference cancellation (IC).

Grant-Free Access

When it comes to granting access to UEs, conventional systems (e.g., LTE) are grant-based. Meaning that, in order for the UE to transmit its data, it has to go through a scheduling-grant procedure with the BS, in which it asks the BS for access and waits until the BS schedules it, and only when its agreed-on scheduling instant arrives, it may transmit [11]. Such an access strategy is efficient when the number

of UEs accessing the network is relatively small, and the transmitted data packet in the end is large enough. However, for massive connectivity targeting machine-type traffic, this can be an issue. On the one hand, large number of UEs having to go through the scheduling-grant procedure can cause a large control overhead at the BS [12]. On the other hand, for machine-type traffic, the data transmitted in the end might be of very short packets (e.g., sensor reading), which can be comparable in size to the control signaling required to setup the connection in the first place. This might result in an inefficient utilization of the network resources [13].

Grant-free access can address these issues [14,15]. It allows the UEs to, more or less, transmit their data on their own, without having to be explicitly scheduled by the BS. However, since the UEs transmit on their own, it can happen that multiple UEs choose to contest the same resources, thus causing a collision of the transmitted packets. At this point, the framework of NOMA comes into action, as it provides the capability to manage the multi-user interference. Therefore, the combination, grant-free NOMA, has received wide attention in the literature and has shown the capability to resolve the collisions and support a large number of UEs accessing the network in a grant-free manner [16–19].

Reconfigurable Intelligent Surfaces

Another technology that is envisioned to have a great role in the evolution of mobile networks are reconfigurable intelligent surfaces (RISs) [20,21]. These surfaces consist of a large number of electronically controllable elements that can modify electromagnetic waves impinged on them. They can produce reflected and/or transmitted waves with modified amplitude, phase, frequency, and/or polarization [22]. The most common adjustment considered is the phase-shifting of the incoming waves; by jointly phase-shifting the waves across the different elements of the surface, it is possible to reshape the propagation environment by focusing (or beamforming) the waves towards a certain UE, extending the coverage area, modifying the spatial structure of the channel in the context of MIMO systems, and more [23,24]. All of that is achieved in a passive manner via phase-shifting across the elements, without requiring active radio-frequency power.

The combination of RISs with NOMA has gathered attention recently, showing the potential of improving the system energy efficiency, sum-rate, and outage performance [25–31]. Perhaps the most relevant aspect in such a combination is that the optimization of the RIS directly impacts the detection order under IC. This is especially pronounced in the uplink, as the RIS would be capable of adjusting the receive powers of the different UEs on the fly.

Scope of Work

Motivated by the benefits and flexibility that NOMA provides, this dissertation investigates various aspects with respect to the communication's chain of the NOMA transmission. Our concern is with enabling massive connectivity in the uplink, and

therefore our focus in this dissertation will be on uplink code-domain NOMA. The considered system then corresponds to the case where multiple UEs try to access the network over the same time-frequency resources, with their data spread by short non-orthogonal signatures (or sequences). The dissertation consists of three parts:

- In the first part, we focus on the transmitter side, i.e., the UEs. We consider the problem of designing the spreading signatures, or codebook, employed by the UEs, where we target designs having low cross-correlation between the different signatures.
- In the second part, we turn our attention to the receiver side, i.e., the BS. Our focus is then to describe practical receiver implementations allowing low-complexity activity detection in the context of grant-free access, and reducing the data detection complexity in heavily overloaded systems with multiple receive antennas at the BS.
- Finally, in the last part, we focus on the channel, and investigate its controllability via RISs. We attempt to characterize the statistical behavior of RIS-assisted NOMA systems and then investigate the optimization of RISs in combination with code-domain NOMA.

1.2 Structure and Contribution

In the following, we describe the structure of the dissertation and the corresponding contributions in detail. We also refer to the publications in which the contributions have been first developed or investigated. Lists of the abbreviations and notation used throughout this work can be found at the end of dissertation.

Chapter 2 – System Model and Methodology

The second chapter provides a description of the system model utilized throughout the dissertation. We derive the input-output relationship of a K -UE code-domain NOMA uplink, with the assumption that the BS is equipped with multiple receive antennas. In our work, we utilize a frame-structure that is similar to LTE/5G, and therefore in this chapter we illustrate it in detail. Finally, we describe how the simulations are carried out and elaborate on some of the employed metrics.

Multiple parts of the framework for link-level simulations have been first developed and implemented as part of the Vienna 5G Link-Level Simulator [32]:

- 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,” *EURASIP Journal on Wireless Communications and Networking*, vol. 2018, no. 1, p. 226, Sep. 2018.

Chapter 3 – Transmit-Side Optimization: NOMA Codebook Design

We start with the transmit-side optimization. We consider the codebook design problem, that is, how to design the spreading signatures for such an uplink code-domain system. Under the assumption of no feedback-loop between the BS and UEs, which would be the case in grant-free access, the design we end up with is of a Grassmannian nature. We propose an iterative construction algorithm that is capable of constructing Grassmannian codebooks with a fast convergence rate compared to other algorithms. We then extend the codebook construction to the case where the UEs can be assigned into groups, such as UEs in different cells or spatial clusters. Under such a setup, we propose to jointly design the codebooks across the different groups, with the aim of reducing their cross-correlation, while simultaneously preserving their internal correlation structure. We formulate a method for finding such codebooks by employing an alternating projection algorithm.

The proposed construction methods in this part have been published in [33, 34]:

- B. Tahir, S. Schwarz, and M. Rupp, “Constructing Grassmannian Frames by an Iterative Collision-Based Packing,” *IEEE Signal Processing Letters*, vol. 26, no. 7, pp. 1056–1060, 2019.
- B. Tahir, S. Schwarz, and M. Rupp, “Joint Codebook Design for Multi-Cell NOMA,” in *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 4814–4818.

Chapter 4 – Receive-Side Processing: NOMA Activity and Data Detection

In this chapter, the focus is on the processing of the received signal at the BS. Our goal is primarily to describe a detection procedure that is manageable in practice. The first part deals with activity detection in the context of grant-free access. We start by giving a brief overview of grant-based and grant-free access, and then we formulate the activity detection via subspace methods, under the assumption of an LTE/5G-like frame-structure. We then propose to improve the activity detection under strong time-frequency correlation of the channel by applying masking sequences. We also consider the other extreme of having strong time-frequency selectivity, and investigate possible pilot reallocation strategies. In the second part, we turn our attention to reducing the detection complexity of the data part of the transmission. By utilizing the time-frequency correlation of the channel, we show how it is possible to greatly reduce the number of calculated filters (equalizers) across the time-frequency frame. Then, finally, with the aid of the multiple receive antennas at the BS, we show how it is possible to reduce the calculation complexity of the individual filters themselves.

The framework developed throughout this chapter has been investigated in-full or in-part in [35–37]:

- B. Tahir, S. Schwarz, and M. Rupp, “Low-Complexity Detection of Uplink NOMA by Exploiting Properties of the Propagation Channel,” in *ICC 2020 - 2020 IEEE International Conference on Communications (ICC)*, 2020, pp. 1–6.
- B. Tahir, S. Schwarz, and M. Rupp, “Collision Resilient V2X Communication via Grant-Free NOMA,” in *2020 28th European Signal Processing Conference (EUSIPCO)*, 2021, pp. 1732–1736.
- B. Tahir, S. Schwarz, and M. Rupp, “Impact of Channel Correlation on Subspace-Based Activity Detection in Grant-Free NOMA,” *submitted to the 2022 IEEE 95th Vehicular Technology Conference (VTC2022-Spring)*, 2022. Preprint available: <https://arxiv.org/abs/2202.11161>.

Chapter 5 – Controlling the Channel: RIS-Assisted Uplink NOMA

The last chapter deals with the remaining part of the communication chain, the channel. We investigate its controllability by the deployment of RISs, and how they can be optimized to boost the performance of NOMA systems. We begin with a two-UE power-domain NOMA uplink assisted by a RIS. Namely, we consider the case where the RIS elements are split between the two UEs, i.e., part of the RIS elements are used to boost the signal of the first UE, while the remaining part is used for the second one. Under such a setup, we characterize the outage performance by approximating the received powers of the UEs as gamma random variables. This allows us to arrive at closed-form expressions for the outage under IC, which helps us understand how the RIS impacts the performance of the system. We then consider the combination of RISs with code-domain NOMA in the context of a cluster-based massive MIMO setup, in which we investigate how the RIS can be configured in order to increase the number of supported UEs.

The analysis and results of this chapter have appeared first in [38–40]:

- B. Tahir, S. Schwarz, and M. Rupp, “Analysis of Uplink IRS-Assisted NOMA Under Nakagami- m Fading via Moments Matching,” *IEEE Wireless Communications Letters*, vol. 10, no. 3, pp. 624–628, 2021.
- B. Tahir, S. Schwarz, and M. Rupp, “Outage Analysis of Uplink IRS-Assisted NOMA under Elements Splitting,” in *2021 IEEE 93rd Vehicular Technology Conference (VTC2021-Spring)*, 2021, pp. 1–5.
- B. Tahir, S. Schwarz, and M. Rupp, “RIS-Assisted Code-Domain MIMO-NOMA,” in *2021 29th European Signal Processing Conference (EUSIPCO)*, 2021, pp. 821–825.

2

System Model and Methodology

This chapter introduces the main system model utilized throughout the rest of this work. A NOMA uplink consisting of K UEs is considered showing the signal flow starting from the transmit-side at the UEs, ending with the receive-side processing at the BS. Based on that system model, we derive expressions for the received signal at the BS and the required processing stages for signal detection and decoding. We illustrate the utilized frame-structure for the transmissions, which follows the LTE/5G standards. Furthermore, we elaborate on how the simulations are carried out and what metrics are used.

2.1 K -UE NOMA Uplink

We consider an uplink consisting of K UEs transmitting via orthogonal frequency-division multiplexing (OFDM) over the same time-frequency resources. We assume the transmissions to be synchronized, even in the case of grant-free access for which the BS does not necessarily coordinate the transmissions of the UEs. This can be achieved via a periodically broadcast signal by the BS containing general system information. Instead of mapping the symbols directly to the time-frequency grid, we consider a code-domain NOMA system, where each data symbol is spread via a short spreading signature (sequence) over multiple resource-elements (REs). This is illustrated in Figure 2.1, where the k -th UE maps its i -th data symbol $x_{k,i}$ by multiplying it with a spreading signature \mathbf{s}_k . Note that direct bits-to-signature mapping is also possible, as commonly done in sparse-code multiple access (SCMA) [41], allowing for a shaping gain. However, in this work, we stick to dense-spreading applied to the quadrature amplitude modulation (QAM) symbols, due to its lower detection complexity.

The increase in signal dimensionality due to the spreading allows for multi-user interference suppression, as we will see later. However, this comes at the cost of decreased spectral efficiency, since each data symbol now occupies multiple REs. In the example figure below, the spreading length is $L = 4$, resulting in a decrease of rate by a factor of four. Therefore, it is important to balance between the reduction in spectral efficiency of the single-UE itself, and the benefit that the spreading brings

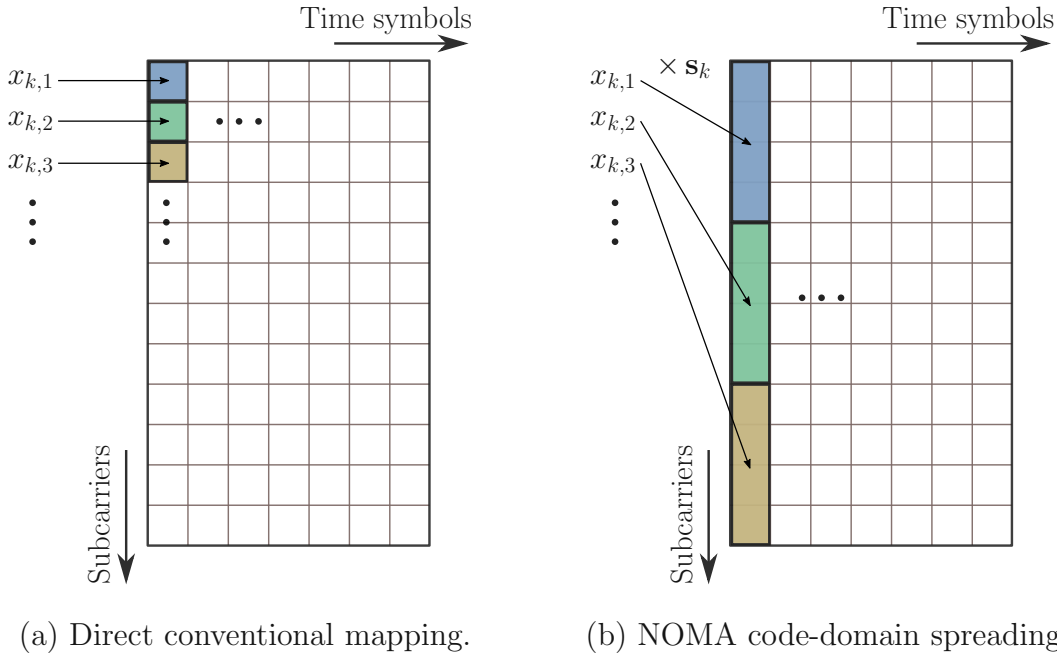


Figure 2.1: Illustration of the data symbols' mapping in conventional systems compared to NOMA code-domain spreading. A resources' region consisting of 12 subcarriers and 7 OFDM time symbols is shown.

with the multi-user interference suppression. In other words, if the number of UEs K is very small, then it makes sense to employ short spreading, while if we expect K to be high, then longer spreading is required, since the performance in that case would be limited by the multi-user interference. Also, in the figure, the spreading is done along the frequency-direction, and this is what we will assume in this work. Nonetheless, performing it along the time-direction, or both (i.e., 2D spreading), is also a possibility.

We adopt a frame-structure that is based on the LTE/5G standards. The data transmission is done on a subframe-basis, consisting of two resource-blocks (RBs) in time. Each RB consists of 12 subcarriers and 7 OFDM symbols. In the middle of each RB, the UE inserts a pilot sequence of length L_p . This is used for channel estimation, and as we will see later, for possibly performing NOMA activity detection in grant-free access. The frame-structure is illustrated in Figure 2.2. We assume that, generally, the data and pilots utilize different spreading signatures. Therefore, the subframe consists of data- and pilots-blocks with different spreading lengths. In the figure, the data spreading length is $L = 4$, while the length of the pilot sequences is $L_p = 12$. Long pilot sequences are required in order to perform channel estimation and robust activity detection for a large number of simultaneously active UEs. For the data part, the sequences would be too short for that purpose, in order to allow for higher spectral efficiency, as discussed before. Therefore, we utilize them solely for data detection.

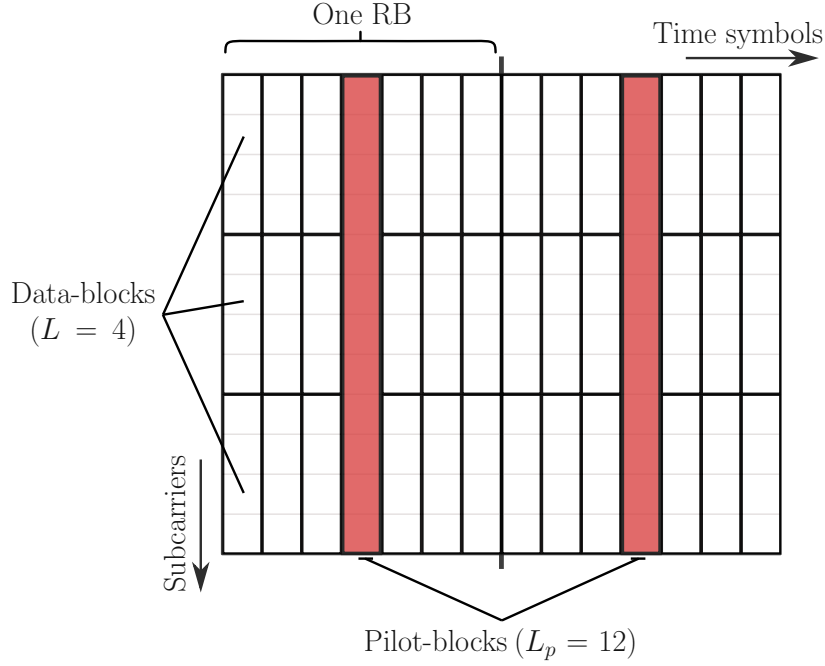


Figure 2.2: Considered frame-structure with 12 subcarriers and 14 OFDM symbols.

After the K single-antenna UEs spread their data symbols and insert their pilots, OFDM modulation is performed and the signal is transmitted over the wireless channel. The BS, equipped with N_R receive antennas, observes the superposition (i.e., the sum) of the signals from the different UEs. After applying fast-Fourier-transform (FFT) at the BS, the received baseband signal at the i^{th} data-block and r^{th} antenna is given by

$$\mathbf{y}_i^{(r)} = \sum_{k=1}^K \sqrt{LP_k} h_{k,i}^{(r)} \mathbf{s}_k x_{k,i} + \mathbf{n}_i^{(r)}, \quad i = 1, 2, \dots, B_d, \quad (2.1)$$

where $P_k \in \mathbb{R}^+$ is the average transmit power of the k^{th} UE, $h_{k,i}^{(r)} \in \mathbb{C}$ is the fading coefficient of the k^{th} UE at the i^{th} data-block and r^{th} antenna, $\mathbf{s}_k \in \mathbb{C}^{L \times 1}$ is the unit-norm spreading signature of the k^{th} UE, $x_{k,i} \in \mathbb{C}$ is the i^{th} transmit symbol of the k^{th} UE, $\mathbf{n}_i^{(r)} \in \mathbb{C}^{L \times 1}$ is the Gaussian noise at the i^{th} block and r^{th} antenna, and B_d is the total number of data-blocks. In the model above, it is assumed that the fading remains constant along the spreading interval, which holds well as an assumption for short spreading, as we will see later in Chapter 4. Moreover, we assume the noise to be circularly symmetric with covariance $\sigma_n^2 \mathbf{I}_L$, where \mathbf{I}_L is the identity matrix of size L . Stacking the signals from the different receive antennas

into a single vector \mathbf{y}_i , we get

$$\mathbf{y}_i = \begin{bmatrix} \mathbf{y}_i^{(1)} \\ \mathbf{y}_i^{(2)} \\ \vdots \\ \mathbf{y}_i^{(N_R)} \end{bmatrix} = \begin{bmatrix} \sum_{k=1}^K \sqrt{LP_k} h_{k,i}^{(1)} \mathbf{s}_k x_{k,i} \\ \sum_{k=1}^K \sqrt{LP_k} h_{k,i}^{(2)} \mathbf{s}_k x_{k,i} \\ \vdots \\ \sum_{k=1}^K \sqrt{LP_k} h_{k,i}^{(N_R)} \mathbf{s}_k x_{k,i} \end{bmatrix} + \begin{bmatrix} \mathbf{n}_i^{(1)} \\ \mathbf{n}_i^{(2)} \\ \vdots \\ \mathbf{n}_i^{(N_R)} \end{bmatrix}. \quad (2.2)$$

Let $\mathbf{h}_{k,i} = [h_{k,i}^{(1)}, h_{k,i}^{(2)}, \dots, h_{k,i}^{(N_R)}]^T$ be channel vector across the receive antennas of the k^{th} UE at block i , and \mathbf{n}_i be the stacked noise vector, then (2.2) can be equivalently written as

$$\mathbf{y}_i = \sum_{k=1}^K \sqrt{LP_k} (\mathbf{h}_{k,i} \otimes \mathbf{s}_k) x_{k,i} + \mathbf{n}_i. \quad (2.3)$$

where \otimes is the Kronecker product. Let $\mathbf{H} = [\mathbf{h}_{1,i}, \mathbf{h}_{2,i}, \dots, \mathbf{h}_{K,i}]$ and $\mathbf{S} = [\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_K]$ be the matrices of the spatial and code-domain signatures of the UEs, respectively. Moreover, let $\mathbf{P} = \text{diag}(LP_1, LP_2, \dots, LP_K)$. We can now write (2.3) in matrix-vector notation as

$$\mathbf{y}_i = (\mathbf{H}_i * \mathbf{S}_i) \mathbf{P}^{1/2} \mathbf{x}_i + \mathbf{n}_i, \quad (2.4)$$

where $*$ is the column-wise Khatri–Rao product and $\mathbf{x}_i = [x_{1,i}, x_{2,i}, \dots, x_{K,i}]^T$. Let $\mathbf{G}_i = (\mathbf{H}_i * \mathbf{S}_i) \mathbf{P}^{1/2}$, we finally have

$$\mathbf{y}_i = \mathbf{G}_i \mathbf{x}_i + \mathbf{n}_i, \quad i = 1, 2, \dots, B_d. \quad (2.5)$$

As can be seen, our final system matrix, $\mathbf{G}_i \in \mathbb{C}^{LN_R \times K}$, shows a structure that depends on the spatial properties of the channels across the different UEs, together with the code-domain spreading signatures they are utilizing.

Note that we sometimes drop the index i from the expressions, when the data-block index is not relevant for the topic considered, such as in Chapters 3 and 5. As for the pilot-blocks, a similar expression to (2.5) exists. We will consider it later in detail in Chapter 4. Moreover, the channel model will be extended in Chapter 5 in order to support RISs. For now, we will keep the RISs out.

2.2 Received Signal Processing

The receiver at the BS generally needs to perform three tasks: first, identify the active set of UEs in the context of grant-free access; second, perform channel estimation; third, detect and decode the data transmission. This is summarized in Figure 2.3 below, where the data detection is shown combined with IC. We will consider grant-free access and activity detection later in Chapter 4, and rather focus

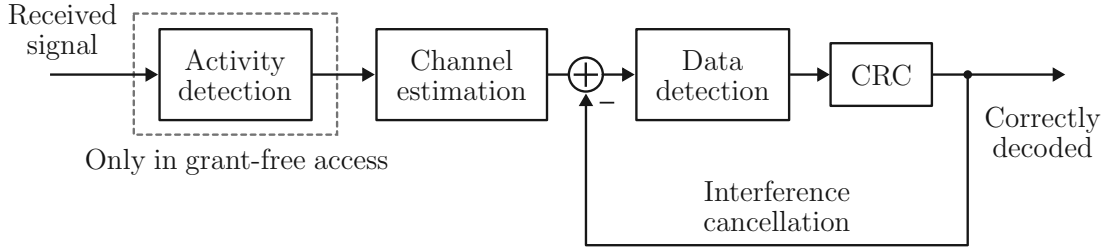


Figure 2.3: Receiver chain at the BS.

here on the data detection part. As for the channel estimation, using the received signal at the pilot-blocks at each of the receive antennas, the channel coefficients are estimated via least-squares (LS). The channel for data-blocks is then obtained by interpolation between the estimated channel coefficients at the pilot-blocks. We will also consider this later in Chapter 4.

Having the spreading signatures known and the channel coefficients estimated, we then have an access to \mathbf{G}_i and the system model in (2.5) is in action. Under minimum mean square error (MMSE) equalization, the estimated data symbols of the UEs for data-block i is then given by

$$\mathbf{x}_i^{\text{MMSE}} = \mathbf{G}_i^H (\mathbf{G}_i \mathbf{G}_i^H + \sigma_n^2 \mathbf{I}_{LN_R})^{-1} \mathbf{y}_i, \quad i = 1, 2, \dots, B_d. \quad (2.6)$$

Once the symbols' estimate for the whole subframe is obtained, channel decoding is performed and the UEs' transmissions are checked by a cyclic-redundancy-check (CRC). All the UEs that pass the CRC have their signals reconstructed and canceled from the received signal. The data detection in (2.6) is then repeated again over the cleared-up signal. Here, all the UEs that pass the CRC are canceled, and therefore the IC scheme mostly considered in this dissertation is parallel IC (PIC). The IC is repeated until no more UEs are left or a maximum number of iterations is reached. The use of PIC allows for a faster detection, as multiple UEs may have sufficient signal-to-interference-plus-noise ratio (SINR) for decodability after the MMSE filtering, and therefore forcing successive IC may result in higher detection latency. This holds especially true when the BS is equipped with many receive antennas, or many digital chains (in the case of hybrid analog/digital architectures), since further interference suppression is achieved by the spatial-domain. Under a suitable configuration of the UEs' transmission rate, it is also possible that all UEs are detected correctly with just a single pass of MMSE filtering, if the number of UEs K is comparable to the product of L and N_R .

2.3 Simulation Setup

In this work, we employ two types of simulations: link-level and outage SINR-based simulations. With link-level simulations, the whole transceiver operation is simulated, starting from the generation of the bits, channel encoding/decoding, symbols mapping/demapping, modulation/demodulation, channel convolution/equalization, etc. In other words, we simulate everything up to the signal sample-level. This type of simulation is used when we would like to investigate block error ratio (BLER) performance, or when we would like to simulate in detail the performance of a certain signal-processing stage. For this type of simulation, many of the signal processing capabilities, such as channel coding, modulation, and channel generation is done through the Vienna 5G Link-Level Simulator [32]¹. For network-wide simulation relying on link-level abstraction, the Vienna 5G System-Level Simulator [42] is also offered. As for the SINR-based simulations, they are used when we would like to characterize the outage probability or outage performance given a certain outage threshold, without having to simulate all sample-level operations. This is beneficial if we would like to run complex simulations that require a large number of simulation repetitions and we are more interested in the behavior of the system rather than exact link-level performance. This type of simulation is utilized throughout Chapter 5, when we attempt to characterize and optimize operation with RISs.

Finally, in the following, we list some of the used parameters and metrics in the simulations and elaborate on their meaning.

Average SNR (over log-domain)

This denotes the signal-to-noise ratio (SNR) of the received signal at the BS averaged over all UEs in dB. For example, if there are two UEs, one with an SNR of 20 dB and the other one with 10 dB, then the average SNR is 15 dB.

SNR spread or pathloss spread

The SNR or pathloss spread is used to mimic the differences in the average receive power of the UEs due to large-scale fading. If the average SNR at the BS is 15 dB and the SNR spread is ± 5 dB, then this indicates that the SNR of the UEs at the BS is uniformly distributed in the range of [10, 20] dB. For such an example, the strongest and weakest UEs can have an SNR gap of up to 10 dB at the BS.

Average BLER

It denotes the BLER of the transmissions averaged over all UEs. This metric is beneficial in order to gauge the overall performance of the system. If the average BLER curve saturates and does not decline as the average SNR increases, then this can be an indicator of the system suffering from non-resolvable multi-user interference.

¹The simulator can be found here (available for free under an academic use license): <https://www.nt.tuwien.ac.at/research/mobile-communications/vccs/vienna-5g-simulators/>.

2.3. Simulation Setup

Correctly detected/decoded UEs

This denotes the number of UEs that have their data successfully detected. In the context of link-level simulations, this is the number of UEs that pass the CRC. For outage-based simulations, this is the number of UEs with SINRs exceeding the outage threshold.

3

Transmit-Side Optimization: NOMA Codebook Design

The design of the spreading codebook (or transmit signatures) in code-domain NOMA not only impacts the detection performance at the BS, but also determines the type of receiver employed. Certain codebook designs permit the use of low-complexity detection algorithms, which can be advantageous when the transmission occupies many RBs, and when a large number of UEs access the resources at the same time. Moreover, as the UEs generally experience different channel propagation conditions, a fair and robust design is required in order to cope with the lack of channel state information (CSI)-based adaptation, which is the case in grant-free access schemes.

In this chapter, we consider the codebook design problem, focusing on dense-spreading signatures, i.e., the sequences are not sparse. Specifically, in the first part, we consider designing codebooks according to a Grassmannian criterion, which is able to tackle the aforementioned issues of performance and complexity. We show a new iterative construction algorithm that is capable of producing codebooks with close-to-optimal correlation properties, while achieving a fast convergence rate compared to other algorithms. In the second part, we consider the problem of jointly designing codebooks across multiple groups of UEs, with the goal of optimizing their cross-group interference. This could be in the form of jointly designed codebooks for neighbouring cells, or for spatial clusters in massive MIMO. The algorithms developed in this chapter are published in [33, 34].

3.1 The Codebook Design Problem

Let us consider the received signal at the BS (with a single antenna) given by

$$\mathbf{y} = \sum_{k=1}^K \sqrt{LP_k} h_k \mathbf{s}_k x_k + \mathbf{n}, \quad (3.1)$$

3.1. The Codebook Design Problem

where the data-block index i was dropped for simplicity. It is clear that the detection performance of such a system depends on the design of the spreading signatures \mathbf{s}_k . The best-case scenario is to have an orthogonal set of signatures, i.e., $\mathbf{s}_l^H \mathbf{s}_k = 0, \forall l \neq k$. In this case, the UEs do not interfere with each other, and the detection can be performed using a low-complexity matched filter (MF). For our NOMA system, the number of UEs typically exceeds the spreading length, i.e., $K > L$, for which we have to deal with non-orthogonal signature sets. Moreover, since the UEs are received with different powers, due to uncorrelated fading realizations, pathloss differences, and possibly different transmit powers, the optimal codebook construction needs to take these constraints into account. This is important, since for the NOMA detection, the filtering is combined with IC. For example, if two UEs are received with equal power, then the BS should assign them near-orthogonal signatures, allowing to eliminate the interference between them via filtering. On the other hand, UEs with a large power difference should be assigned near-collinear signatures, since the interference between these can be managed via IC by first detecting the stronger UE, canceling it, and then detecting the weaker one. In our system, the number of UEs is much higher than just two, and therefore the optimization considering those aspects have to be done jointly over a larger set.

Throughout this work, we assume that the BS employs a fixed codebook, that is only constructed once, and we assume that the BS is not capable of optimizing the signature assignment, but rather the signatures are assigned randomly to the UEs. On the one hand, this has the implication that the BS does not have to construct codebooks and optimize them in an online-fashion, which reduces the overhead and operation complexity. On the other hand, such an assumption is more realistic when lacking CSI at the BS, which is the typical case in more general grant-free systems. Under such a lack of knowledge regarding the channel conditions of the UEs, it becomes natural to go for a robust min-max approach for the signature design. In other words, we design the codebook according to

$$\mathbf{S}_{\text{robust}} = \arg \min_{\mathbf{S}} \max_{\substack{\mathbf{s}_l, \mathbf{s}_k \in \mathbf{S} \\ \forall l \neq k}} |\mathbf{s}_l^H \mathbf{s}_k|. \quad (3.2)$$

That is, every signature in the codebook is as uncorrelated as possible, or as far away as possible from the other signatures, which at certain dimensionality results in a perfect equiangular separation of the signatures. Such a design criterion ensures fairness across the UEs, since each UE would experience similar level of interference from the other UEs, and therefore no advantage is given to a certain UE in terms of the detection performance.

The design criterion in (3.2) is known as the Grassmannian criterion, and the resulting codebook is called a Grassmannian codebook, as it is directly related to the Grassmannian line-packing problem [43] and frame theory [44], from which such codebooks get their name as Grassmannian frames. An important subset of these frames (or codebooks) are equiangular tight frames (ETFs), or better known in the telecommunications field as Welch-bound-equality (WBE) sequences. As the name

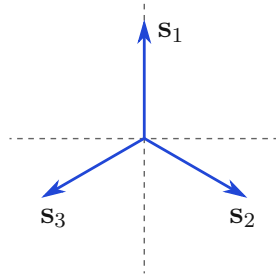


Figure 3.1: An ETF for $K = 3$ over \mathbb{R}^2 , showing the maximum angular separation between the signatures.

suggests, in addition to being a minimizer of (3.2), the resulting signatures have equal cross-correlation, i.e., $|\mathbf{s}_l^H \mathbf{s}_k| = \mu_{\text{Welch}}, \forall l \neq k$, where μ_{Welch} is the Welch bound [45]. Figure 3.1 shows an example ETF for $K = 3$ over \mathbb{R}^2 . ETFs can be maximizers of the UEs' SINR under MF [46, 47], and therefore under relatively low overloading factor K/L , they might allow the MF to provide a sufficient SINR for decodability, without having to rely on the more complex MMSE filter. Unfortunately, ETFs only exist for certain combinations of K and L , and therefore the more general Grassmannian criterion in (3.2) provides an extension. Note that such a signature design has been considered for previous CDMA systems, as in [48].

3.2 Construction of Grassmannian Codebooks

The construction of Grassmannian codebooks is difficult in general, and although there exist closed-form constructions, such as [43, 49–52], those are usually restricted to certain combinations of the space dimension L and the number of packed vectors K . Numerical methods that iteratively minimize the codebook maximum cross-correlation in (3.2), known as coherence, come in as a natural alternative. In [51], the authors view the problem as sphere vector quantization and employ a generalized Lloyd algorithm. The alternating projection (AP) algorithm, which iterates between certain spectral and structural constraints, has been used in [53, 54]. An algorithm based on the search for best complex antipodal spherical codes was proposed in [55], which maximizes distances between charged particles. A faster variant of that algorithm, known as the coherence-based Grassmannian codebook (CBGC) algorithm followed shortly in [56]. We present next an iterative algorithm based on a collision-based packing of equal-radius hyperspheres on the surface of a unit-norm hypersphere. As shown by the results, the algorithm is capable of producing codebooks with very low coherence levels, obtained together at a fast convergence.

3.2.1 Preliminaries

Grassmannian Frames

Let $\mathbf{S} = \{\mathbf{s}_k\}_{k=1}^K$ be the set of K vectors lying in the L -dimensional Hilbert space \mathbb{C}^L . This set is called a frame for \mathbb{C}^L , if there exist positive constants A and B such that

$$A\|\mathbf{c}\|^2 \leq \sum_{k=1}^K |\mathbf{c}^H \mathbf{s}_k|^2 \leq B\|\mathbf{c}\|^2, \quad (3.3)$$

for all $\mathbf{c} \in \mathbb{C}^L$, and where $\|\cdot\|$ denotes the Euclidean norm. Given a unit-norm frame (i.e., $\|\mathbf{s}_k\| = 1, \forall k$), the coherence is defined as

$$\mu(\mathbf{S}) = \max_{\substack{\mathbf{s}_l, \mathbf{s}_k \in \mathbf{S} \\ \forall l \neq k}} |\mathbf{s}_l^H \mathbf{s}_k|. \quad (3.4)$$

In other words, it is the maximum cross-correlation between any two distinct vectors in the set. The set \mathbf{S} is called a Grassmannian frame, if it is a minimizer of the coherence, i.e.,

$$\mathbf{S}_{\text{Grass.}} = \arg \min_{\mathbf{S} \in \mathbb{C}^{L \times K}} \mu(\mathbf{S}). \quad (3.5)$$

Unfortunately, except for certain combinations of L and K and depending on the space being \mathbb{R}^L or \mathbb{C}^L , the minimum coherence is generally unknown.

Lower Bounds on the Coherence

Under certain conditions, it is possible to derive lower bounds on the coherence. In the following, we list some of them for the complex case.

Theorem 3.1 (Welch bound [44, 45]). Let $\mathbf{S} \in \mathbb{C}^{L \times K}$ be a unit-norm frame, then

$$\mu(\mathbf{S}) \geq \mu_{\text{Welch}}(L, K) = \sqrt{\frac{K-L}{L(K-1)}}. \quad (3.6)$$

Equality can be achieved for the range $K \leq L^2$ when the frame \mathbf{S} is tight (i.e., $\mathbf{S}\mathbf{S}^H = \frac{K}{L}\mathbf{I}$) and for each $l \neq k$ we have

$$|\mathbf{s}_l^H \mathbf{s}_k| = \sqrt{\frac{K-L}{L(K-1)}}. \quad (3.7)$$

A frame (or codebook) with such properties is known as an ETF. It follows that all ETFs are Grassmannian, since they are minimizers of the coherence.

Theorem 3.2 (Orthoplex bound [50, 57]). Let $\mathbf{S} \in \mathbb{C}^{L \times K}$ be a unit-norm frame,

then

$$\mu(\mathbf{S}) \geq \sqrt{\frac{1}{L}}. \quad (3.8)$$

The bound holds for the range $L^2 < K \leq 2(L^2 - 1)$. For large K , the bounds by [58]

$$\mu(\mathbf{S}) \geq \sqrt{\frac{2K - L^2 - L}{(L + 1)(K - L)}}, \quad (3.9)$$

and [51, 59]

$$\mu(\mathbf{S}) \geq 1 - 2K^{-\frac{1}{L-1}}, \quad (3.10)$$

can be adopted. Putting these bounds together, the following composite bound is obtained in a fashion similar to [51, 55]

$$\mu(\mathbf{S}) \geq \begin{cases} \text{for } K \leq L^2 : \\ \sqrt{\frac{K - L}{L(K - 1)}}, \\ \text{for } L^2 < K \leq 2(L^2 - 1) : \\ \max \left\{ \sqrt{\frac{1}{L}}, \sqrt{\frac{2K - L^2 - L}{(L + 1)(K - L)}}, 1 - 2K^{-\frac{1}{L-1}} \right\}, \\ \text{for } K > 2(L^2 - 1) : \\ \max \left\{ \sqrt{\frac{2K - L^2 - L}{(L + 1)(K - L)}}, 1 - 2K^{-\frac{1}{L-1}} \right\} \end{cases} \quad (3.11)$$

$$= \mu_{\text{bound}}(L, K).$$

We will use this bound when we present our results later on.

3.2.2 Minimizing Coherence by Euclidean Distance Maximization

As pointed out in previous works such as [55, 60], the coherence can be minimized by performing Euclidean distance maximization between the vectors of the codebook. In the following, we identify key points for such an approach.

Consider the two vectors $\mathbf{a}, \mathbf{b} \in \mathbb{R}^L$ constrained to the surface of a unit-norm hypersphere (i.e., they are unit-norm vectors). How should they be placed (or packed) on that surface in such a way that the magnitude of their inner product $|\mathbf{b}^T \mathbf{a}|$ is minimized? The quantity $|\mathbf{b}^T \mathbf{a}|$ is related to the Euclidean distance between the two vectors. Therefore, one can attempt to pack the vectors in such a way that they become as far apart as possible in the Euclidean sense. However, going for a pure distance maximization will result in the two vectors being packed in an an-

3.2. Construction of Grassmannian Codebooks

tipodal fashion. Obviously, this is not our goal, since the two vectors would be just reflections of each other, i.e., $|\mathbf{b}^T \mathbf{a}| = |(-\mathbf{a})^T \mathbf{a}| = 1$. This problem can be tackled by including the reflection of the vector into the distance maximization. In other words, the vector \mathbf{a} needs to maximize its distance not only against \mathbf{b} , but also against its reflection $-\mathbf{b}$. By doing so, it is possible to obtain an orthogonal configuration leading to $|\mathbf{b}^T \mathbf{a}| = 0$.

When the two vectors are complex-valued, then not only the vector and its reflection has the same magnitude of the inner product to the other vector, but rather all of its complex-plane rotations. This can be easily seen through

$$|(e^{j\phi} \mathbf{b})^H \mathbf{a}| = |\mathbf{b}^H \mathbf{a}|, \quad (3.12)$$

for any $\phi \in [0, 2\pi]$, which includes the reflection as the special case of $\phi = \pi$. For \mathbf{a} to minimize $|\mathbf{b}^H \mathbf{a}|$, it has to stay away not only from \mathbf{b} , but also from all its complex-plane rotations $e^{j\phi} \mathbf{b}$. Our approach to this problem is to consider only the rotation that is closest to \mathbf{a} . Then, if \mathbf{a} maximizes its distance to that point, it will automatically maximize its distance to all of the complex-plane rotations of \mathbf{b} .

Proposition 3.1. Given two unit-norm vectors $\mathbf{a}, \mathbf{b} \in \mathbb{C}^L$, the point \mathbf{p} obtained by a rotation of \mathbf{b} in its complex-plane which has the smallest Euclidean distance to \mathbf{a} , is given by

$$\mathbf{p} = \frac{\mathbf{b}^H \mathbf{a}}{|\mathbf{b}^H \mathbf{a}|} \mathbf{b}. \quad (3.13)$$

Proof. We seek to perform the following minimization

$$\min_{\phi \in [0, 2\pi]} \|e^{j\phi} \mathbf{b} - \mathbf{a}\|^2. \quad (3.14)$$

By expanding the squared distance, we get

$$\min_{\phi \in [0, 2\pi]} \left(\|\mathbf{b}\|^2 + \|\mathbf{a}\|^2 - 2\Re\{e^{-j\phi} \mathbf{b}^H \mathbf{a}\} \right). \quad (3.15)$$

Since \mathbf{a} and \mathbf{b} are unit-norm, the minimization is equivalent to the maximization of the third term

$$\max_{\phi \in [0, 2\pi]} \Re\{e^{-j\phi} \mathbf{b}^H \mathbf{a}\}. \quad (3.16)$$

This is maximized if the imaginary part of the input argument is zero, which is achieved when $e^{j\phi} = \mathbf{b}^H \mathbf{a} / |\mathbf{b}^H \mathbf{a}|$. \square

The generalization of such a packing approach to more than two vectors is straightforward: every vector needs to stay away from all other vectors and all of their complex-plane rotations.

3.2.3 Collision-Based Packing

Iterative Distance Maximization

How can we perform such a distance maximization? In this subsection, we describe a collision-based algorithm that iteratively maximizes the distances between the vectors under the constraint derived in the previous subsection. We start by giving a general description of the algorithm:

1. Generate randomly K unit-norm vectors of dimension L . Consequently, these vectors will be lying on the surface of a unit-norm hypersphere.
2. At each point (vector), generate a hypersphere of equal radius r .
3. Adjust r by a step γ_r .
4. In case of a collision (i.e., the hyperspheres overlapping), the hyperspheres repel each other. This is achieved by moving their corresponding center vectors. The amount of repulsion is proportional to how much they were colliding by.
5. Until a stopping criterion is met, repeat 3) and 4). As the vectors keep getting pushed away from each other, we soon reach a configuration satisfying (3.2), or at least close to it.

It is important to note that the collisions here are not only with respect to the vectors themselves, but also with respect to their whole complex-plane rotations, as pointed out in Section 3.2.2.

A collision between two hyperspheres with centers \mathbf{s}_k and \mathbf{s}_l occurs if

$$\|\mathbf{d}_{k,l}\| < 2r, \quad (3.17)$$

where $\mathbf{d}_{k,l}$ is defined as

$$\mathbf{d}_{k,l} = \mathbf{s}_k - \frac{\mathbf{s}_l^H \mathbf{s}_k}{|\mathbf{s}_l^H \mathbf{s}_k|} \mathbf{s}_l. \quad (3.18)$$

That is, we are checking the collision between \mathbf{s}_k and the rotation of \mathbf{s}_l in its complex-plane which is closest to \mathbf{s}_k . To clear out the collision, \mathbf{s}_k is pushed away from that rotation by the amount of collision using the movement vector

$$\mathbf{u}_k = (2r - \|\mathbf{d}_{k,l}\|) \frac{\mathbf{d}_{k,l}}{\|\mathbf{d}_{k,l}\|}. \quad (3.19)$$

However, it can happen that \mathbf{s}_k collides with more than one vector. Therefore, the movement vector \mathbf{u}_k is extended to

$$\mathbf{u}_k = \sum_{l \in S_k} (2r - \|\mathbf{d}_{k,l}\|) \frac{\mathbf{d}_{k,l}}{\|\mathbf{d}_{k,l}\|}, \quad (3.20)$$

3.2. Construction of Grassmannian Codebooks

where \mathcal{S}_k is the index-set of all hyperspheres that are colliding with \mathbf{s}_k (i.e., satisfying (3.17)). The centers of the hyperspheres are then adjusted according to

$$\mathbf{s}_k \leftarrow \frac{\mathbf{s}_k + \mathbf{u}_k}{\|\mathbf{s}_k + \mathbf{u}_k\|}. \quad (3.21)$$

The normalization is performed to make sure that the result lies again on the surface of the unit-norm hypersphere.

The Packing Radius

The packing (or collision) radius is directly related to the coherence of the codebook.

Proposition 3.2. Given two unit-norm vectors $\mathbf{s}_k, \mathbf{s}_l \in \mathbb{C}^L$ and $\mathbf{d}_{k,l}$ as defined in (3.18), then

$$\|\mathbf{d}_{k,l}\| = \sqrt{2(1 - |\mathbf{s}_l^H \mathbf{s}_k|)}. \quad (3.22)$$

Proof. Expand the squared distance

$$\|\mathbf{d}_{k,l}\|^2 = \|\mathbf{s}_k\|^2 + \left\| \frac{\mathbf{s}_l^H \mathbf{s}_k}{|\mathbf{s}_l^H \mathbf{s}_k|} \mathbf{s}_l \right\|^2 - 2\Re \left\{ \left(\frac{\mathbf{s}_l^H \mathbf{s}_k}{|\mathbf{s}_l^H \mathbf{s}_k|} \mathbf{s}_l \right)^H \mathbf{s}_k \right\}. \quad (3.23)$$

The first two terms belong to unit-norm vectors and therefore they are equal to one. For the third term, we have

$$\left(\frac{\mathbf{s}_l^H \mathbf{s}_k}{|\mathbf{s}_l^H \mathbf{s}_k|} \mathbf{s}_l \right)^H \mathbf{s}_k = \frac{(\mathbf{s}_l^H \mathbf{s}_k)^*}{|\mathbf{s}_l^H \mathbf{s}_k|} \mathbf{s}_l^H \mathbf{s}_k = |\mathbf{s}_l^H \mathbf{s}_k|. \quad (3.24)$$

Plugging the result in (3.23), we get

$$\|\mathbf{d}_{k,l}\|^2 = 2(1 - |\mathbf{s}_l^H \mathbf{s}_k|). \quad (3.25)$$

Finally, we apply the square root to both sides. \square

When the hyperspheres are barely touching each other (no collision yet), then the distance between them is $\|\mathbf{d}_{k,l}\| = 2r$. Using (3.22), we obtain

$$r = \sqrt{\frac{1}{2}(1 - |\mathbf{s}_l^H \mathbf{s}_k|)}. \quad (3.26)$$

We can now translate the lower bound of the coherence in (3.11) into an upper bound on the maximum possible packing radius, i.e.,

$$r_{\text{bound}}(L, K) = \sqrt{\frac{1}{2}(1 - \mu_{\text{bound}}(L, K))}. \quad (3.27)$$

Consequently, it is not possible to have K vectors spaced larger than $2r_{\text{bound}}$ from each other, as this is not achievable at all in the L -dimensional ambient space.

The ICBP Algorithm

The implementation of the proposed algorithm is shown below. We call it iterative collision-based packing (ICBP) ¹. The starting points $\mathbf{S}_{\text{initial}}$ are random complex-valued unit-norm Gaussian vectors. The algorithm operates by setting a target packing radius, and attempts to approach it iteratively. The initial target radius is set equal to r_{bound} , which is the largest possible. Collision clearance is then performed according to (3.21). When the packing radius at the current iteration r_{current} is worse than the packing radius achieved in the previous iterations r_{best} , then we decrease the target radius by a step size γ_r , because it might be that the original target radius is not achievable in \mathbb{C}^L . If it improves on the next iterations, then we increase it back again. This way, we approach the best possible packing in the ambient space. In certain cases, the algorithm might exhibit an oscillatory behavior. To address that, a damping factor β is introduced in the update of (3.21). In our test cases, we found that $\beta = 0.8$ yielded good results most of the time, and therefore we adopt it here.

Algorithm 1: Iterative Collision-Based Packing (ICBP)

```

input :  $L, K, \mathbf{S}_{\text{initial}}, \gamma_r, I_{\text{max}}$ 
output: The codebook  $\mathbf{S}$ 
 $\mathbf{S} \leftarrow \mathbf{S}_{\text{initial}}, r \leftarrow r_{\text{bound}}(L, K), r_{\text{best}} \leftarrow 0, \beta \leftarrow 0.8$ 
for  $i \leftarrow 1$  to  $I_{\text{max}}$  do
     $r_{\text{current}} \leftarrow \sqrt{0.5 (1 - \mu(\mathbf{S}))}$ 
    if  $r_{\text{current}} < r_{\text{best}}$  then
         $r \leftarrow r - \gamma_r$ 
    else
         $r_{\text{best}} \leftarrow r_{\text{current}}$ 
         $r \leftarrow \min\{r + \gamma_r, r_{\text{bound}}(L, K)\}$ 
    for  $k \leftarrow 1$  to  $K$  do
         $\mathbf{u}_k \leftarrow \sum_{l \in \mathcal{S}_k} (2r - \|\mathbf{d}_{k,l}\|) \frac{\mathbf{d}_{k,l}}{\|\mathbf{d}_{k,l}\|}$ 
         $\mathbf{s}_k \leftarrow \frac{\mathbf{s}_k + \beta \mathbf{u}_k}{\|\mathbf{s}_k + \beta \mathbf{u}_k\|}$ 
return  $\mathbf{S}$ 

```

Although our discussion was focused on the complex space, the algorithm can be used to generate real-valued codebooks as well. This depends on whether the initial set $\mathbf{S}_{\text{initial}}$ is real or complex-valued.

¹The algorithm is available for download at <https://www.nt.tuwien.ac.at/christian-doppler-laboratory/cd-download/>

3.2.4 Algorithm Performance

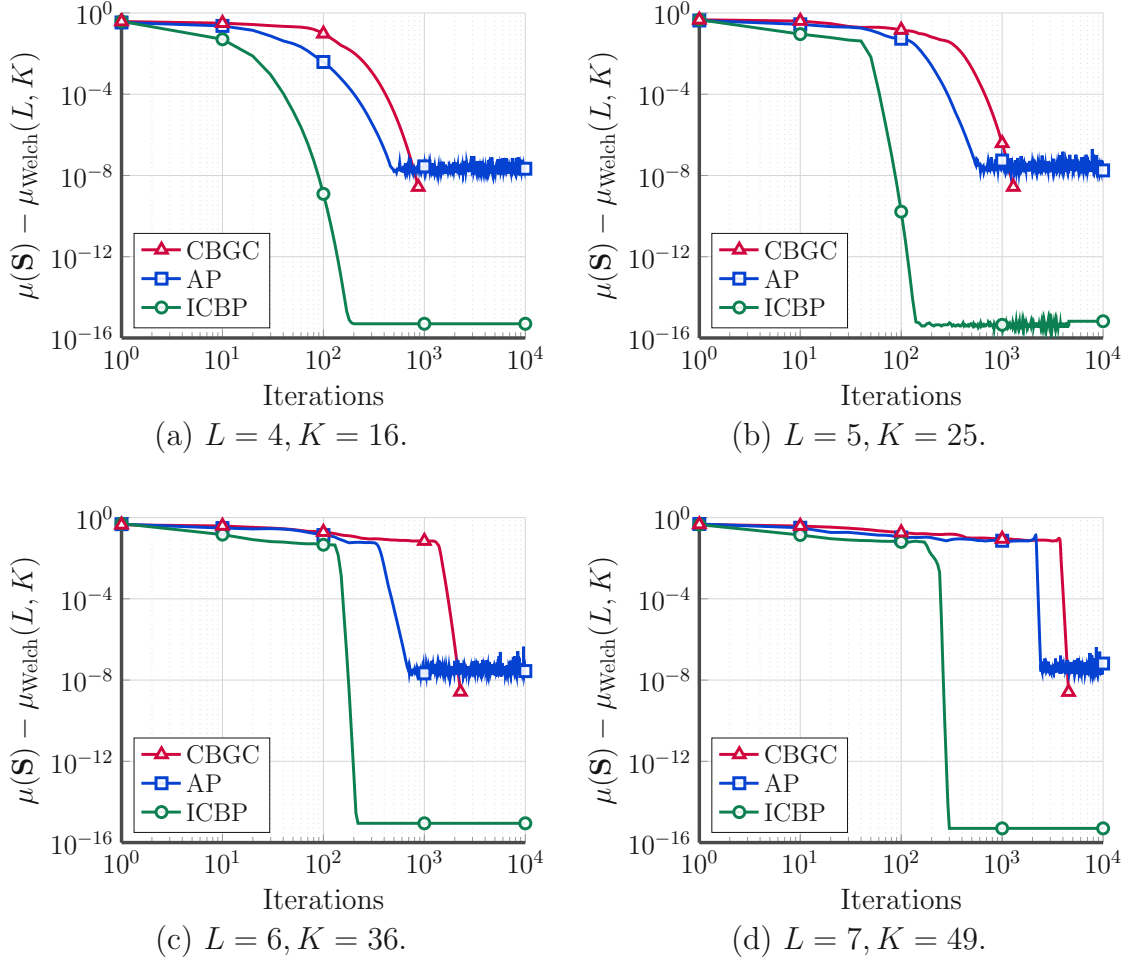
Convergence Performance and Achieved Coherence

We investigate the convergence performance of our algorithm and the achieved coherence of the constructed codebooks. As a reference, we also show the results obtained by the CBGC and AP algorithms. The CBGC algorithm shares two aspects with our algorithm. First, it is in the category of distance maximization, and second, the update step of their algorithm, although derived in a different way, uses a similar approach as the one we used via Proposition 3.1. The AP algorithm on the other hand, operates on the Gramian of the codebook, and is more flexible when it comes to enforcing additional constraints on the vectors.

In our tests, it became clear that the starting set $\mathbf{S}_{\text{initial}}$ has a significant impact not only on the final obtained coherence, but also on the convergence behavior of the algorithms. To address that, we ran all three algorithms ten times, in each run the starting points were generated randomly. We then took the best result out of those runs. For all of our tests in this subsection, a step-size of $\gamma_r = 10^{-4}$ is used for the ICBP algorithm.

First, we attempt to construct Grassmannian codebooks with $K = L^2$ for $L = \{4, 5, 6, 7\}$. For such a combination of L and K , an ETF exists, and therefore the minimum coherence is equal to the Welch bound (refer to [61] for tables on existing ETFs). In Figure 3.2 we plot the difference to the Welch bound versus the number of iterations. The reason why the CBGC curve stops suddenly is because the algorithm uses a convergence parameter of 10^{-10} , which we did not change, as it affects the performance within the sub-problems it tries to solve. Our algorithm managed not only to hit the bound with high accuracy, but also achieved it with a much faster convergence compared to the others. For the $L \times K = 5 \times 25$ configuration, we measured the time each algorithm takes to hit an accuracy higher than 10^{-7} . The recorded run-times in seconds were 0.04, 0.38, 0.15 for the ICBP, CBGC, and AP algorithms, respectively. This demonstrates the time saving that our algorithm can offer, by requiring less iterations.

Next, we aim to construct Grassmannian codebooks at arbitrary combinations of L and K . The coherence results for $L = \{4, 8, 12\}$ are shown in Figure 3.3, where we plot the difference to the composite bound of (3.11). The CBGC algorithm, in general, requires significantly more iterations to achieve low coherence levels compared to our algorithm. In the shown results, our algorithm achieves very close coherence levels to the CBGC algorithm with a maximum number of iterations I_{max} that is much less. Increasing the number of iterations for our algorithm did yield better results. However, the CBGC algorithm kept outperforming it by a slight margin when its number of iterations is set to 10^6 . The AP algorithm showed the worse performance compared to the other algorithms, and since it did not show a considerable improvement beyond 10^4 iterations, we did not show the result for a higher number of iterations. The reason for the worse performance is possibly due to the lack of a target coherence adjustment during the iterations.

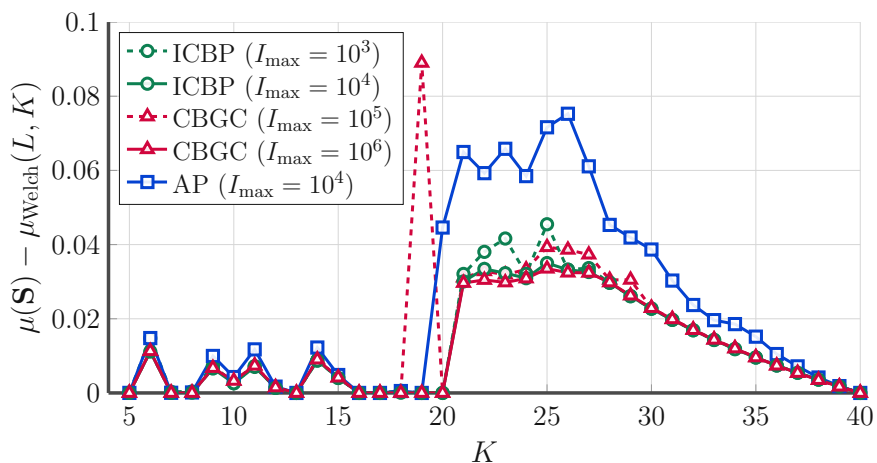
Figure 3.2: Convergence of the algorithms for the construction of ETFs with $K = L^2$.

The algorithms seem to struggle at large L and K . For example, they were not able to locate the ETFs of 8×64 and 12×45 . Nonetheless, the obtained results across all combinations (including ETFs) are close to the lower bound.

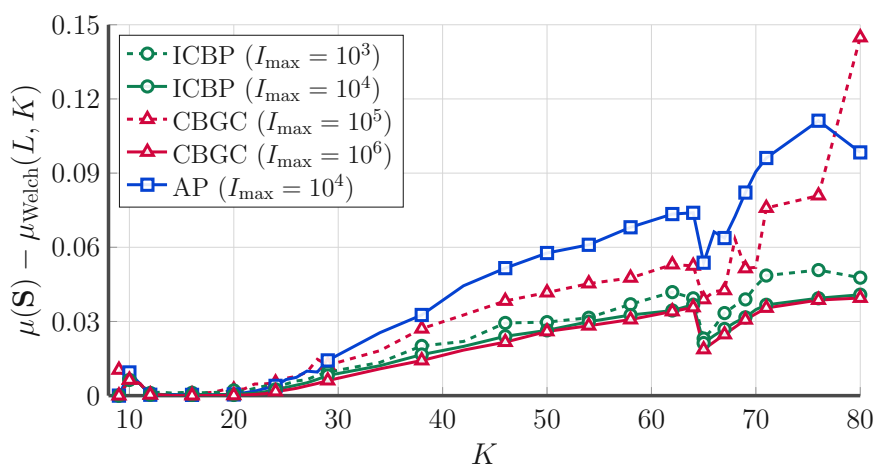
Detection Performance under MF and MMSE Filtering

Next, we evaluate the detection performance of a NOMA uplink utilizing those signatures. We consider a setup consisting of $K = 8$ UEs, transmitting via spreading signatures of length $L = 4$. All the UEs transmit with the same power over Rayleigh channels, using 4-QAM and turbo-coding with a code-rate of $1/2$. Figure 3.4 shows the average BLER of all the UEs versus their average SNR at the BS. We compare the detection performance using Grassmannian and randomly constructed codebooks under MF with IC against the performance under MMSE filtering with IC. It can be observed that the MMSE filter is able to provide the necessary interference suppression capabilities for both codebooks, with only a small performance gap

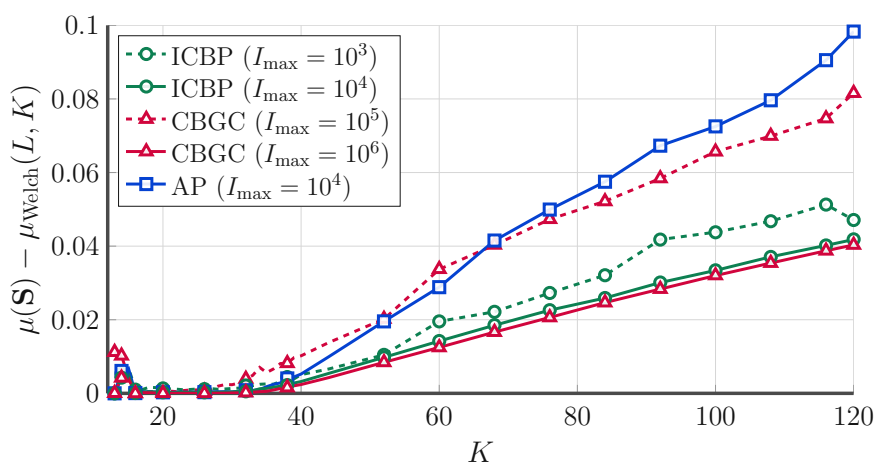
3.2. Construction of Grassmannian Codebooks



(a) $L = 4$.



(b) $L = 8$.



(c) $L = 12$.

Figure 3.3: Coherence results of the constructed codebooks for arbitrary L and K , with different maximum number of iterations.

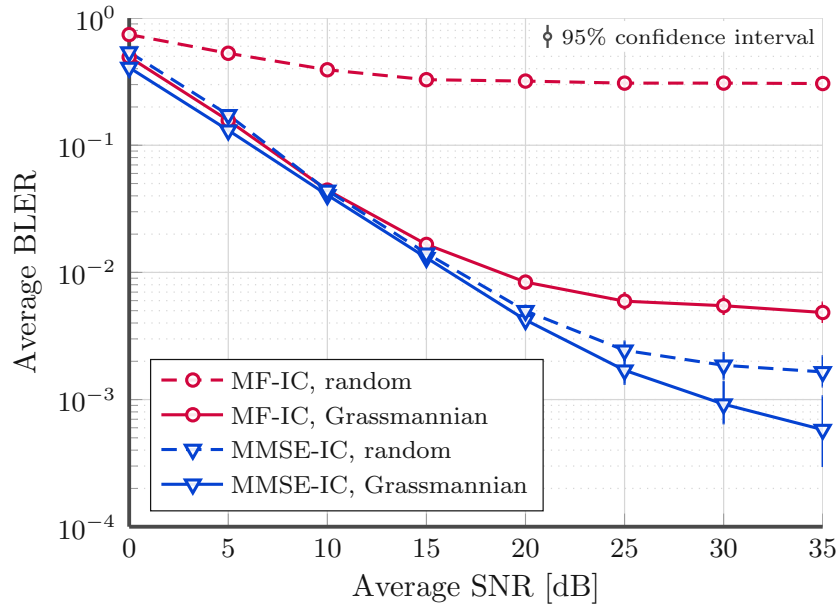
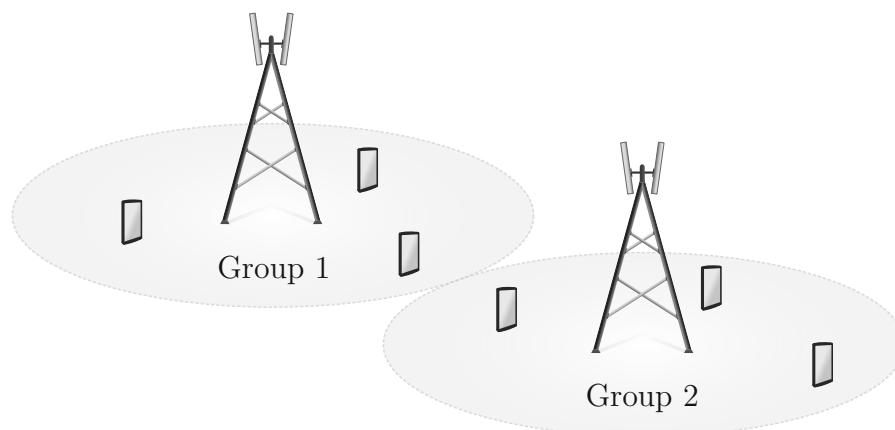


Figure 3.4: Codebook performance under MF and MMSE detection.

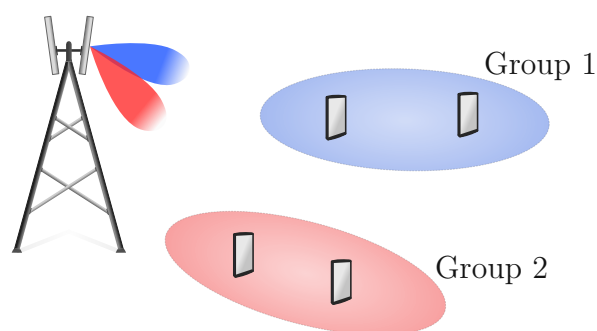
between them. As for the MF, a substantial improvement can be observed for the transmission utilizing the Grassmannian codebook. This can help in reducing the detection complexity by replacing the MMSE filter with a MF, especially when the transmission is over many RBs, which would require the calculation of a large number of MMSE weights.

3.3 Multi-Group Joint Codebook Design

So far, we have considered the codebook design problem for a single group of UEs. However, the UEs are usually available in multiple groups. One example is cellular deployments as depicted in Figure 3.5a, where each cell has a group of UEs, and each BS only attempts to detect the transmissions from its own cell, while transmissions from neighbouring cells cause inter-cell interference. The question here is, can we do better by jointly designing the codebooks across the two cells, instead of reusing the same codebook in each of them? This is especially important for cell-edge UEs, as these would suffer the most from inter-cell interference. Another example is the availability of spatial clusters, where a massive MIMO BS would be capable of serving these clusters by forming beams towards them, as shown in Figure 3.5b. Since residual inter-cluster interference might be present after beamforming, the question is whether a joint design across the clusters can increase the robustness to the inter-cluster interference. Therefore, in the remainder of this chapter, we will consider the problem of designing codebooks that are optimized jointly across multiple user groups.



(a) Multi-cell deployment.



(b) Spatial clusters.

Figure 3.5: Example multi-group setups.

3.3.1 Cross-Codebook Optimization

Instead of reusing the codebook \mathbf{S} across the interfering groups, we use different codebooks that have the same correlation properties as of \mathbf{S} , however, they are designed jointly, in an attempt to reduce the impact of one codebook on another. The reason for requiring them to have the same correlation properties is because we do not want to reduce the inter-group interference at the cost of an increased intra-group interference, and this is achieved by preserving the internal correlation structure of the codebooks. The correlation properties of the codebook are fully captured by the Gramian $\mathbf{S}^H \mathbf{S}$. Therefore, we first have to determine what freedom do we have in designing these codebooks, given that their Gramian is the same. The answer to that is, codebooks with the same Gramian are equivalent up to an isometry. In our case, since all the signatures have unit-norm, then the isometry is just a rotation (including reflections). To understand why this is true, consider the

3.3. Multi-Group Joint Codebook Design

two codebooks \mathbf{S}_1 and \mathbf{S}_2 with identical Gramian, i.e.,

$$\mathbf{S}_1^H \mathbf{S}_1 = \mathbf{S}_2^H \mathbf{S}_2. \quad (3.28)$$

We can introduce two unitary matrices $\mathbf{U}_1, \mathbf{U}_2 \in \mathbb{C}^{L \times L}$ ($\mathbf{U}_1^H \mathbf{U}_1 = \mathbf{U}_2^H \mathbf{U}_2 = \mathbf{I}_L$) without altering the equality

$$\begin{aligned} \mathbf{S}_1^H \mathbf{U}_1^H \mathbf{U}_1 \mathbf{S}_1 &= \mathbf{S}_2^H \mathbf{U}_2^H \mathbf{U}_2 \mathbf{S}_2 \\ (\mathbf{U}_1 \mathbf{S}_1)^H (\mathbf{U}_1 \mathbf{S}_1) &= (\mathbf{U}_2 \mathbf{S}_2)^H (\mathbf{U}_2 \mathbf{S}_2), \\ \Rightarrow \mathbf{U}_1 \mathbf{S}_1 &= \mathbf{U}_2 \mathbf{S}_2, \end{aligned} \quad (3.29)$$

from which follows that

$$\mathbf{S}_2 = \mathbf{U}_2^{-1} \mathbf{U}_1 \mathbf{S}_1. \quad (3.30)$$

The quantity $\mathbf{U}_2^{-1} \mathbf{U}_1$ is just another unitary matrix, and therefore the two codebooks are related by a unitary transformation. In other words, \mathbf{S}_2 is a rotation (including reflections) of \mathbf{S}_1 .

Next, we need to answer two questions; what determines a good rotation, is there a metric for it? And then, given the target metric, how to perform such a rotation in the first place? Ultimately, we would like to have $\mathbf{S}_1^H \mathbf{S}_2 = \mathbf{0}_{K \times K}$, where $\mathbf{0}_{K \times K}$ is the all-zeros matrix of size $K \times K$; that is, the codebooks are orthogonal and do not interfere with each other. However, in our case, this is not possible at all, because we usually have $K \geq L$, and according to (3.2), we design the signatures such that they are as far apart as possible in the ambient space \mathbb{C}^L . Therefore, we end up having $\text{span}\{\mathbf{S}_1\} = \text{span}\{\mathbf{S}_2\} = \mathbb{C}^L$. Since the codebooks cannot be orthogonal with respect to each other, we turn our attention to getting close to orthogonality by means of some metric $\|\mathbf{S}_1^H \mathbf{S}_2\|$. Metrics such as the Frobenius or spectral norm cannot be used here. To see why, consider the trace definition of the Frobenius norm

$$\|\mathbf{S}_1^H \mathbf{S}_2\|_F^2 = \text{tr}(\mathbf{S}_1^H \mathbf{S}_2 \mathbf{S}_2^H \mathbf{S}_1). \quad (3.31)$$

As mentioned before, a class of optimal solutions to (3.2) are ETFs satisfying $\mathbf{S} \mathbf{S}^H = \frac{K}{L} \mathbf{I}_L$. Therefore, if our codebooks are ETFs, then $\mathbf{S}_1 \mathbf{S}_1^H = \mathbf{S}_2 \mathbf{S}_2^H = \frac{K}{L} \mathbf{I}_L$, and thus

$$\|\mathbf{S}_1^H \mathbf{S}_2\|_F^2 = \frac{K}{L} \text{tr}(\mathbf{S}_1^H \mathbf{S}_1) = \frac{K}{L} \|\mathbf{S}_1\|_F^2. \quad (3.32)$$

This is a constant that does not depend on how \mathbf{S}_2 is rotated with respect to \mathbf{S}_1 . In a similar fashion, using the eigenvalue definition of the spectral norm, we can show that it also takes a constant value for such codebooks. Moreover, following our assumptions in Section 3.1, for which the BSs do not adapt their codebooks in an online-fashion, our design metric should provide a robust and fair criterion.

3.3. Multi-Group Joint Codebook Design

Consider the element-wise maximum norm defined as

$$\|\mathbf{A}\|_{\max} = \max_{k,l} |[\mathbf{A}]_{kl}|, \quad (3.33)$$

where $[\mathbf{A}]_{kl}$ is the element at the k^{th} row and l^{th} column of the matrix \mathbf{A} . By applying it to our problem, we obtain

$$\|\mathbf{S}_1^H \mathbf{S}_2\|_{\max} = \max_{k,l} |[\mathbf{S}_1^H \mathbf{S}_2]_{kl}| = \max_{\mathbf{a} \in \mathbf{S}_1, \mathbf{b} \in \mathbf{S}_2} |\mathbf{a}^H \mathbf{b}|. \quad (3.34)$$

Minimizing this Grassmannian-like metric would then guarantee that a certain separation (or angle) between the codebooks is maintained. In other words, for two UEs belonging to two different groups and lying at the edge, they would never get to use the same transmit signature, thus reducing their interference to each other. Therefore, we adopt this metric in our joint design approach. Note that (3.34) is a min-max problem that is similar to the Grassmannian design problem, which is difficult to solve. Instead, we seek to bring the maximum cross-correlation between the codebooks below a certain level μ .

Let the number of codebooks (groups) be J ; the goal is to find codebooks $\mathbf{S}_1, \mathbf{S}_2, \dots, \mathbf{S}_J \in \mathbb{C}^{L \times K}$ with the following conditions

$$\begin{aligned} \mathbf{S}_i^H \mathbf{S}_i &= \mathbf{S}^H \mathbf{S}, \quad \forall i, \\ \|\mathbf{S}_i^H \mathbf{S}_j\|_{\max} &\leq \mu, \quad \forall i \neq j. \end{aligned} \quad (3.35)$$

The first condition defines the internal structure of the codebooks, e.g., according to (3.2), while the second condition enforces the cross-correlation between the signatures of the different codebooks to go below a specific level μ . For two $K = 3$ ETFs over \mathbb{R}^2 , this design approach is illustrated in Figure 3.6. Note that the cross-correlation level μ cannot be arbitrarily small, but it is rather limited by the maximum packing possible in the ambient space. We are unaware of lower bounds to the packings of codebooks with non-orthogonal vectors, but μ is certainly larger than the Grassmannian bounds of Section 3.2.1, such as the Welch bound, and it cannot take values larger than 1. Our choice of μ is rather experimental, as explained later in the results subsection.

The construction of such rotated codebooks can be performed using the iterative algorithm of alternating projection [62, 63], in a fashion similar to the problem of subspace packing on the Grassmannian manifold [64]. In our case, we do not have subspaces, but rather codebooks that span the whole ambient space with generally non-orthogonal vectors (signatures), and we pack those codebooks with respect to the $\|\cdot\|_{\max}$ norm.

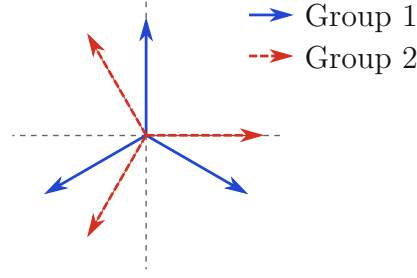


Figure 3.6: Optimization of two $K = 3$ ETFs over \mathbb{R}^2 according to their maximum cross-correlation. Notice how the internal correlation structure of the codebooks is preserved.

3.3.2 Alternating Projection

Let $\mathbf{\Sigma} = [\mathbf{S}_1 \ \mathbf{S}_2 \ \dots \ \mathbf{S}_J]$ be the matrix containing the codebooks. The Gramian of the codebooks is given by

$$\mathbf{G} = \mathbf{\Sigma}^H \mathbf{\Sigma} = \begin{bmatrix} \mathbf{S}_1^H \mathbf{S}_1 & \mathbf{S}_1^H \mathbf{S}_2 & \dots & \mathbf{S}_1^H \mathbf{S}_J \\ \mathbf{S}_2^H \mathbf{S}_1 & \mathbf{S}_2^H \mathbf{S}_2 & \dots & \mathbf{S}_2^H \mathbf{S}_J \\ \vdots & \dots & \ddots & \vdots \\ \mathbf{S}_J^H \mathbf{S}_1 & \mathbf{S}_J^H \mathbf{S}_2 & \dots & \mathbf{S}_J^H \mathbf{S}_J \end{bmatrix}. \quad (3.36)$$

The properties of the Gramian matrix \mathbf{G} meeting the conditions in (3.35) are

- \mathbf{G} is Hermitian.
- The diagonal blocks satisfy $\mathbf{G}_{ii} = \mathbf{S}^H \mathbf{S}$.
- Every off-diagonal block satisfies $\|\mathbf{G}_{ij}\|_{\max} \leq \mu$.
- \mathbf{G} is positive semi-definite.
- \mathbf{G} has a rank of L .
- \mathbf{G} has a trace equal to KJ .

The first three properties are structural properties, while the last three are spectral. Define the structural constraints' set as

$$\mathcal{H} = \{\mathbf{H} \in \mathbb{C}^{KJ \times KJ} : \mathbf{H} = \mathbf{H}^H, \mathbf{H}_{ii} = \mathbf{S}^H \mathbf{S}, \|\mathbf{H}_{ij}\|_{\max} \leq \mu, \forall i \neq j\}, \quad (3.37)$$

and the spectral constraints' set as

$$\mathcal{G} = \{\mathbf{G} \in \mathbb{C}^{KJ \times KJ} : \mathbf{G} \succeq 0, \text{rank}(\mathbf{G}) = L, \text{tr}(\mathbf{G}) = KJ\}. \quad (3.38)$$

We use the alternating projection algorithm to find a matrix \mathbf{G} satisfying both constraint sets. Let the maximum number of iterations be T , the algorithm is summarized as follows

3.3. Multi-Group Joint Codebook Design

1. Start with a random Hermitian $\mathbf{G}^{(0)} \in \mathbb{C}^{KJ \times KJ}$.
2. Set the iteration number $t = 0$.
3. Solve a nearest matrix problem to the set \mathcal{H}

$$\mathbf{H}^{(t)} = \arg \min_{\mathbf{H} \in \mathcal{H}} \|\mathbf{H} - \mathbf{G}^{(t)}\|_F. \quad (3.39)$$

4. Solve a nearest matrix problem to the set \mathcal{G}

$$\mathbf{G}^{(t+1)} = \arg \min_{\mathbf{G} \in \mathcal{G}} \|\mathbf{G} - \mathbf{H}^{(t)}\|_F. \quad (3.40)$$

5. Break if $t = T$. Otherwise, increase t and go to 3.

Let the eigendecomposition of $\mathbf{G}^{(T+1)}$ be $\mathbf{U}\mathbf{\Lambda}\mathbf{U}^H$, the codebooks' matrix $\mathbf{\Sigma}$ is then given by the rows of $\mathbf{\Lambda}^{1/2}\mathbf{U}^H$ corresponding to the L largest eigenvalues.

Two nearest matrix problems need to be solved. Let us start with (3.39). From the constraints set \mathcal{H} , the diagonal blocks are forced to be $\mathbf{H}_{ii} = \mathbf{S}^H\mathbf{S}$. Therefore, the nearest matrix problem is then concerned with the off-diagonal blocks only ($i \neq j$)

$$\mathbf{H}_{ij} = \arg \min_{\mathbf{A} \in \mathbb{C}^{K \times K}} \|\mathbf{A} - \mathbf{G}_{ij}\|_F^2, \quad \|\mathbf{A}\|_{\max} \leq \mu. \quad (3.41)$$

Both the objective and constraint functions are convex. Therefore, a unique solution exists.

Proposition 3.3. For the optimization problem in (3.41), the solution is given by

$$[\mathbf{H}_{ij}]_{kl} = \begin{cases} [\mathbf{G}_{ij}]_{kl}, & [\mathbf{G}_{ij}]_{kl} \leq \mu, \\ \mu[\mathbf{G}_{ij}]_{kl} / |[\mathbf{G}_{ij}]_{kl}|, & [\mathbf{G}_{ij}]_{kl} > \mu. \end{cases} \quad (3.42)$$

That is, every element of the off-diagonal block $\mathbf{G}_{ij}^{(t)}$ which has its magnitude larger than μ , is scaled to have a magnitude of exactly μ .

Proof. The key observation here is that the constraint $\|\mathbf{A}\|_{\max} \leq \mu$ applies element-wise. It does not enforce a relationship across the entire structure of the matrix \mathbf{A} from the feasible set. Expand the Frobenius norm in (3.41)

$$\|\mathbf{A} - \mathbf{G}_{ij}\|_F^2 = \sum_k \sum_l |[\mathbf{A}]_{kl} - [\mathbf{G}_{ij}]_{kl}|^2. \quad (3.43)$$

Since there is no relationship between the elements, and since each of the sum terms is non-negative, then the minimum of $\|\mathbf{A} - \mathbf{G}_{ij}\|_F^2$ is obtained, when every sum term $|[\mathbf{A}]_{kl} - [\mathbf{G}_{ij}]_{kl}|^2$ is minimized individually. We can then reformulate (3.41) equivalently in terms of the elements

$$[\mathbf{H}_{ij}]_{kl} = \arg \min_{[\mathbf{A}]_{kl} \in \mathbb{C}} |[\mathbf{A}]_{kl} - [\mathbf{G}_{ij}]_{kl}|^2, \quad |[\mathbf{A}]_{kl}| \leq \mu, \quad (3.44)$$

3.3. Multi-Group Joint Codebook Design

from which directly follows that if $|\mathbf{G}_{ij}|_{kl} \leq \mu$, the solution is $[\mathbf{H}_{ij}]_{kl} = [\mathbf{G}_{ij}]_{kl}$. Otherwise, we look for a solution that satisfies the constraint and at the same time is the closest to $[\mathbf{G}_{ij}]_{kl}$, and that would be $[\mathbf{H}_{ij}]_{kl} = \mu[\mathbf{G}_{ij}]_{kl} / |[\mathbf{G}_{ij}]_{kl}|$. \square

As for the second nearest matrix problem in (3.40), it follows the same solution as in [64], which is solved by applying the Karush-Kuhn-Tucker (KKT) conditions [65].

3.3.3 Case Study: Multi-Cell Deployment

We investigate a scenario consisting of two interfering cells, with a spreading length of $L = 8$. Our focus is on the uplink performance of the cell-edge UE, as it is the one that suffers the most from inter-cell interference. First, we construct a NOMA codebook \mathbf{S} of dimensions 8×24 according to our construction algorithm in Section 3.2. From the Welch bound, we know that at such L and K , there is a possibility that a Grassmannian codebook exists with a coherence of 0.2949. We were able to construct a codebook with a coherence of 0.2972. Next, we apply our joint optimization method to find the rotated codebooks for the two cells. For the choice of μ , we start with a large value, say 0.9. If the algorithm succeeds in finding such a packing, then we reduce μ to 0.8. If it again succeeds, then we further reduce it, and so on (i.e., via a bisection method). The best packing we found was for $\mu = 0.52$. Going below that level caused the algorithm to fail in maintaining the first condition of (3.35). The reason for that is due to the limit on the maximum possible packing in the ambient space, which then would prevent finding a valid Gramian in the step of (3.40). Furthermore, we also performed the packing for an OMA system, by taking \mathbf{S} to be a unitary matrix of dimensions 8×8 . We then applied our method and obtained a best packing of the two codebooks for $\mu = 1/\sqrt{8}$.

In the primary cell, the SNR of the UEs at the BS is uniformly distributed in the range [4, 20] dB. This corresponds to the variation between the received power of the UEs due to their position within the cell. The cell-edge UE under consideration has its SNR fixed to 4 dB. The SNR of the interfering UEs from the interfering cell at the primary BS is distributed in the range $[-12, 4] - P_{\text{cell-edge}}/I_{\text{strongest}}$ dB, where $P_{\text{cell-edge}}/I_{\text{strongest}}$ is the ratio between the average received power of our cell-edge UE to the average received power of the strongest interferer. When this ratio is zero, then the strongest interferer from the interfering cell can be as strong as our cell-edge UE. All the UEs transmit using 4-QAM and turbo-coding with a code rate of 1/3, with Rayleigh-fading assumed. At the BS, an MMSE-IC receiver is employed.

In Figure 3.7, the BLER of the cell-edge UE at different interference levels is shown for a fully loaded OMA system; that is, $N_P = 8$ and $N_I = 8$, where N_P and N_I are the number of active UEs in the primary and interfering cells, respectively. The ‘reuse’ curve corresponds to the strategy where a single codebook is reused in each cell, while the ‘joint’ curve corresponds to the jointly designed codebooks obtained using our method. We observe that the jointly designed codebooks are able to sustain higher interference power, while providing the same BLER as the reused codebooks.

3.3. Multi-Group Joint Codebook Design

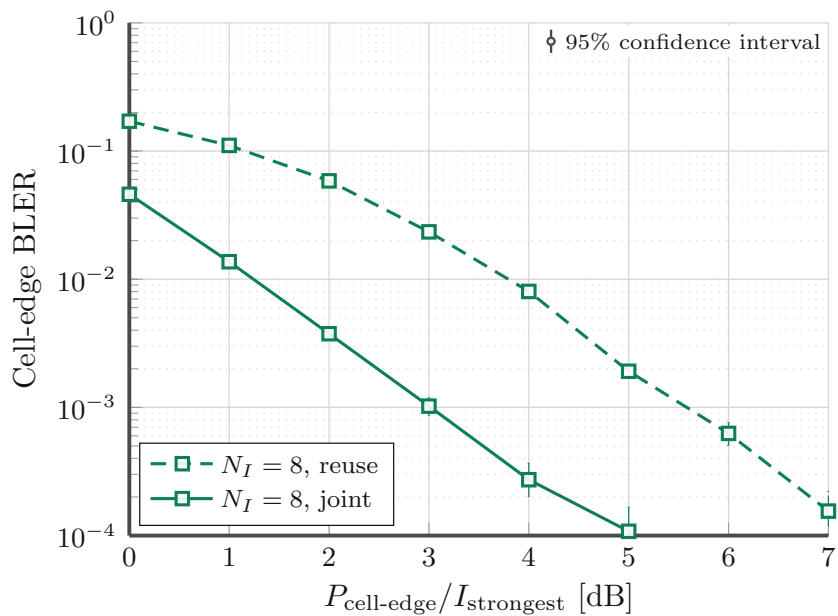


Figure 3.7: Performance of the OMA cell-edge UE for $N_P = 8$ (100%), and $N_I = 8$ (100%).

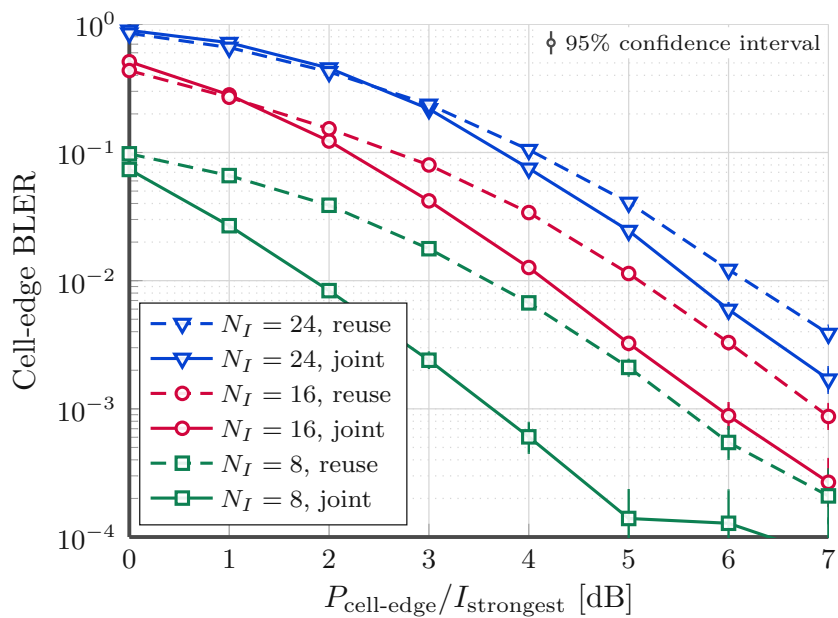


Figure 3.8: Performance of the NOMA cell-edge UE for $N_P = 24$ (300%), and $N_I = 8, 16, 24$ (100, 200, 300%).

Next, we consider a NOMA system of $N_P = 24$, and different number of interfering UEs $N_I = 8, 16$, and 24 . For a spreading length of 8, the 8, 16, and 24 activity corresponds to an overloading of 100%, 200%, 300%, respectively. We see a similar trend in Figure 3.8 compared to the OMA system; the jointly designed NOMA codebooks outperform the reuse strategy. We observe that the gain becomes smaller, as the number of interfering UEs increases. This is to be expected due to the high amount of interference experienced at such high overloading levels in both the primary and interfering cells.

We conclude that such a joint codebook design is beneficial for group-edge UEs, if the number of strong interferes from the other groups are relatively not very high.

3.4 Final Remarks

Our focus in this chapter was on the design of spreading signatures for uplink code-domain NOMA. We considered both the optimization of a single codebook, and also the joint optimization of multiple codebooks. Our main design approach was Grassmannian or Grassmannian-like, focusing on the correlation properties of the codebooks. However, correlation is not the only relevant metric when it comes to designing such codebooks. For example, when it comes to energy consumption, transmitting signals with low peak-to-average-power ratio (PAPR) can be crucial for UEs with strict battery-consumption constrains. Nonetheless, the designs proposed here can be modified to support that. For the ICBP algorithm, we can replace the line of

$$\mathbf{s}_k \leftarrow \frac{\mathbf{s}_k + \beta \mathbf{u}_k}{\|\mathbf{s}_k + \beta \mathbf{u}_k\|}$$

with the following two lines

$$\begin{aligned} \mathbf{s}_k &\leftarrow \mathbf{s}_k + \beta \mathbf{u}_k, \\ [\mathbf{s}_k]_l &\leftarrow \sqrt{\frac{1}{L}} \frac{[\mathbf{s}_k]_l}{\|[\mathbf{s}_k]_l\|}, \quad l = 1, 2, \dots, L. \end{aligned}$$

That is, we normalize each entry of \mathbf{s}_k by its magnitude. This gives us the closest vector to \mathbf{s}_k with unit-modulus entries. The algorithm then keeps going with this normalization step in each iteration. This would result in signatures that are well-spaced, but at the same time, have unit-modulus magnitude. For the joint codebook design, it is also possible to produce signatures with low PAPR. This is done by enforcing additional constraints on the alternating projection algorithm, as done in [53].

Finally, although our focus is on code-domain NOMA based on dense spreading, we would like to mention that other code-domain signatures also exist, and can provide good multi-user suppression capabilities, such as sparse spreading sequences, user-specific interleaving, user-specific scrambling, and many more [8, 18].

4

Receive-Side Processing: NOMA Activity and Data Detection

Perhaps the main challenge when it comes to supporting uplink NOMA transmissions is the detection procedure at the BS. The presence of multi-user interference requires the joint detection across multiple UEs, which increases the detection complexity, compared to the case where each UE is transmitting alone. Combining this with multiple receive antennas can lead to a further complication, as the spatial domain should generally be utilized in a joint manner with the NOMA code-domain for an improved detection performance. In the context of grant-free access, where the BS is not aware of which UE is active at a certain time, activity detection is required as a first step to identify the active set of UEs, before proceeding with the data detection. Therefore, in such systems, the entire procedure of activity detection, channel estimation, and data detection, needs to be carried out in a low-complexity fashion.

The focus in this chapter will be on describing the NOMA detection procedure in the context of a practical frame-structure, identifying possible issues and investigating low-complexity implementations. In the first part, we consider activity detection in grant-free systems based on subspace methods. We propose the application of masking sequences to ensure a robust performance under correlated time-frequency fading, while for highly selective channels, we investigate different pilots' allocation strategies. In the second part, we focus on reducing the data detection complexity in two ways: first, by utilizing the time-frequency correlation of the channel, we reduce the number of calculated filters across the resources' grid; and second, we utilize the spatial domain to reduce the calculation complexity of the filters themselves. Our goal is showing that the overall NOMA detection procedure can be carried out with a practically viable complexity.

The framework developed in this chapter is based on our publications in [35–37].

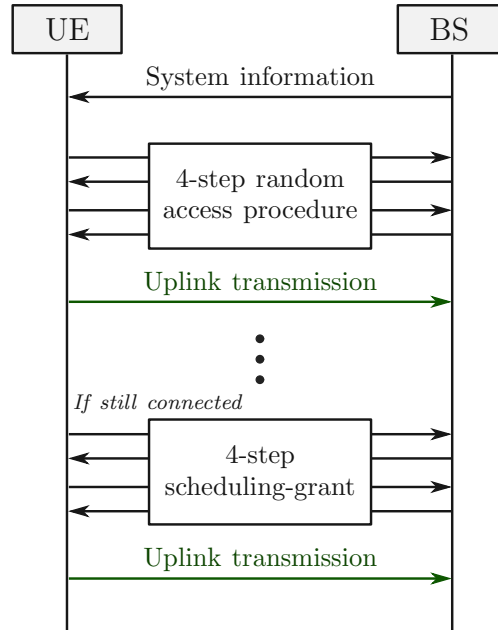


Figure 4.1: Grant-based access.

4.1 Activity Detection in Grant-Free NOMA

Conventional wireless systems are grant-based, meaning that in order for the UE to access the network, it has to be explicitly scheduled by the BS for each of its data transmissions [11]. This is illustrated in Figure 4.1, where after establishing a connection to the network via a 4-step random access procedure, it has to go through a 4-step scheduling-grant every time it has data to transmit [66]. Such an access scheme is efficient when the number of UEs accessing the network is small enough, allowing each UE to occupy its separate time-frequency resources, and when the goal is to support high data-rate per UE. However, for future systems, massive MTC will require the support of a large number of devices that can activate and transmit at an arbitrary time. The MTC traffic is typically a low-rate transmission, consisting of relatively small packets. For example, it could be a traffic-monitoring sensor, temperature meter, etc. Such a combination of massive connectivity and low-rate transmission can render current access systems inefficient; on the one hand, a large number of UEs attempting to access the network may result in a large scheduling or queuing delays at the BS before the UE gets the chance to access the network [18]. On the other hand, due to the short-packet transmission, the actual data that is transmitted in the end can be comparable in size to the control signaling required to perform the scheduling-grant in the first place, leading to an inefficient utilization of the network resources [13].

The framework of grant-free NOMA offers the possibility to tackle these issues. With grant-free access, the BS does not have to coordinate the data transmissions of the UEs, but rather the UEs transmit directly on their own, once they have data to

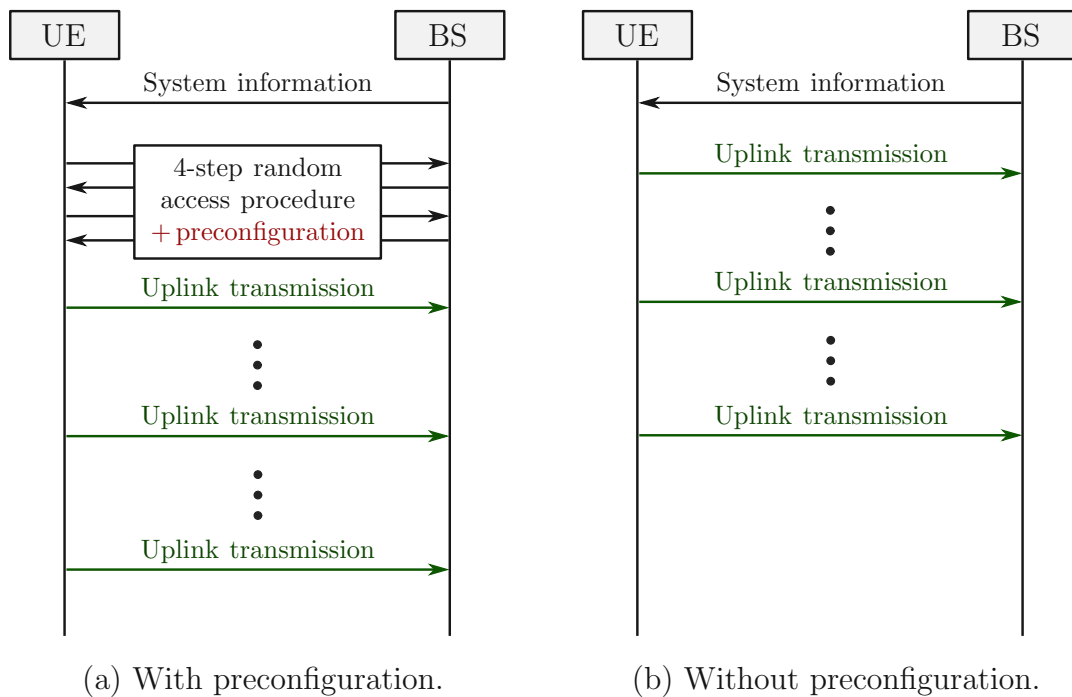


Figure 4.2: Possible grant-free access setups.

transmit. In this case, multiple UEs may choose to transmit at the same time, and therefore would contest the same resources. To manage the interference between the UEs, grant-free access can be combined with NOMA, which naturally supports multiple UEs transmitting simultaneously, and is able to manage the multi-user interference effectively via code-domain processing. However, since the BS is not aware of which UE is active at a certain time, the BS now has to perform activity detection (e.g., by utilizing the pilot sequences of the active UEs), before it proceeds with the channel estimation and data detection.

Figure 4.2 shows two possible setups, where such a framework is applied [67]. In Figure 4.2a, the scheduling grant procedure is removed, while the initial random access procedure is still maintained, but now involves a preconfiguration by the BS where the UE is informed about which resources to contest, and what pilot sequence to employ. Then, for all future transmissions, the UE transmits directly with the preassigned settings. This ensures that for different UEs contesting the same resources, they utilize different pilot sequences, which avoids pilots' collision and ensures that channel estimation can be carried out. For the other case in Figure 4.2b, the random access procedure is removed as well, and therefore there is no prior configuration of the UEs. In this case, collisions of the pilot sequences can happen, since the UEs pick the pilot sequences randomly on their own. In either case, a sufficiently large pilots' codebook, with possibly non-orthogonal sequences, is required in order to handle a large number of active UEs.

Note that in 5G Release-15, the concept of mini-slots has been introduced to reduce the transmission latency by using shorter frames [68]. Also, in the recent Release-16, the concept of 2-step random access is introduced, in which data is transmitted directly together with the random access preamble [69], similarly to what we just discussed. Both of these techniques can help reduce the access latency and support connectivity for a large number of devices in 5G. The combination we consider next with the NOMA framework is applicable to these techniques as well.

4.2 Optimizing Activity Detection via MUSIC

Activity detection algorithms aim to find the active set of UEs at a certain time. Given the short-packet nature of MTC traffic and its sporadicity, the subset of UEs active at a certain time is typically smaller than the total number of UEs available. This sparsity has motivated the application of compressed sensing (CS)-based methods, where the orthogonal matching pursuit (OMP) algorithm and its extensions have been proposed in the literature, such as in [70–73]. Another category of algorithms are those based on the estimated sample autocorrelation matrix, where subspace methods, such as MULTiple SIGNAL Classification (MUSIC) can be applied to find the active subsets, as in [74, 75]. In [76] joint activity and data detection using approximate message passing (AMP) and expectation maximization (EM) is proposed. Other algorithms such as expectation propagation (EP) has been applied in [77]. Deep learning was also considered for this problem, as in [78].

4.2.1 Considered Model for Subspace Detection

In this work, we focus on activity detection using subspace methods; namely, the MUSIC algorithm. Compared to CS-based methods, AMP, and EP, it is non-iterative and can be implemented in a parallel fashion with low-complexity. Different from the majority of the mentioned works, we assume that the data and pilots employ different spreading sequences (or signatures). Typically, relatively long sequences are employed for the activity detection and channel estimation in order to support a large number of devices; however, applying these long sequences to the data part can be highly inefficient, as it can substantially reduce the spectral efficiency of the data transmission. Therefore, for the data part we employ short spreading (e.g., $L = 4$), and the long sequences are only employed for the pilots. Of course, this can impact the activity detection performance, since now only the pilots can be utilized for it, while the data part can not. But, as we will see later, using the MUSIC algorithm and the pilots only is sufficient to achieve good performance. Moreover, we apply the activity detection in the context of a frame-structure that is similar to the LTE uplink, and therefore the formulated solution is more applicable to practice. In Figure 4.3, the frame-structure is re-illustrated again for completeness, consisting of two RBs. In the middle OFDM symbol of each RB, a pilot signature of length $L_p = 12$ is assumed (similar to LTE). The data part consists of spread blocks, each of length $L = 4$.

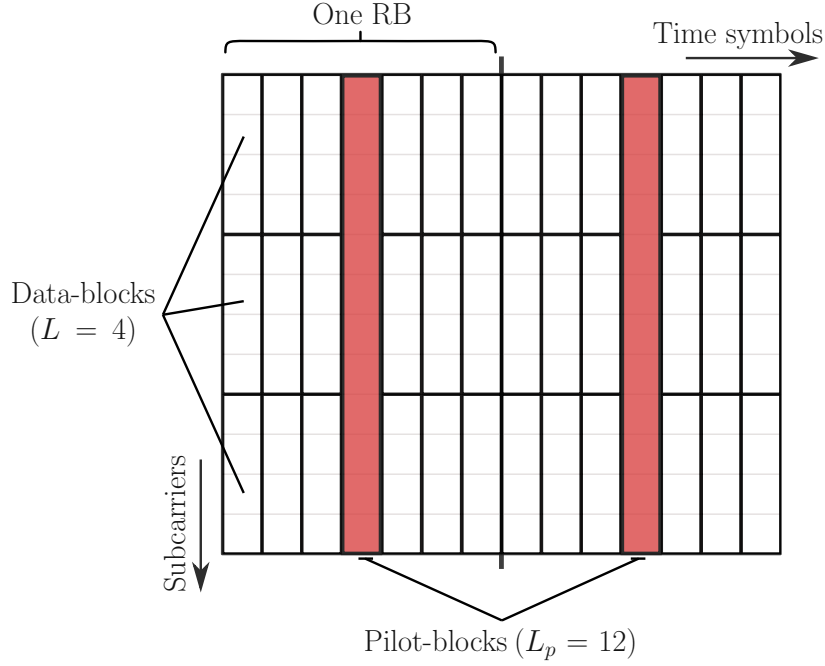


Figure 4.3: Considered frame-structure with 12 subcarriers and 14 OFDM symbols.

Formulation of the Subspace Detector

Let the number of active UEs at a certain instant be K_a out of total K , the received signal at the BS at pilot-block i is given by

$$\mathbf{y}_{p,i} = \sum_{k=1}^{K_a} \sqrt{L_p P_k} \mathbf{a}_k h_{k,i} + \mathbf{n}_{p,i}, \quad i = 1, 2, \dots, B_p, \quad (4.1)$$

where \mathbf{a}_k is the pilot signature of UE k , $h_{k,i}$ and $\mathbf{n}_{p,i}$ are channel coefficient and noise at pilot-block i , respectively, and B_p is the total number of pilot-blocks. Here, for the start, it is assumed that the fading is flat across the pilot-block; hence, it is given as a scalar. Let $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_{K_a}]$, $\mathbf{h}_i = [\sqrt{L_p P_1} h_{1,i}, \sqrt{L_p P_2} h_{2,i}, \dots, \sqrt{L_p P_{K_a}} h_{K_a,i}]^T$, (4.1) can be written equivalently as

$$\mathbf{y}_{p,i} = \mathbf{A} \mathbf{h}_i + \mathbf{n}_{p,i}. \quad (4.2)$$

Under such a system model, the autocorrelation matrix of the received signal is

$$\mathbf{R}_{\mathbf{y}_p} = \mathbb{E}\{\mathbf{y}_p \mathbf{y}_p^H\} = \mathbf{A} \mathbf{R}_{\mathbf{h}} \mathbf{A}^H + \sigma_{\mathbf{n}}^2 \mathbf{I}, \quad (4.3)$$

where $\mathbf{R}_{\mathbf{h}} = \mathbb{E}\{\mathbf{h} \mathbf{h}^H\}$. As can be seen, the autocorrelation matrix consists of two components, corresponding to the signal and noise parts. The idea behind subspace methods is, as long as $0 < K_a < L_p$, then the eigenspace of the autocorrelation matrix can be divided into signal-plus-noise and noise-only subspaces. The active

4.2. Optimizing Activity Detection via MUSIC

signatures live in the signal subspace, and theoretically, they have zero contribution to the noise subspace. Based on this, it is possible to tell which signature is active based on how much energy it has in the noise subspace. To elaborate further, consider the eigenvalue decomposition of $\mathbf{R}_{\mathbf{y}_p}$ given by

$$\mathbf{R}_{\mathbf{y}_p} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^H, \quad (4.4)$$

with $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_{L_p})$ being the matrix of eigenvalues with order $\lambda_1 > \lambda_2 > \dots > \lambda_{L_p}$, and $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{L_p}]$ is the matrix of the corresponding eigenvectors. Given that K_a UEs are active, the noise subspace \mathbf{U}_n is given by the collection of eigenvectors corresponding to the smallest $L_p - K_a$ eigenvalues, i.e.,

$$\mathbf{U}_n = [\mathbf{u}_{K_a+1}, \mathbf{u}_{K_a+2}, \dots, \mathbf{u}_{L_p}]. \quad (4.5)$$

The MUSIC spectrum is then calculated as

$$M(k) = \frac{1}{\|\mathbf{U}_n^H \mathbf{a}_k\|^2}, \quad k = 1, 2, \dots, K, \quad (4.6)$$

which simply measures the energy of all possible pilot signatures in the noise-subspace. The K_a signatures with highest $M(k)$ are then declared active. All the other signatures that are not active occupy a portion of the noise-subspace and therefore will result in small $M(k)$. In practice, we do not have access to the true autocorrelation matrix, but rather we estimate it from the received pilot-blocks. That is, we utilize the sample autocorrelation matrix calculated as

$$\hat{\mathbf{R}}_{\mathbf{y}_p} = \frac{1}{B_p} \sum_{i=1}^{B_p} \mathbf{y}_{p,i} \mathbf{y}_{p,i}^H. \quad (4.7)$$

If multiple receive antennas are available at the BS, then this estimate can be further improved by averaging over the pilot-blocks from all of the receive antennas.

Estimation of K_a

In order to determine the noise subspace in (4.5) and consequently pick the strongest K_a signatures from (4.6), we need to know K_a a priori; however, that is not possible in grant-free access, since the BS is not aware of how many UEs are active at a certain time. A similar issue also exists in other activity detection methods, such as in CS, where the sparsity level of the problem needs to be known. Therefore, K_a has to be estimated from the received signal as well, and here, we utilize the Bayesian information criterion (BIC) [79, 80], which estimates the dimensionality of a fitting model from measured samples. Other information criteria also exist, such as the Akaike information criterion (AIC) [81]; but BIC generally is preferred, as it is a consistent selector [82], i.e., as the number of available samples (pilot-blocks) increases, the probability of choosing the correct model (with correct K_a) approaches

4.2. Optimizing Activity Detection via MUSIC

100%. However, as argued in [82], this inconsistency might not necessary be a flaw in AIC, as it can provide an advantage for certain problems compared to BIC. In our own testing in the context of activity detection in grant-free NOMA, we found that BIC performs better, and therefore this is what we adopt next. Note that BIC has been applied already for NOMA activity detection in [74]. Compared to that work, we formulate it here for the case where the noise power is known at the BS, in a fashion similar to [83]. The BIC is defined as

$$\text{BIC}(K_a) \triangleq -\log f(\mathbf{y}_p; \hat{\boldsymbol{\theta}}(K_a)) + \frac{1}{2} |\hat{\boldsymbol{\theta}}(K_a)| \ln B_p, \quad K_a = 1, 2, \dots, L_p - 1, \quad (4.8)$$

where $f(\mathbf{y}_p; \hat{\boldsymbol{\theta}}(K_a))$ is the likelihood function of the received pilot-blocks, evaluated at the maximum likelihood (ML) estimate of the underlying distribution parameters $\hat{\boldsymbol{\theta}}(K_a)$ under the assumption of K_a active components, and $|\hat{\boldsymbol{\theta}}(K_a)|$ is the number of parameters. Under the assumption of Gaussian noise, zero-mean received signal, and the samples being i.i.d, (4.8) becomes [83]

$$\text{BIC}(K_a) = B_p \sum_{k=1}^{L_p} \left(\frac{\hat{\lambda}_k}{\bar{\lambda}_k} + \ln \bar{\lambda}_k \right) + B_p L_p \ln \pi + \frac{K_a}{2} (2L_p - K_a) \ln B_p, \quad (4.9)$$

where $\hat{\lambda}_k$ is the k -th eigenvalue of the estimated sample autocorrelation matrix $\hat{\mathbf{R}}_{\mathbf{y}_p}$, and we set $\bar{\lambda}_k$ as

$$\bar{\lambda}_k = \begin{cases} \hat{\lambda}_k, & k = 1, 2, \dots, K_a, \\ \sigma_{\mathbf{n}}^2, & k = K_a + 1, \dots, L_p. \end{cases} \quad (4.10)$$

That is, for the current K_a under investigation, the eigenvalues up to K_a are set equal to their ML estimate which we simply obtain from an eigendecomposition of the sample autocorrelation matrix, while the remaining eigenvalues are supposed to belong to the noise-only subspace and thus are set equal to the known noise power. Our estimate of K_a is then the one that minimizes BIC, i.e.,

$$\hat{K}_a = \arg \min_{K_a} \text{BIC}(K_a). \quad (4.11)$$

It can happen that at a certain time, no UE is active (i.e., $K_a = 0$), which we would not be able to tell from (4.8). Typically, BSs have a certain threshold for decodability; signals with powers below that threshold are considered too weak to be detected anyway. Therefore, one way to tell whether there is an activity at all, is to compare the strongest estimated eigenvalue to that threshold. Let the threshold be σ_{\min}^2 , then K_a is estimated as

$$\hat{K}_a = \begin{cases} 0, & \hat{\lambda}_1 < \sigma_{\min}^2, \\ \arg \min_{K_a} \text{BIC}(K_a), & \text{otherwise.} \end{cases} \quad (4.12)$$

4.2.2 Impact of Strong Time-Frequency Correlation

The subspace activity detector was formulated under the assumption that the autocorrelation matrix can be divided into distinguished signal-plus-noise and noise-only subspaces. However, in order for this to work, the autocorrelation matrix of the channel coefficients, i.e., \mathbf{R}_h has to be full rank in order for the whole product of $\mathbf{A}\mathbf{R}_h\mathbf{A}^H$ in (4.3) to have a rank equal to the number of active UEs. To get a full rank \mathbf{R}_h , the channel coefficients across the different pilot-blocks need to be i.i.d., which corresponds to the case where the channel is highly selective from one pilot-block to another; however, typically, the channel is correlated in time and frequency, and therefore neighbouring pilot-blocks, especially in time, are likely to experience a similar channel. This is utilized in practical wireless systems, where only a few pilots are transmitted with the data, which are then used to obtain channel estimates at their positions, and then for the data positions, the channel is simply obtained by interpolation between the pilots. In other words, in many situations, the signal received at the pilot-blocks can be highly correlated, and this can be an issue for our detector. To further elaborate on this, let us consider the worst-case scenario where all the UEs undergo flat-fading in both time and frequency. Under such an assumption, we will have $\mathbf{h}_1 = \mathbf{h}_2 = \dots = \mathbf{h}_{B_p} = \mathbf{h}$, i.e., a constant, resulting in $\mathbf{R}_h = \mathbb{E}\{\mathbf{h}\mathbf{h}^H\} = \mathbf{h}\mathbf{h}^H$. The term $\mathbf{h}\mathbf{h}^H$ is simply an outer product of a vector with itself, i.e., a rank-1 matrix. Our resultant autocorrelation matrix is then given by

$$\mathbf{R}_{y_p} = \mathbf{A}\mathbf{h}\mathbf{h}^H\mathbf{A}^H + \sigma_n^2\mathbf{I}. \quad (4.13)$$

Consequently, the signal part is rank-1; hence, we can only detect one active UE. This can also be observed by looking at the sample autocorrelation matrix in (4.7). If the samples (pilot-blocks) used are correlated, or in a worst-case scenario identical, then it would be a sum of the same outer product, resulting in a rank-1 estimate of the signal part of the autocorrelation matrix. Note that the BIC expression in (4.9) is derived under an i.i.d. factorization of the likelihood function. Therefore, the correlation between the pilot-blocks also impacts this simplified BIC calculation.

Here, to tackle this issue, we propose the use of masking sequences, applied on top of the pilot-blocks. Figure 4.4 illustrates the application of masking sequences for UEs transmitting over four RBs, where UE1 transmits with the red pilot signature, while UE2 transmits with the blue one. Instead of transmitting the pilot signatures directly, they are overlaid with a masking sequence. In the considered example, UE1 uses the masking sequences $[+1, +1, -1, -1]$, while UE2 uses $[+1, -1, +1, -1]$. The system model now becomes

$$\mathbf{y}_{p,i} = \sum_{k=1}^{K_a} \sqrt{L_p P_k} \mathbf{a}_k h_{k,i} m_{k,i} + \mathbf{n}_{p,i}, \quad (4.14)$$

where $m_{k,i}$ is the i -th coefficient of UE- k masking sequence applied at the i -th pilot-block. Let \mathbf{m}_i be the collection of the masking coefficients across the different UEs

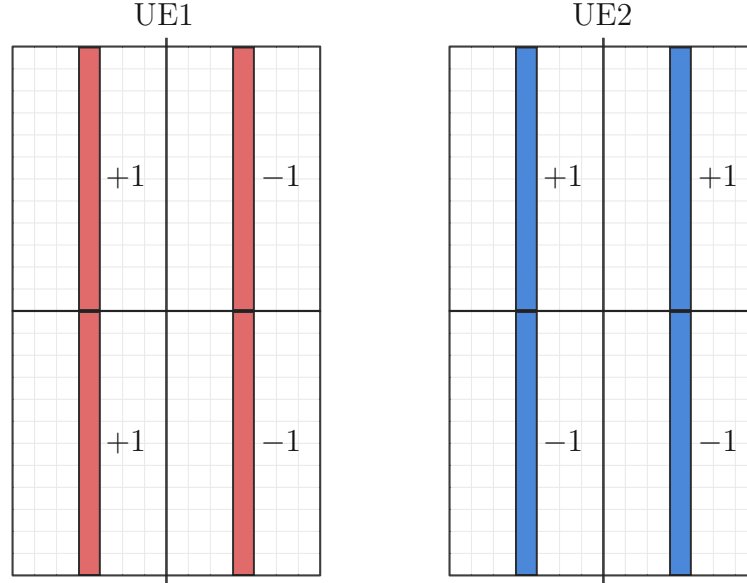


Figure 4.4: Illustration of the masking sequences for two UEs over 4 RBs.

at block i , then

$$\mathbf{y}_{p,i} = \mathbf{A}(\mathbf{h}_i \circ \mathbf{m}_i) + \mathbf{n}_{p,i}, \quad (4.15)$$

where \circ denotes the element-wise Hadamard product. Treating the \mathbf{m}_i as a random process across the pilot-blocks, it can be easily shown that the autocorrelation matrix then is given by

$$\mathbf{R}_{\mathbf{y}_p} = \mathbf{A}(\mathbf{R}_{\mathbf{h}} \circ \mathbf{R}_{\mathbf{m}})\mathbf{A}^H + \sigma_{\mathbf{n}}^2\mathbf{I}, \quad (4.16)$$

where $\mathbf{R}_{\mathbf{m}} = \mathbb{E}\{\mathbf{m}\mathbf{m}^H\}$. Therefore, by designing the masking sequences to be like an i.i.d. process across the UEs, then it would be possible to recover a correct estimation of the signal part of the autocorrelation matrix. For example, even if we have $\mathbf{R}_{\mathbf{h}} = \mathbf{h}\mathbf{h}^H$, the product $\mathbf{h}\mathbf{h}^H \circ \mathbf{R}_{\mathbf{m}}$ can still be full rank if $\mathbf{R}_{\mathbf{m}}$ is full-rank, which is achieved by designing the masking sequences random-like. One way to design these sequences is to pick them randomly from the alphabet $\{-1, +1\}$, as illustrated in the example of Figure 4.4. Other designs are also possible, which might enjoy more structure, such as binary Golay sequences [84]. Those masking sequences would be defined in the standard in a similar fashion as the pilot signatures, preferably having a one-to-one correspondence with them. That is, for each pilot signature, there is a corresponding masking sequence.

In the next stage where channel estimation is performed, the effect of employing the masking sequence is removed from the channel estimate. Let the matrix of

estimated active signatures be $\hat{\mathbf{A}}$, the masking-free channel estimate is given by

$$\hat{\mathbf{h}}_i = ((\hat{\mathbf{A}}^H \hat{\mathbf{A}})^{-1} \hat{\mathbf{A}}^H \mathbf{y}_{p,i}) \oslash \mathbf{m}_i, \quad (4.17)$$

where \oslash denotes the element-wise division. Once the channel estimates are obtained over all pilot-blocks, linear interpolation (and also extrapolation) is performed to obtain the channel for the whole time-frequency grid.

Example Scenario with Eight Active UEs

We now evaluate the performance of the considered activity detection framework with masking sequences under a grant-free system with preconfiguration (i.e., as in Figure 4.2a). We consider a scenario where there are $K = 32$ possibly active UEs; however, at a certain time, only $K_a = 8$ UEs are active. They contest a resources' region of 72 subcarriers \times 14 time symbols corresponding to 12 RBs, and the BS employs 2 receive antennas. This brings the number of pilot-blocks that are used to calculate the sample autocorrelation matrix to 24 (12 for each antenna). The task of the receiver is then find the active pilot set, perform channel estimation, and finally equalize the data part and decode it. We consider a scenario with high time-frequency correlation, by assuming a Pedestrian-A channel model, which has a low root-mean-square (RMS) delay spread of 45 ns, and assume the UEs are static and therefore the channel is time-invariant. We assume that there is a pathloss spread of ± 5 dB between the UEs, meaning that the strongest and weakest UEs can have a gap in the receive power of up to 10 dB. Both the data and pilot signatures are from Grassmannian codebooks. As for the masking sequences, we construct them randomly once from $\{+1, -1\}$, and then fix them for all the simulation repetitions. The simulation parameters are summarized in Table 4.1.

Figure 4.5 shows the average number of correctly decoded UEs using the proposed method versus their average SNR at the BS. The goal here is to have all of the 8 active UEs identified and decoded correctly. The perfect activity detection denotes the case where the BS knows exactly which UEs are active, and therefore it only has to perform channel estimation and data detection. This serves as a baseline for our results. As can be seen in the figure, combining the subspace activity detector with the masking sequences results in a performance that is very close to the case with perfect activity detection. When no masking sequences are employed, then due to the correlation of the channel, it is difficult for the signal part to be constructed properly, thus greatly deteriorating the performance.

Note that if the data part was also spread with the same signatures as the pilots, as commonly done in the literature, then in that case, the received signal from the data part can also be utilized in the construction of the sample autocorrelation matrix. As the data symbols can be assumed i.i.d., they would serve a similar function as the masking sequences, allowing the signal part of the autocorrelation matrix to reach the required rank. However, as we mentioned, spreading the data symbols with long sequences can be spectrally inefficient, and therefore we utilize only the pilots for this task, necessitating the application of the masking sequences.

4.2. Optimizing Activity Detection via MUSIC

Parameter	Value
Active UEs	$K_a = 8$ out of total $K = 32$
Contention region	72 subcarriers \times 14 time symbols
Data signatures	$L = 4$ (4×16 Grassmannian codebook)
Pilots signatures	$L_p = 12$ (12×32 Grassmannian codebook)
Masking sequences	Randomly generated from $\{+1, -1\}$
Receive antennas	$N_R = 2$
Center frequency	2 GHz
Subcarrier spacing	15 kHz
Pathloss spread	± 5 dB
Channel model	Rayleigh, Ped-A (45 ns RMS), velocity = 0
Modulation	4-QAM
Channel coding	Turbo, code-rate 2/3
Activity detection	BIC + MUSIC
Channel estimation	LS with linear inter/extrapolation
Data detection	MMSE-PIC (max. 6 iters)

Table 4.1: Simulation parameters for activity detection with correlated channels.

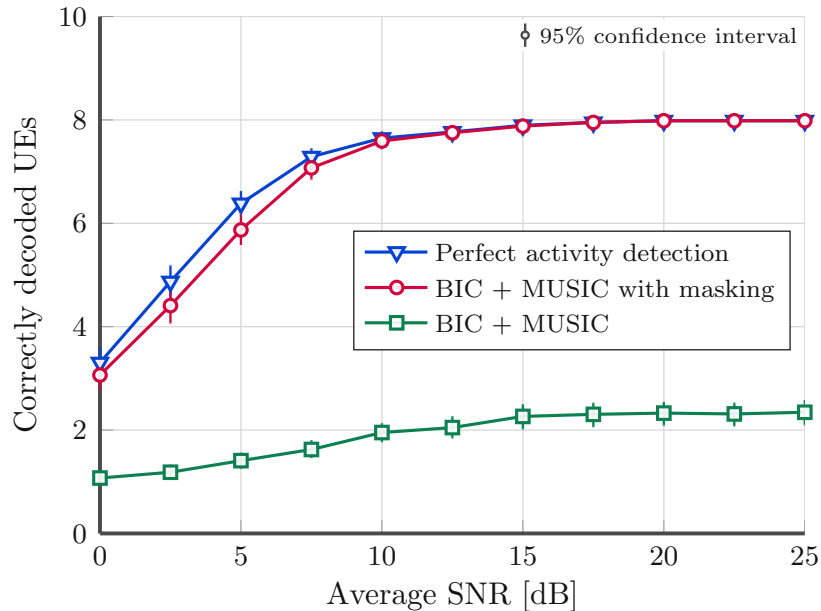


Figure 4.5: Activity detection via the detector with masking sequences.

4.2.3 Impact of Strong Time-Frequency Selectivity

In the previous subsection, we discussed the benefit of having a selective channel, as it can reduce the correlation between the received signal at the pilot-blocks, which in turn can improve the detection performance. However, since our transmission is sequence-based, we require the channel to be flat within the spreading interval. Therefore, if the channel is too selective, then our system model would not hold anymore, since different parts of the sequence would experience different channel conditions. This is illustrated in Figure 4.6 for a tapped-delay-line-C (TDL-C) channel [85] with a 300 ns RMS delay spread, shown over two RBs. As can be seen in the considered example with $L_p = 12$, the channel can vary substantially within the pilot-block, thus possibly destroying the structure of the pilot signatures.

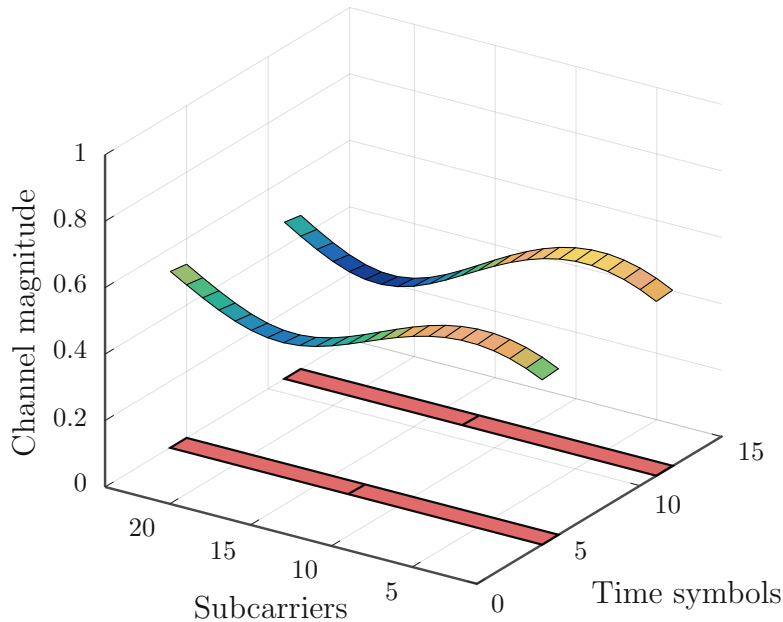
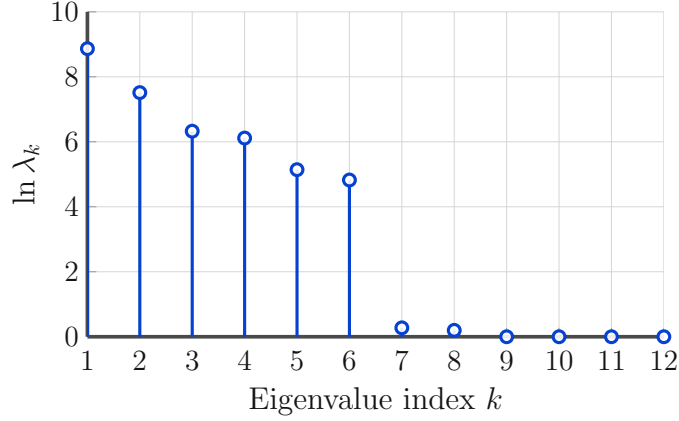
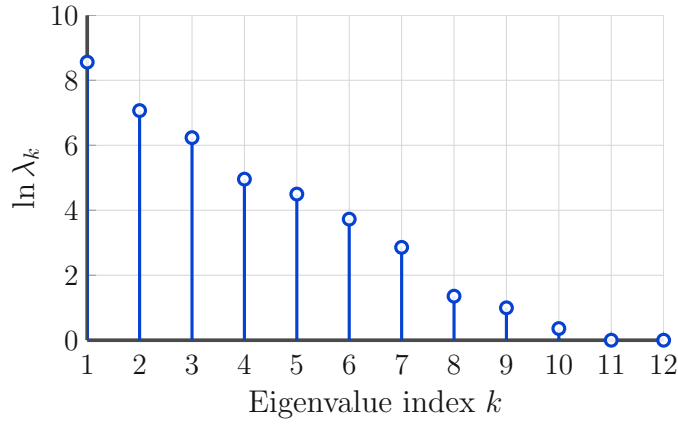


Figure 4.6: Illustration of a frequency-selective channel over the pilot-blocks.

This modification of the pilot signature changes from one pilot-block to another, and also it is different for the different UEs. Under such conditions, the eigenstructure of the autocorrelation matrix no longer reflects the activity of the pilot signatures, but rather a modified version of them. The more selective the channel is, the more corrupt the estimated autocorrelation matrix will be. In Figure 4.7a, the eigenvalues of a realization of the autocorrelation matrix for a TDL-C channel with no selectivity is shown for $K_a = 6$ and $L_p = 12$. A clear distinction can be seen between the signal part and the noise part in terms of the magnitude of the eigenvalues. On the other hand, Figure 4.7b shows the case with 300 ns RMS delay spread. As can be seen, the subspaces are not very clearly separable, which makes



(a) 0 ns RMS delay spread.



(b) 300 ns RMS delay spread.

Figure 4.7: Eigenvalues of one realization of $\hat{\mathbf{R}}_{\mathbf{y}_p}$ over a TDL-C channel for $K_a = 6$ and $L_p = 12$.

it challenging for our subspace detector. The easiest solution to the selectivity problem, is to use shorter pilot sequences; however, the shorter the pilot sequences are, the smaller is the number of UEs that can be supported. Therefore, we have to keep the pilot sequences long enough, in order to support a sufficient number of connections, while at the same time reduce the impact of selectivity. To that end, we consider here repositioning the pilot-blocks around the RBs. Instead of inserting the pilot sequence along a single RB over 12 consecutive subcarriers, it is split over two neighbouring RBs in time. This is illustrated in Figure 4.8, where a pilot sequence of length $L_p = 12$ is split over two blocks, each with a length of 6. With such a setup, the pilot length in the frequency direction is halved, and therefore less frequency-selectivity is experienced per sub-sequence. The drawback of such an allocation is the worse time-resolution for the channel estimation, since now only a single effective sequence is used to estimate the channel in time.

We investigate the benefit of such an allocation strategy using the same setup of Table 4.1, but we change the channel model to TDL-C with varying RMS delay

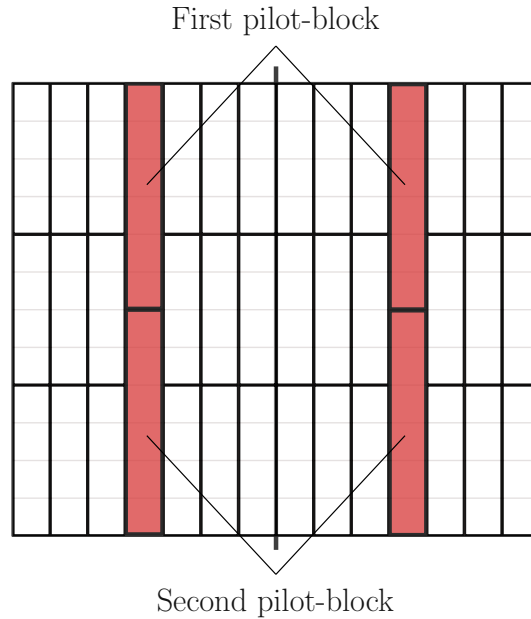


Figure 4.8: Splitting the pilot sequences over two blocks in time for $L_p = 12$.

spread and fix the average SNR of the UEs to 20 dB. The result is shown in Figure 4.9. There are two factors impacting the performance here. On the one hand, due to the increased frequency-selectivity, the performance of the activity detection deteriorates, as we just explained before. On the other hand, the frequency-selectivity impacts the channel estimation performance as well, i.e., as the channel gets more selective, denser pilots are required in order to be able to track the variations in the channel. As can be seen in the figure, splitting the pilot sequence provides gains for both of these issues. The gain for the activity detection can be seen by comparing the performance under BIC + MUSIC, which shows better robustness against the selectivity with the splitting strategy. The channel estimation gain can be clearly seen when comparing the performance under perfect activity detection. In this case, the split strategy offers better frequency resolution, which helps to provide a better estimation of the channel. Of course, this comes at the cost of decreased robustness to time-selectivity. If time-selectivity is an issue, then one solution is to allocate more OFDM symbols to the pilots. For example, the 5th and 12th OFDM symbols would hold pilot-blocks as well, and then the split can occur over neighbouring OFDM symbols only, leading to less time-variation within the split-block. The disadvantage of doing so is a lower data transmission rate.

Finally, it is important to mention that the setup we considered here is a worst-case scenario in which the channels of all the UEs are suffering from high delay spreads. This might not be the case in practice. Also, per 3GPP [85], 300 ns is already considered a long delay spread, and they declare the nominal delay spread to be 100 ns. Therefore, above 400 ns would represent extreme cases. Under moderate conditions, the deterioration in performance may not be severe then, especially when combined with those splitting strategies.

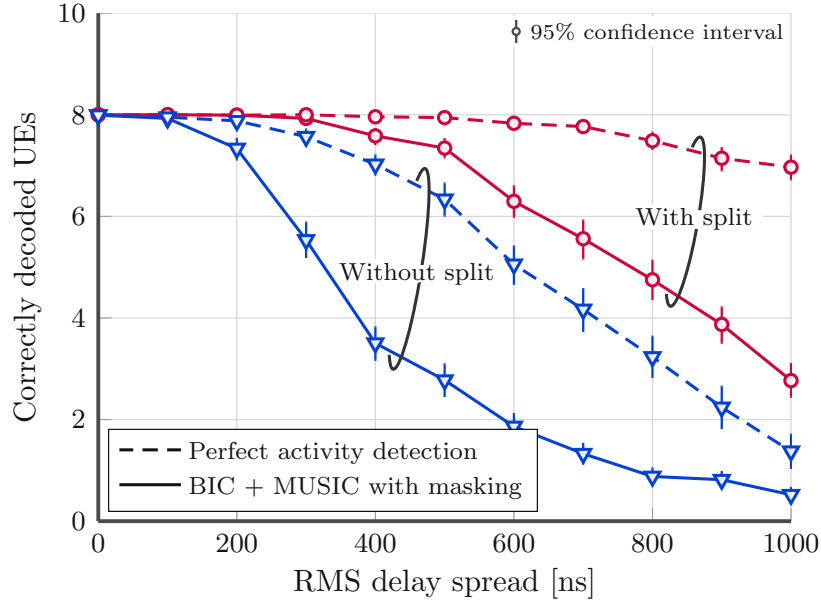


Figure 4.9: Detection performance over a frequency-selective channel with $K_a = 8$.

4.3 Reducing Data Detection Complexity

The previous part of the chapter focused on the pilots' part of the received signal, which is utilized for the activity detection in the case of grant-free systems, and for channel estimation. In the remainder of this chapter, we focus on the data part. We continue here with the same frame-structure as the one we considered before in Figure 4.3. For each of the data-blocks, multi-user detection is performed in order to recover the transmit symbols of the UEs. More specifically, from (2.6), the BS calculates an MMSE estimate for each spreading-block i , i.e.,

$$\mathbf{x}_i^{\text{MMSE}} = \mathbf{G}_i^H (\mathbf{G}_i \mathbf{G}_i^H + \sigma_n^2 \mathbf{I}_{LN_R})^{-1} \mathbf{y}_i, \quad i = 1, 2, \dots, B_d. \quad (4.18)$$

Once the B_d symbol estimates of the entire subframe are obtained, channel decoding is performed and the resulting messages are checked for CRC. All the UEs that pass the CRC get their signals removed via IC, and the remaining ones are equalized again through a next iteration of IC using the cleared-up signal from the previous iteration.

Depending on whether the right- or left-pseudoinverse is used when calculating the MMSE filter, its complexity can scale with the number of UEs, the spreading length, and/or the number of receive antennas. Moreover, the calculation in (4.18) is done B_d times, i.e., for each of the data-blocks, which can be an issue for large B_d . In the following, we investigate possible ways to reduce the detection complexity, either by reducing the number of calculated filters, or by reducing the calculation complexity of the filter itself.

4.3.1 Exploiting Time-Frequency Correlation

Following the same structure as in Figure 4.3, each RB consists of 12 subcarriers and 7 symbols, with a pilot-block in the middle, and we adopt a data spreading of length $L = 4$. Therefore, a RB would consist of $12 \times 6 / 4 = 18$ data-blocks. In this case, the BS needs to calculate an MMSE filter 18 times in every RB, since there are 18 of them. However, by inspecting \mathbf{G}_i , we notice that the only changing quantity from one data-block to another are the channel fading coefficients $\mathbf{h}_{k,i}$. The spreading signature \mathbf{s}_k and the transmit power P_k are fixed during the entire subframe transmission. This motivates us to look into the dynamic behavior of $\mathbf{h}_{k,i}$, because if the whole subframe has a constant $\mathbf{h}_{k,i}$, then \mathbf{G}_i would be constant in i , and the BS would only need to calculate a single MMSE filter and reuse it for the entire subframe. Such an assumption does not hold in general, due to the frequency-selectivity and time-selectivity resulting from multi-path propagation and motion, respectively. However, as we have seen in the previous sections, as long as the selectivity is not so severe, the wireless channel will exhibit correlation in time and frequency, at least between neighboring blocks.

Let us analyze the effect of reusing filters on the post-filtering/equalization SINR. Consider the detection of data-block i using a filter calculated at another block j . Let $\mathbf{g}_{k,i}$ be the k^{th} column of \mathbf{G}_i , then from (4.18), the post-filtering SINR of the k^{th} UE at block i is given by

$$\text{SINR}_{k,i}^{\text{MMSE}} = \frac{|\mathbf{g}_{k,j}^H \mathbf{Z}_j^{-1} \mathbf{g}_{k,i}|^2}{\sum_{l \neq k} |\mathbf{g}_{k,j}^H \mathbf{Z}_j^{-1} \mathbf{g}_{l,i}|^2 + \sigma_{\mathbf{n}}^2 \|\mathbf{Z}_j^{-1} \mathbf{g}_{k,j}\|^2}, \quad (4.19)$$

where

$$\mathbf{Z}_j = \mathbf{G}_j \mathbf{G}_j^H + \sigma_{\mathbf{n}}^2 \mathbf{I}. \quad (4.20)$$

where we write \mathbf{I} short for \mathbf{I}_{LN_R} . Under the assumption that the fading is independent across the different UEs, the interference power term in (4.19) will not be affected by the choice of j with respect to i . The noise power term is also unaffected by the choice of j . Therefore, we only need to consider the desired signal power term. Since \mathbf{Z}_j^{-1} is Hermitian positive semi-definite, we can apply the Cholesky decomposition $\mathbf{Z}_j^{-1} = \mathbf{U}_j^H \mathbf{U}_j$, where \mathbf{U}_j is an upper-triangular matrix. The numerator of (4.19) becomes

$$|\mathbf{g}_{k,j}^H \mathbf{Z}_j^{-1} \mathbf{g}_{k,i}|^2 = |\mathbf{g}_{k,j}^H \mathbf{U}_j^H \mathbf{U}_j \mathbf{g}_{k,i}|^2. \quad (4.21)$$

Using the Cauchy-Schwarz inequality, we formulate an upper bound on the desired signal power as follows

$$|\mathbf{g}_{k,j}^H \mathbf{U}_j^H \mathbf{U}_j \mathbf{g}_{k,i}|^2 \leq \|\mathbf{U}_j \mathbf{g}_{k,j}\|_2^2 \|\mathbf{U}_j \mathbf{g}_{k,i}\|_2^2, \quad (4.22)$$

with the maximum achieved in the case of proportionality, i.e., for some constant α ,

4.3. Reducing Data Detection Complexity

we have

$$\begin{aligned}\mathbf{U}_j \mathbf{g}_{k,j} &= \alpha \mathbf{U}_j \mathbf{g}_{k,i} \\ \mathbf{g}_{k,j} &= \alpha \mathbf{g}_{k,i},\end{aligned}\tag{4.23}$$

where in the last step, the inverse of \mathbf{U}_j was applied to both sides, as \mathbf{U}_j is always invertible. One maximizer is the choice $j = i$, i.e., we filter block i using a filter calculated at block i , which is the trivial case. However, we can also see that as long as $\mathbf{g}_{k,j}$ is similar to $\mathbf{g}_{k,i}$, that is, they are correlated, then the loss in the post-filtering SINR might be acceptable. In general, we would like to choose an optimal j^* , such that the minimum post-filtering SINR across all the blocks and all UEs, is maximized, i.e.,

$$j^* = \arg \max_j \min_{k,i} |\mathbf{g}_{k,j}^H \mathbf{Z}_j^{-1} \mathbf{g}_{k,i}|^2,\tag{4.24}$$

which is difficult to solve, especially when considering correlated fading. Of course, we can also perform this search on a per-RB basis, instead of the entire subframe, which is then expected to perform better in the case of high selectivity. Still, it does not make it easier to solve. We thus turn our attention to low-complexity suboptimal strategies based on the correlation properties of the channels. We consider the following four strategies, which are depicted in Figure 4.10:

- Receiver (a): assumes the channels are highly selective in frequency, but correlated in time. It only calculates MMSE filters for the spreading-blocks in the middle of the subframe (7th OFDM symbol), and then these filters are reused for the other time symbols.
- Receiver (b): assumes the channels are moderately selective in frequency, but correlated in time. It only calculates MMSE filters for the middle subcarrier group of each RB at the middle time symbol (7th OFDM symbol).
- Receiver (c): assumes the channels are highly selective in frequency, but moderately selective in time. It calculates MMSE filters for all spreading-blocks of the 4th and 11th OFDM symbols. Then, for the first slot, the 4th symbol filters are reused, and for the second slot, the filters from the 11th symbol are reused.
- Receiver (d): assumes the channels are moderately selective in both frequency and time. It calculates a single MMSE filter in the middle of every RB, and then the entire RB is filtered with that filter.

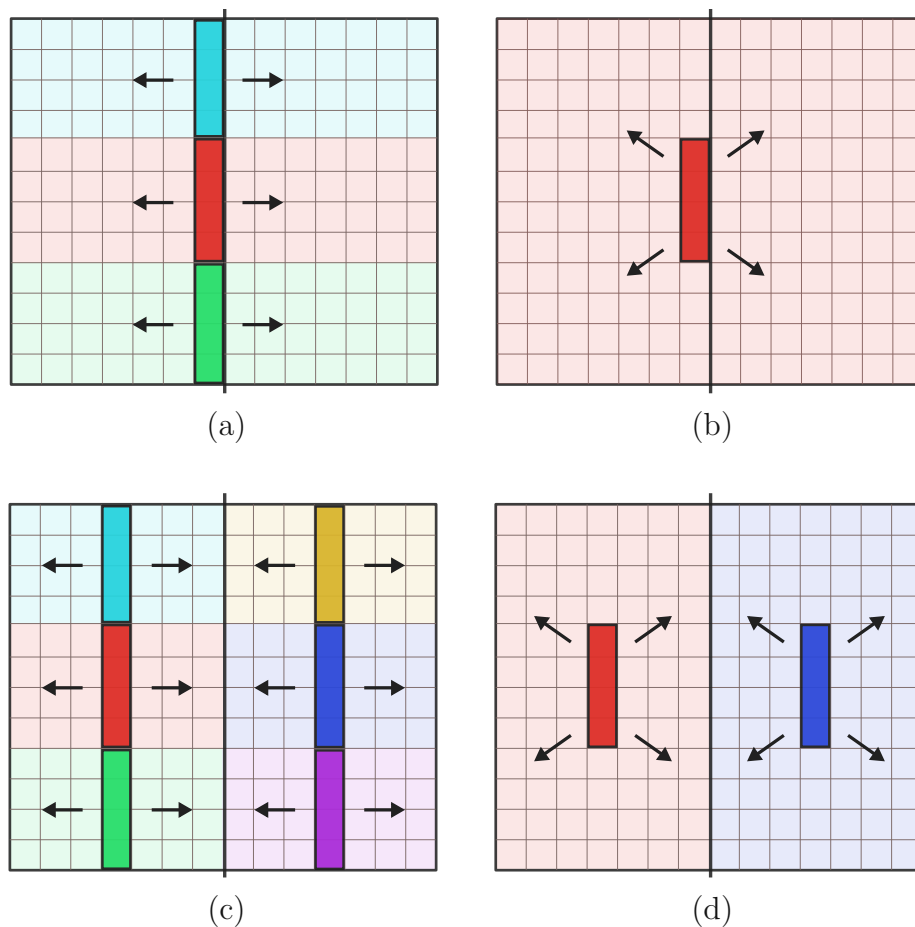


Figure 4.10: Receiver strategies under consideration. The solid blocks indicate the place where the MMSE filters are calculated. The shades indicate where they are reused.

Simulation Scenario Setup

We put the different receiver strategies to test, and compare them against a full receiver that calculates MMSE filters for all spreading-blocks. The setup is as follows: $K_a = 12$ of $K = 16$ UEs are simultaneously transmitting with a spreading length of $L = 4$ in a grant-free manner. We assume here perfect activity detection and channel estimation at the BS, and also the transmit signal is multiplied directly with the channel coefficients, without actual OFDM modulation and demodulation. The reason for doing so is to assess the performance loss coming only from the reuse of the calculated filters between the blocks, and not due to activity detection and channel estimation errors, or due to the impact of delay and Doppler spread on the OFDM transmission. The other simulation parameters are similar to what we considered before and are summarized in Table 4.2 below.

All the results shown next are accompanied by the performance of a single user

4.3. Reducing Data Detection Complexity

Parameter	Value
Active UEs	$K = 12$
Contention region	72 subcarriers \times 14 time symbols
Data signatures	$L = 4$ (4×16 Grassmannian codebook)
Receive antennas	$N_R = 2$
Center frequency	2 GHz
Subcarrier spacing	15 kHz
Average SNR	6 dB
Pathloss spread	± 5 dB
Channel model	Rayleigh, TDL-C
Modulation	4-QAM
Channel coding	New-radio LDPC, code-rate 1/2
Activity detection	Perfect
Channel estimation	Perfect
Data detection	MMSE-PIC (max. 6 iters)

Table 4.2: Simulation parameters for the filters' reuse scenario.

occupying the time-frequency resources alone, detected without any of the approximations that we propose here. This represents the performance of a perfect OMA spreading system that does not suffer from multi-user interference, and thus serves as a baseline to our NOMA results. Note that for $L = 4$, an OMA system can only support up to four UEs.

Performance over Delay Spread

First, we investigate the performance of such a system in terms of the average BLER at various levels of the RMS delay spread. All the UEs are moving at a fixed velocity of 50 km/h. The results are shown in Figure 4.11. We observe that receivers (b) and (d) exhibit a bad performance at very high delay spreads. This is expected since the MMSE filters for those receivers are only calculated at the middle subcarrier groups, and therefore when the channels are highly frequency-selective, the calculated filters fail to represent the other subcarrier groups. Also, in the simulation, we set all the UEs to have the same delay spread, i.e., the results represent a worst-case scenario. In practice, some UEs will experience short delay spread, others will experience a long one. Based on that, we conclude that for moderate to long delay spreads, any of the receiver strategies is applicable with a small difference in performance compared to a full receiver that calculates an MMSE filter for every spreading-block. Also, we notice that receivers (a) and (b) have a slightly worse performance than (c) and (d) at low delay spreads. The difference is due to the UEs moving at 50 km/h, i.e., the channels are time-variant, for which the later receivers are better suited, as we see next.

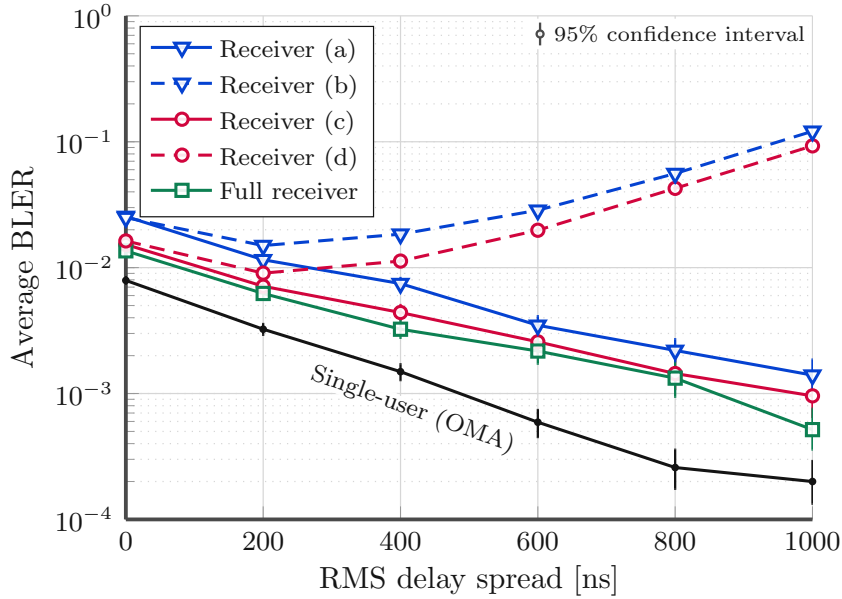


Figure 4.11: Performance of the receiver strategies at different levels of delay spread.

Performance over Velocity

Now, we fix the delay spread of the UEs to be 300 ns, and we investigate the BLER performance at different velocities. The results are shown in Figure 4.12. We observe that receivers (c) and (d) are more robust to the time-selectivity, especially at high velocities. This is due to the MMSE filters being calculated in each time-slot of the subframe, which allow them to better approximate the filters in each slot, compared to the case where a single filter is used to approximate the entire subframe (two slots) in the case of receivers (a) and (b).

Discussion

From the previous results, we propose employing receiver strategy (d). It exhibits only a small performance loss in moderate to long delay spreads, and it is very robust to high Doppler spreads. In terms of complexity saving, it calculates a single filter in the middle of each RB, which is then applied to the entire RB. Since each RB has 18 data-blocks as mentioned before, then receiver (d) would provide a complexity reduction of $18/1 = 18$ times, compared to a full blown receiver that calculates an MMSE filter at every data-block. If the environment is expected to be highly frequency-selective, then receiver (c) would be a better choice, but it will have a reduction in complexity of only $18/3 = 6$ times, compared to a full receiver.

It can also be noticed that in both figures, the performance of the full receiver gets better as the channels get more selective. The reason for that is due to the diversity gain harvested by the channel code at the bit-level [86], and also by the short spreading signatures when the channel becomes highly selective in the direction of spreading (frequency in our case). However, under practical activity detection,

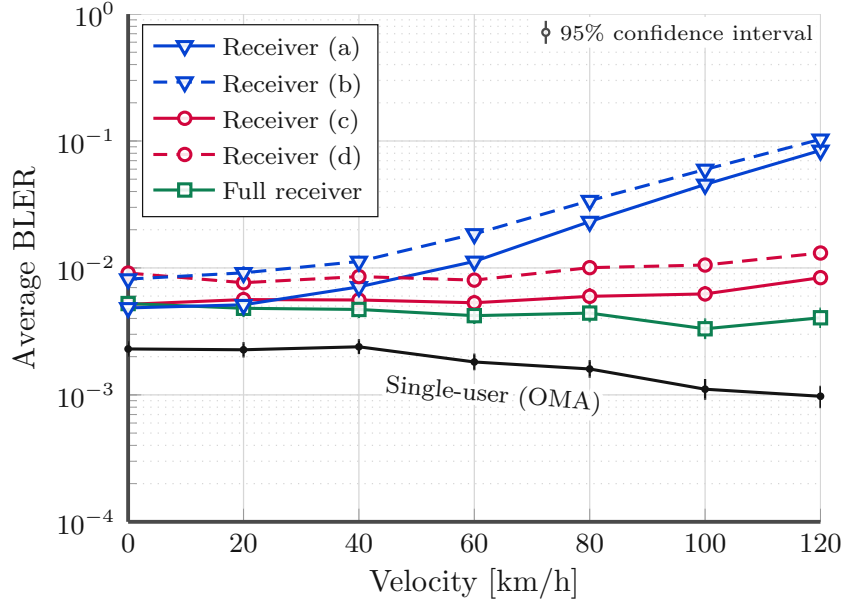


Figure 4.12: Performance of the receiver strategies at different velocities.

channel estimation, and realistic OFDM modulation, we expect the performance to become worse at high delay and Doppler spreads, as the impact of these imperfections would dominate the overall detection performance.

4.3.2 Exploiting the Spatial Domain

Adding more receive antennas at the BS makes the columns of \mathbf{G}_i in (2.2) longer, i.e., the effective signatures of the UEs get longer. Combined with the assumption that the antennas experience sufficiently uncorrelated fading, these resultant signatures will in turn have a lower cross-correlation, i.e., less interference between the UEs. This behavior is well-known in the context of massive MIMO, which results in linear receivers performing close to optimal, and even a MF can provide a good performance [87]. In our case, we do not have a massive MIMO assumption, and we employ iterative receivers. The question then is, as we add more antennas, do we really need to use an exact MMSE filter for the NOMA detection in (4.18), or is it sufficient to employ an approximation, such as a MF? To answer that, consider the MMSE filter given by

$$\mathbf{V}_i^{\text{MMSE}} = \mathbf{G}_i^H (\mathbf{G}_i \mathbf{G}_i^H + \sigma_n^2 \mathbf{I})^{-1}.$$

Let $\mathbf{Z}_i = \mathbf{G}_i \mathbf{G}_i^H + \sigma_n^2 \mathbf{I}$ as in (4.20). Using a Neumann series expansion, its inverse can be approximated as [87–89]

$$\mathbf{Z}_i^{-1} \approx \beta_i \sum_{n=0}^{N-1} (\mathbf{I} - \beta_i \mathbf{Z}_i)^n, \quad (4.25)$$

4.3. Reducing Data Detection Complexity

where N is the number of terms used for the approximation, and β_i is a parameter that should be chosen such that the sum converges as $N \rightarrow \infty$. For guaranteed convergence, the spectral radius of $\mathbf{I} - \beta_i \mathbf{Z}_i$ must satisfy [88]

$$\rho(\mathbf{I} - \beta_i \mathbf{Z}_i) < 1, \quad (4.26)$$

where $\rho(\cdot)$ denotes the spectral radius. Since \mathbf{Z}_i is positive semi-definite, then condition (4.26) is equivalent to

$$|1 - \beta_i \lambda_{\max}(\mathbf{Z}_i)| < 1, \quad (4.27)$$

where $\lambda_{\max}(\mathbf{Z}_i)$ denotes the maximum eigenvalue of \mathbf{Z}_i . Solving (4.27) yields the following range for β_i

$$0 < \beta_i < 2/\lambda_{\max}(\mathbf{Z}_i). \quad (4.28)$$

The idea is to choose a β_i that results in a fast convergence, that is, with as few sum terms as possible. However, this can be computationally problematic, since it involves the calculation of eigenvalues. Therefore, we use a low-complexity approximation given by the trace, i.e., we choose β_i as

$$\beta_i = 2/\text{tr}(\mathbf{Z}_i). \quad (4.29)$$

Now, if we take a zeroth-order approximation, we end up with the MF. Let us instead consider a first-order approximation

$$\mathbf{Z}_i^{-1} \approx \beta_i(\mathbf{I} + (\mathbf{I} - \beta_i \mathbf{Z}_i)) = \beta_i(2\mathbf{I} - \beta_i \mathbf{Z}_i). \quad (4.30)$$

Our approximate filter then is given by

$$\mathbf{V}_i^{\text{Approx}} = \beta_i \mathbf{G}_i^H (2\mathbf{I} - \beta_i \mathbf{Z}_i). \quad (4.31)$$

Note that for modulation orders higher than 4-QAM, both the MF and approximate filter outputs need to be scaled down by the effective filter-channel gains to get correct amplitude scaling for the demapping operation. In our case, we only consider 4-QAM, and therefore we ignore the exact scaling.

Example Scenario with Two and Four Antennas

We consider a simulation scenario with a similar setup as in Table 4.2. All the UEs' channels have an RMS delay spread of 300 ns, and all the UEs are moving at a velocity of 50 km/h. We employ receiver strategy (d), in which only a single filter is calculated in every RB. We compare the case where the BS is equipped with two antennas (as in previous simulations) against the case where it is equipped with four. In each case, we compare exact MMSE filtering against the first-order approximation in (4.31) and against a MF. The x-axis this time is the average SNR of the UEs, and similarly to the previous section, a pathloss spread of ± 5 dB is applied.

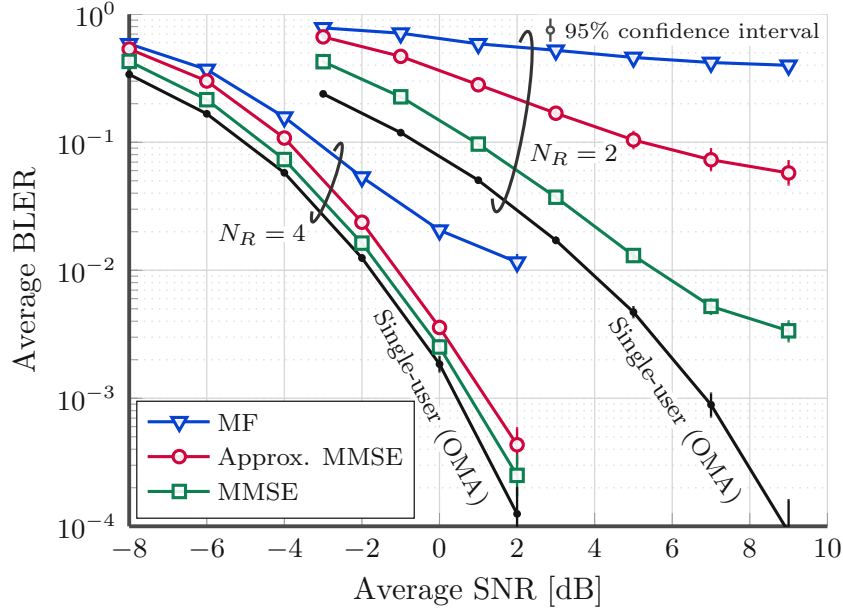


Figure 4.13: Performance of the approximate filter and the MF for the cases of $R = 2$ and 4. Strategy (d) is employed here.

Figure 4.13 shows the BLER performance of the different filters. It can be observed that for the two antennas case, a substantial gap exists between the approximate and exact filters. However, the gap becomes almost non-existent once we switch to four antennas. Similarly, for the MF, the gap is much smaller compared to the two antennas case. The interesting point here is that by adding only two additional antennas, we almost closed the gap to the exact MMSE filter. We did not need to go for a large number of antennas. This trend suggests that for even a higher number of receive antennas, a MF would be a proper replacement of the exact MMSE filter, with only a small impact on performance. Also, we see that the gap to the OMA user is reduced as well, suggesting that receiver strategy (d) now incurs even less performance loss compared to a full receiver.

Complexity-wise, it might seem that the approximate filter has a cubic complexity order, as it involves a matrix-matrix multiplication. However, since we are only interested in the filter output, it can be implemented as a cascade of two MFs [87], i.e., first apply $(2\mathbf{I} - \beta_i \mathbf{Z}_i)$, and then filter the output with $\beta_i \mathbf{G}_i^H$. This requires only matrix-vector multiplications which have a complexity order of $\mathcal{O}(L^2 N_R^2 + K L N_R)$. On the other hand, exact MMSE filtering can also be implemented through cascade filters, however, one of the cascade filters requires a matrix inverse computation, and therefore the exact filter has a complexity order of $\mathcal{O}(L^3 N_R^3 + L^2 N_R^2 + K L N_R)$. As can be seen, it is much higher compared to the first-order approximate filter, especially in the range where $K \leq L N_R$.

4.4 Final Remarks

We considered in this chapter receive-side aspects of the NOMA transmission, with the goal of describing low-complexity implementations. In the first part, we dealt with the problem of user activity detection in the context of grant-free access, where we utilized subspace methods, and investigated the influence of the channel selectivity on the performance. Although showing promising results in terms of providing reliable detection capabilities, there are still aspects that have not been considered in detail. The length of the pilot signatures considered was $L_p = 12$, similar to LTE. However, if the number of simultaneously active UEs is higher than that, then such a setup would not work anymore, and the length of the pilot sequences has to be increased. This increase, however, should be done in a way which does not cause further channel variability along the spreading interval. Spreading pilots in 2D might offer good robustness in this case. The other aspect is with respect to the number of RBs occupied. Recall that subspace methods are based on the estimated sample autocorrelation matrix. Therefore, it is important that the number of pilot-blocks over the time-frequency grid and across the different receive antennas, exceed the number of active UEs. This can be a problem, if the transmission is allocated a very small number of RBs, which could be the case for MTC. However, as long as the BS is equipped with a sufficient number of antennas, this issue can be circumvented.

In the second part, we utilized the correlation of the channel to reuse the calculated filters between neighbouring spreading blocks over the time-frequency frame. We also utilized the increase in the dimensionality stemming from the spatial domain to replace exact MMSE filtering with a lower-complexity approximation requiring no inverse calculation. Implementing this in practice requires further investigation, because how much the link performance is susceptible to those approximations depends on the used rate parameters, such as the coding rate, modulation order, and spreading length. For higher transmission rates, the residual errors due to those approximations might prevent correct decodability.

5

Controlling the Channel: RIS-Assisted Uplink NOMA

For beyond fifth-generation (B5G) wireless networks, RISs have been identified as a key technology to enhance the spectral- and energy-efficiency at low-cost. Consisting of configurable nearly-passive elements, these surfaces are capable of altering the propagation of the electromagnetic waves impinged on them, allowing to perform passive beamforming of the waves from and to a certain point, suppress interference, extend the coverage area, and more. The combination of RISs with NOMA has attracted attention recently, showing promising gains in terms of the energy efficiency, sum-rate, and outage performance. As these surfaces are capable of adjusting the received powers of the UEs, optimizing them has a direct impact on the IC detection order at the BS. Hence, identifying proper RIS configurations for such a combination with NOMA is important.

In the previous two chapters, we considered optimizing the transmitter and receiver for NOMA operation, accepting whatever the channel provides and dealing with its randomness. In this chapter, we take control of the channel by employing RISs, focusing mainly on adjusting the received powers of the UEs at the BS. First, we consider a RIS-assisted power-domain NOMA uplink with two UEs only. We attempt to characterize the statistics of the resulting propagation channel by approximating the received powers as gamma random variables. This allows us to characterize the outage performance of UEs under various RIS configurations, and to identify robust operating points. Then, in the second part, we bring the NOMA code-domain into the mix, and target scenarios supporting massive connectivity. Namely, we consider a massive MIMO cluster-based deployment, where each cluster is served by RIS in combination with code-domain NOMA.

The analysis and results in this chapter have been first published in [38–40].

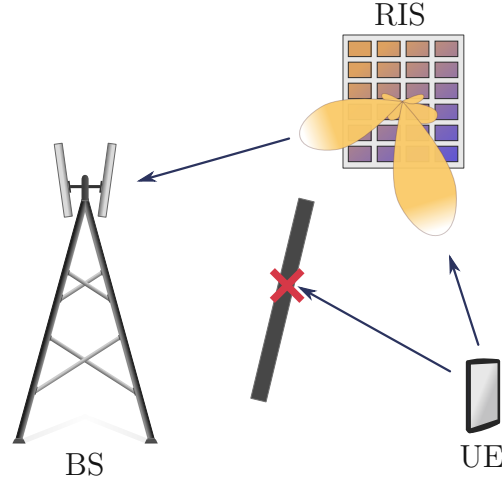


Figure 5.1: RIS-assisted transmission with blocked LOS to the BS.

5.1 Reconfigurable Intelligent Surfaces

RISs, also known as intelligent reflecting surfaces (IRSs), are surfaces consisting of a large number of elements that can actively interact with the electromagnetic waves incident upon them. By electrically controlling these elements, those surfaces can produce reflected waves with modified phase, amplitude, frequency, and/or polarization [22]. Being able to control the propagation channel can lead to favorable wireless transmission with improved spectral and energy efficiency [20,21]. The most considered form of wave-front modification is phase-shifting of the impinging waves upon the different elements. In this case, a form of passive beamforming can be achieved, which is beneficial in the context of extending the coverage area, focusing the energy towards a certain UE, reducing interference, etc [22,23]. An example application would be to provide a better coverage for a UE that has a blocked line-of-sight (LOS) to the BS, as illustrated in Figure 5.1. This is achieved by aligning the channels across the different elements in phase, resulting in a focusing of the energy from the transmitter to the receiver.

To gain further insights, let us consider the system model for single-antenna transmitter and receiver with blocked LOS, assisted by a RIS consisting of N elements. The received signal at the BS transmitted by UE- k is the sum of the signals reflected along the individual elements, i.e.,

$$y = \sqrt{\ell_{\text{BS}}\ell_{h_k}P_k} \mathbf{h}_{\text{BS}}^T \mathbf{\Phi} \mathbf{h}_k x_k + n, \quad (5.1)$$

where $\mathbf{h}_{\text{BS}} \in \mathbb{C}^N$, and $\mathbf{h}_k \in \mathbb{C}^N$ are the small-scale fading vectors of the BS-RIS and RIS-UE links, respectively. The parameters ℓ_{BS} and ℓ_{h_k} are the corresponding pathlosses, P_k and x_k are the transmit power and symbol of the k^{th} UE, and n is the zero-mean Gaussian noise with power σ_n^2 . The matrix $\mathbf{\Phi} \in \mathbb{C}^{N \times N}$ is a diagonal

5.1. Reconfigurable Intelligent Surfaces

matrix defined as

$$\Phi = \text{diag}(\eta_1 e^{j\phi_1}, \eta_2 e^{j\phi_2}, \dots, \eta_N e^{j\phi_N}), \quad (5.2)$$

where η_l and ϕ_l are the amplitude adjustment and phase-shift applied at the l^{th} -element of the RIS, respectively. Note that the RIS term is nothing more than the sum of individual channels across the elements, i.e.,

$$\mathbf{h}_{\text{BS}}^T \Phi \mathbf{h}_k = \sum_{l=1}^N \eta_l e^{j\phi_l} \mathbf{h}_{\text{BS},l} \mathbf{h}_{k,l}, \quad (5.3)$$

where $\mathbf{h}_{\text{BS},l}$ and $\mathbf{h}_{k,l}$ are the l^{th} elements of \mathbf{h}_{BS} and \mathbf{h}_k , respectively. We assume that no amplitude adjustment happens at the surface (not even material loss), and therefore we fix $\eta_1 = \eta_2 = \dots = \eta_N = 1$. We also assume that the elements of the RIS are sufficiently spaced from each other, such that uncorrelated fading across the elements holds at least approximately.

To maximize the receive power of the k^{th} UE, i.e., coherently combine its reflected waves, the phase-shifts are set as

$$\phi_l = -\arg(\mathbf{h}_{\text{BS},l} \mathbf{h}_{k,l}), \quad (5.4)$$

which can be shown by a simple application of the triangular inequality on the received amplitude. This results in the following expression for the RIS term

$$\mathbf{h}_{\text{BS}}^T \Phi \mathbf{h}_k = \sum_{l=1}^N |\mathbf{h}_{\text{BS},l}| |\mathbf{h}_{k,l}|. \quad (5.5)$$

Assuming i.i.d. $\mathbf{h}_{\text{BS},l}$ and $\mathbf{h}_{k,l}$, where $\mathbb{E}\{|\mathbf{h}_{\text{BS},l}|^2\} = \mathbb{E}\{|\mathbf{h}_{k,l}|^2\} = 1$ and $\mathbb{E}\{|\mathbf{h}_{\text{BS},l}| |\mathbf{h}_{k,l}|\} = \mu$, the mean power of the sum can be expressed as

$$\begin{aligned} \mathbb{E}\left\{\left|\sum_{l=1}^N |\mathbf{h}_{\text{BS},l}| |\mathbf{h}_{k,l}|\right|^2\right\} &= \mathbb{E}\left\{\sum_{l=1}^N |\mathbf{h}_{\text{BS},l}| |\mathbf{h}_{k,l}|\right\}^2 + \text{Var}\left\{\sum_{l=1}^N |\mathbf{h}_{\text{BS},l}| |\mathbf{h}_{k,l}|\right\} \\ &= \left(\sum_{l=1}^N \mathbb{E}\{|\mathbf{h}_{\text{BS},l}| |\mathbf{h}_{k,l}|\}\right)^2 + \sum_{l=1}^N \text{Var}\{|\mathbf{h}_{\text{BS},l}| |\mathbf{h}_{k,l}|\} \\ &= N^2 \mu^2 + N(1 - \mu^2). \end{aligned} \quad (5.6)$$

This shows a quadratic scaling of the receive power with N . Similarly, under the random combining case, in which ϕ_l is drawn randomly from the range $[0, 2\pi]$, one can show that the receive power is given by

$$\mathbb{E}\left\{\left|\sum_{l=1}^N e^{j\phi_l} \mathbf{h}_{\text{BS},l} \mathbf{h}_{k,l}\right|^2\right\} = N, \quad (5.7)$$

which shows only a linear scaling with the number of elements. The quadratic and linear gains achieved under coherent and random combining, respectively, are subject to controversy. The reason is that one may argue that the RIS channel should be normalized as part of the environment, and therefore the sum term should be scaled by $1/\sqrt{N}$. This would be the case if the RIS is to be deployed along an already existing structure in the environment, such as a building wall. In this case, a beamforming gain of only N would be observed under coherent combining, while under random combining there would almost be no gain, since the RIS would be producing random reflections, similar to what the wall beneath it would otherwise produce. In the remaining of this chapter, we will assume that the RIS occupies a new space in the environment, e.g., on a building roof-top, adding new propagation paths, and therefore the quadratic gain would make sense.

5.2 RIS-Assisted Two-User NOMA Uplink

The combination of NOMA with RISs has been gaining attention in the literature, with many works showing potential gains in terms of energy efficiency, sum-rate, and outage performance [25–31]. Given the multi-user nature of NOMA, an important aspect is how to configure the elements of the surface across the multiple UEs. In some works, the phase shifts across the different elements are set jointly according to a certain design criterion [26, 27, 90], such as maximizing the sum-rate. Other works consider the case where the entire surface is used to boost one of the NOMA UEs [25, 28]. In this section, we consider a two-UE RIS-assisted NOMA uplink, in which the elements of the RIS are split between the two NOMA UEs, i.e., part of the surface is used to coherently combine the signal of the first UE, while the other part is used to coherently combine the signal of the second one. We assume the communication to take place primarily through the surface, e.g., due to blockage of the direct links to the BS. All the links are assumed to undergo Nakagami- m fading, allowing to flexibly capture line-of-sight (LOS) and non-line-of-sight (NLOS) propagation conditions [91]. Our goal is to analyze the outage probability under NOMA interference cancellation (IC) for different splits of the elements and pathloss differences between the UEs.

5.2.1 Two-UE System Model

We consider a single-antenna two-UE NOMA uplink assisted by an N -elements RIS, as shown in Figure 5.2. At the BS, the received signal is given by

$$y = \sum_{k=1}^2 \sqrt{\ell_{\text{BS}} \ell_{h_k} P_k} \mathbf{h}_{\text{BS}}^T \Phi \mathbf{h}_k x_k + n. \quad (5.8)$$

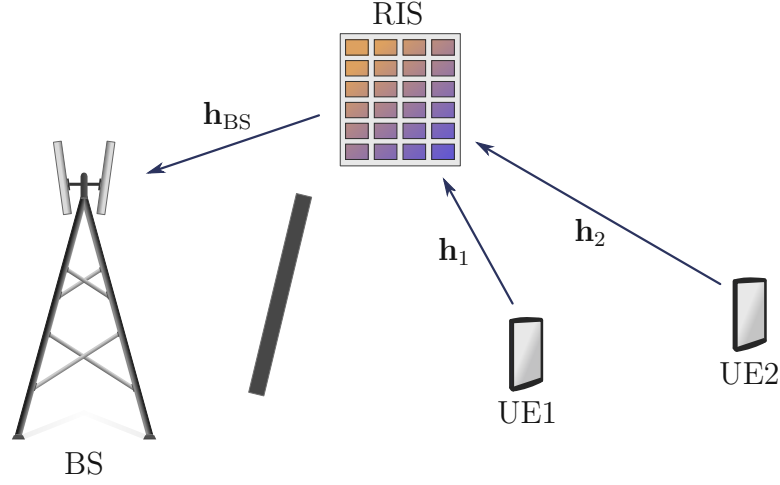


Figure 5.2: The uplink RIS-assisted NOMA setup.

To be flexible in terms of modeling LOS and NLOS propagation conditions, the links are assumed to undergo Nakagami- m fading, i.e.,

$$\begin{aligned} |\mathbf{h}_{\text{BS},l}| &\sim \text{Nakagami}(m_{\text{BS}}, 1), \\ |\mathbf{h}_{k,l}| &\sim \text{Nakagami}(m_{h_k}, 1), \end{aligned} \quad (5.9)$$

where m_{BS} and m_{h_k} are the corresponding distribution parameters.

As mentioned, we consider the case where the elements of the RIS is split between the two UEs, i.e., a total of N_1 elements are configured to coherently combine the signal of UE1, while $N_2 = N - N_1$ elements are configured for UE2. The phases are then set to

$$\phi_l = -\arg(\mathbf{h}_{\text{BS},l} \mathbf{h}_{k,l}), \quad l \in \mathcal{C}_k, \quad (5.10)$$

where \mathcal{C}_k is the set of elements that are configured to boost the k^{th} UE. Therefore, the RIS term can be written as

$$\mathbf{h}_{\text{BS}}^T \Phi \mathbf{h}_k = \underbrace{\sum_{l \in \mathcal{C}_k} |\mathbf{h}_{\text{BS},l}| |\mathbf{h}_{k,l}|}_{\text{coherently combined part of the } k^{\text{th}} \text{ UE}} + \underbrace{\sum_{l \in \bar{\mathcal{C}}_k} e^{j\phi_l} \mathbf{h}_{\text{BS},l} \mathbf{h}_{k,l}}_{\text{randomly combined part of the } k^{\text{th}} \text{ UE}}, \quad (5.11)$$

where the complement set $\bar{\mathcal{C}}_k$ is the set of elements that are not configured for the k^{th} UE, and thus will result in a random combining of its phases. Note that $\mathcal{C}_1 = \bar{\mathcal{C}}_2$, i.e., the part that will coherently combine the signal of one of the UEs, will randomly combine the signal of the other one. This is under the assumption that the channels of the two UEs are uncorrelated.

5.2.2 Outage Analysis

Our goal is to obtain expressions that describe the outage probability of the k^{th} UE in combination with IC. In the presence of the interference from the other j^{th} UE, the SINR outage for the k^{th} UE is defined as

$$p_{\text{out}}^{(k)} = \mathbb{P} \left\{ \frac{Z_k P_k}{Z_j P_j + \sigma_n^2} \leq \epsilon \right\}, \quad (5.12)$$

where Z_k , as defined below in (5.14), is the effective channel power of the k^{th} UE, and ϵ is the outage threshold. If the interference is removed via IC, then the outage is defined for the SNR as

$$p_{\text{out, SNR}}^{(k)} = \mathbb{P} \left\{ \frac{Z_k P_k}{\sigma_n^2} \leq \epsilon \right\}, \quad (5.13)$$

which is simply the cumulative distribution function (CDF) of Z_k evaluated at $\epsilon \sigma_n^2 / P_k$. In order to evaluate these probabilities, an access to the distributions of Z_k and Z_j is required, which are difficult to characterize, let alone obtaining exact closed-form expressions from them. For that reason, we resort to approximating the received powers as gamma random variables (RVs) via moment matching in a fashion similar to [92]. On the one hand, the gamma distribution encompasses many power distributions as special cases, allowing to model various propagation conditions, and on the other hand, it allows for tractability when evaluating the outage. To do so, we need access to the moments of Z_k and Z_j , for which we first need to characterize their statistics. Before going further, we state how the gamma moments matching is performed.

Lemma 5.1. Let X be a non-negative RV with first and second moments given by $\mu_X = \mathbb{E}\{X\}$ and $\mu_X^{(2)} = \mathbb{E}\{X^2\}$, respectively. The gamma RV $Y \sim \Gamma(v, \theta)$ with the same first and second moments has shape v and scale θ parameters

$$v = \frac{\mu_X^2}{\mu_X^{(2)} - \mu_X^2}, \quad \theta = \frac{\mu_X^{(2)} - \mu_X^2}{\mu_X}.$$

Additionally, gamma RVs satisfy the scaling property, in the sense that if $Y \sim \Gamma(v, \theta)$, then $cY \sim \Gamma(v, c\theta)$.

Proof. It can be found in statistics books, such as [93]. □

Statistics of the Received Power

Let Z_k be the channel power of the k^{th} UE, i.e.,

$$Z_k = \ell_{\text{BS}} \ell_{h_k} \left| \sum_{l \in \mathcal{C}_k} |\mathbf{h}_{\text{BS},l}| |\mathbf{h}_{k,l}| + \sum_{l \in \bar{\mathcal{C}}_k} e^{j\phi_l} \mathbf{h}_{\text{BS},l} \mathbf{h}_{k,l} \right|^2. \quad (5.14)$$

Our goal here is to approximate Z_k as a gamma RV via second-order moments matching, which requires an access to its first two moments. To simplify matters, we first address the statistics of the two sum terms inside.

Lemma 5.2. For the two sum terms in (5.14) given by

$$S_{\mathcal{C}_k} = \sum_{l \in \mathcal{C}_k} |\mathbf{h}_{\text{BS},l}| |\mathbf{h}_{k,l}|,$$

$$S_{\bar{\mathcal{C}}_k} = \sum_{l \in \bar{\mathcal{C}}_k} e^{j\phi_l} \mathbf{h}_{\text{BS},l} \mathbf{h}_{k,l},$$

their distributions are approximated as

$$S_{\mathcal{C}_k} \stackrel{\text{approx}}{\sim} \Gamma\left(N_k \frac{\mu_k^2}{1 - \mu_k^2}, \frac{1 - \mu_k^2}{\mu_k}\right),$$

$$S_{\bar{\mathcal{C}}_k} \stackrel{\text{approx}}{\sim} \mathcal{CN}(0, N - N_k),$$

with

$$\mu_k = \frac{\Gamma(m_{\text{BS}} + \frac{1}{2})\Gamma(m_{h_k} + \frac{1}{2})}{\Gamma(m_{\text{BS}})\Gamma(m_{h_k})(m_{\text{BS}} m_{h_k})^{1/2}},$$

where $\Gamma(\cdot)$ is the gamma function.

Proof. For $S_{\mathcal{C}_k}$, all the fading terms are positive in-phase aligned, constituting a sum of identical unit-power double-Nakagami RVs. By the causal form of the central limit theorem (CLT) [94], we can approximate the sum of positive RVs by a gamma RV via Lemma 5.1 (similar approximation for Rayleigh fading can be found in [95]). For that, we need the first and second moments of the sum. Note that the denominator of v and the numerator of θ in Lemma 5.1 are the variance, which is easier to calculate here. The mean and variance of the sum under unit-power i.i.d. conditions are given by

$$\mu_{S_{\mathcal{C}_k}} = \sum_{l \in \mathcal{C}_k} \mathbb{E}\{|\mathbf{h}_{\text{BS},l}| |\mathbf{h}_{k,l}|\} = N_k \mu_k,$$

$$\mu_{S_{\mathcal{C}_k}}^{(2)} - \mu_{S_{\mathcal{C}_k}}^2 = \sum_{l \in \mathcal{C}_k} \text{Var}\{|\mathbf{h}_{\text{BS},l}| |\mathbf{h}_{k,l}|\} = N_k(1 - \mu_k^2),$$

with

$$\mu_k = \mathbb{E}\{|\mathbf{h}_{\text{BS},l}| |\mathbf{g}_{1,n}|\} = \mathbb{E}\{|\mathbf{h}_{\text{BS},l}|\} \mathbb{E}\{|\mathbf{g}_{1,n}|\}$$

being the product of the mean of two independent Nakagami RVs. Substituting the values, we arrive at the final result. For the second term $S_{\bar{\mathcal{C}}_k}$, it consists of out-of-phase complex unit-power i.i.d. RVs, which can be approximated by a complex-

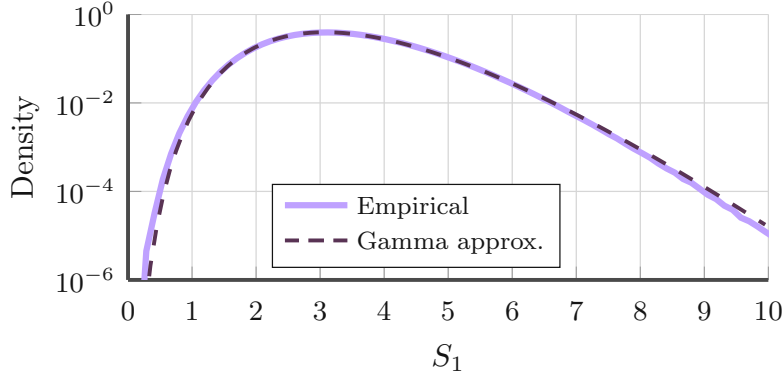
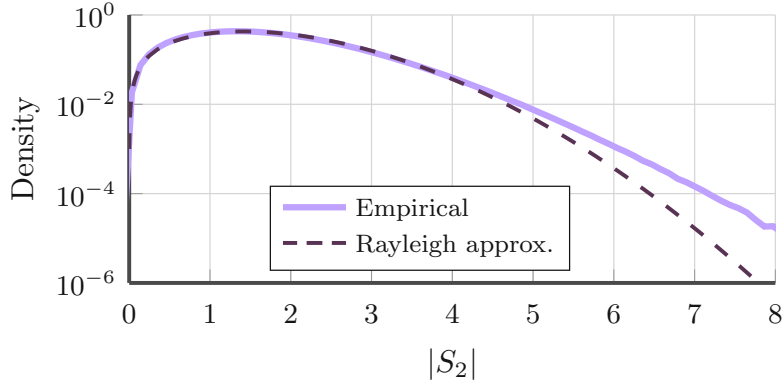

 (a) Density of S_1 .

 (b) Density of $|S_2|$.

Figure 5.3: Density of S_1 and $|S_2|$ for $N_1 = N_2 = 4$, $m_{\text{BS}} = 3$, $m_{g_1} = 1$, and ϕ_l being uniformly distributed.

Gaussian through the conventional CLT. As shown in [29], even under correlation of the real and imaginary parts, the CLT still provides a good approximation. \square

Figure 5.3 shows in log-scale the density of the sum terms obtained empirically by simulation vs. our approximations for $N_1 = N_2 = 4$. As can be seen in Figure 5.3a, the gamma approximation of S_{C_k} holds very well even in the case of only four elements, showing good match at the tails of the distribution as well. As for $S_{\bar{C}_k}$, we compare its magnitude to our Rayleigh fit (magnitude of Gaussian). We see that it does provide a good fit; however, it is not as good at the tails compared to the gamma approximation in the case before. For a larger number of elements, the assumptions based on the CLT will hold tighter and therefore the approximations will further improve.

We can now write the channel power compactly as

$$Z_k = \ell_{\text{BS}} \ell_{h_k} |S_{C_k} + S_{\bar{C}_k}|^2, \quad (5.15)$$

with the first two moments given by the following lemma.

Lemma 5.3. The first two moments of the channel power under elements splitting are given by

$$\begin{aligned}\mu_{Z_k} &= \ell_{\text{BS}} \ell_{h_k} \left(\mu_{S_{C_k}}^{(2)} + \mu_{|S_{\bar{C}_k}|}^{(2)} \right), \\ \mu_{Z_k}^{(2)} &= (\ell_{\text{BS}} \ell_{h_k})^2 \left(\mu_{S_{C_k}}^{(4)} + \mu_{|S_{\bar{C}_k}|}^{(4)} + 4 \mu_{S_{C_k}}^{(2)} \mu_{|S_{\bar{C}_k}|}^{(2)} \right),\end{aligned}$$

where

$$\begin{aligned}\mu_{S_{C_k}}^{(p)} &= \frac{\Gamma\left(N_k \frac{\mu_k^2}{1-\mu_k^2} + p\right) \left(\frac{1-\mu_k^2}{\mu_k}\right)^p}{\Gamma\left(N_k \frac{\mu_k^2}{1-\mu_k^2}\right)}, \\ \mu_{|S_{\bar{C}_k}|}^{(p)} &= \Gamma\left(1 + \frac{p}{2}\right) (N - N_k)^{p/2}.\end{aligned}$$

Proof. Expanding (5.15), we have

$$\mu_{Z_k} = \ell_{\text{BS}} \ell_{h_k} \mathbb{E}\{S_{C_k}^2 + |S_{\bar{C}_k}|^2 + 2S_{C_k} \Re\{S_{\bar{C}_k}\}\}.$$

We make the assumption here that the phase of $S_{\bar{C}_k}$ is zero-mean symmetric. This is valid since it results from an out-of-phase summation of the terms. Therefore, its real part will be zero-mean as well, leading to $\mathbb{E}\{S_{C_k} \Re\{S_{\bar{C}_k}\}\} = \mathbb{E}\{S_{C_k}\} \mathbb{E}\{\Re\{S_{\bar{C}_k}\}\} = 0$, giving the final result. We proceed in a similar manner for $\mu_{Z_k}^{(2)}$. After the expansion we get

$$\begin{aligned}\mu_{Z_k}^{(2)} &= (\ell_{\text{BS}} \ell_{h_k})^2 \mathbb{E}\{S_{C_k}^4 + |S_{\bar{C}_k}|^4 + 2S_{C_k}^2 |S_{\bar{C}_k}|^2 + 4S_{C_k}^3 \Re\{S_{\bar{C}_k}\} \\ &\quad + 4S_{C_k} |S_{\bar{C}_k}|^2 \Re\{S_{\bar{C}_k}\} + 4S_{C_k}^2 \Re\{S_{\bar{C}_k}\}^2\}.\end{aligned}$$

Following the assumptions of independence and zero-mean symmetry, we have

$$\mathbb{E}\{S_{C_k}^3 \Re\{S_{\bar{C}_k}\}\} = \mathbb{E}\{S_{C_k}^3\} \mathbb{E}\{\Re\{S_{\bar{C}_k}\}\} = 0,$$

and

$$\begin{aligned}\mathbb{E}\{S_{C_k} |S_{\bar{C}_k}|^2 \Re\{S_{\bar{C}_k}\}\} &= \mathbb{E}\{S_{C_k}\} \mathbb{E}\{|S_{\bar{C}_k}|^2 \Re\{S_{\bar{C}_k}\}\} \\ &= \mathbb{E}\{S_{C_k}\} \mathbb{E}\{\Re\{S_{\bar{C}_k}\}^3 + \Im\{S_{\bar{C}_k}\} \Re\{S_{\bar{C}_k}\}\} \\ &= 0,\end{aligned}$$

where the final result follows from the fact that the third moment is zero as well (due to symmetry), and independence between the real and imaginary parts. For the last term, we assume that the power is split equally across the real and imaginary parts,

and therefore

$$\mathbb{E}\{S_{\mathcal{C}_k}^2 \Re\{S_{\mathcal{C}_k}\}^2\} = \mathbb{E}\{S_{\mathcal{C}_k}^2\} \mathbb{E}\{|S_{\mathcal{C}_k}|^2\}/2.$$

We get the final results by collecting the terms back and substituting the moments of gamma and Rayleigh (magnitude of Gaussian) RVs. \square

Finally, after scaling with the transmit power, the received power of the k^{th} UE is expressed as

$$Z_k P_k \stackrel{\text{approx}}{\sim} \Gamma(v_k, P_k \theta_k), \quad (5.16)$$

where v_k and θ_k are the gamma parameters matched to the moments in Lemma 5.3.

Outage Probability under Interference Cancellation

Next, we calculate the outage probability for the uplink RIS-NOMA system under IC. First, we evaluate the outage probability without IC.

Proposition 5.1. Let $Z_k P_k \sim \Gamma(v_k, P_k \theta_k)$ be the received power of the k^{th} UE, $Z_j P_j \sim \Gamma(v_j, P_j \theta_j)$ the received power of the j^{th} UE, with σ_n^2 being the noise power. The RIS-NOMA outage probability without IC is given by

$$p_{\text{out}}^{(k)} \approx I\left(\frac{\epsilon \hat{\theta}_j}{\hat{\theta}_k + \epsilon \hat{\theta}_j}; \hat{v}_k, \hat{v}_j\right),$$

where $I(., ., .)$ is the regularized incomplete beta function, and

$$\hat{v}_k = v_k, \quad \hat{\theta}_k = \theta_k P_k,$$

$$\hat{v}_j = \frac{(v_j \theta_j P_j + \sigma_n^2)^2}{v_j (\theta_j P_j)^2}, \quad \hat{\theta}_j = \frac{v_j (\theta_j P_j)^2}{v_j \theta_j P_j + \sigma_n^2}.$$

Proof. Let $X \sim \Gamma(v_X, \theta_X)$ and $Y \sim \Gamma(v_Y, \theta_Y)$ be two independent gamma RVs, then their ratio $R = X/Y$ is known to be beta prime distributed, i.e., $R \sim \beta'(v_X, v_Y, 1, \theta_X/\theta_Y)$, with its CDF given by

$$\mathbb{P}\{R \leq \epsilon\} = I\left(\frac{\epsilon \theta_Y}{\theta_X + \epsilon \theta_Y}; v_X, v_Y\right).$$

However, the denominator in (5.12) is not gamma distributed, due to the presence of the noise term. Therefore, we approximate the interference-plus-noise term by an equivalent gamma RV, again, via moments matching. By doing so, and using the gamma scaling property, we arrive at the final results. \square

The detection scheme we consider here is parallel, in the sense that UE1, UE2, or both can be detected correctly at the first iteration and removed from the received

signal. Whatever remains can be detected in the second iteration after IC. Such a formulation allows us to assume an arbitrary cancellation order and saves us the hassle of order statistics as would be required under successive IC. This is formulated in the following proposition.

Proposition 5.2. The RIS-NOMA outage probability of the k^{th} UE under IC is given by

$$p_{\text{out,IC}}^{(k)} \approx 1 - \min(p_{\text{succ}}^{(k)} + p_{\text{succ}}^{(j)} p_{\text{succ,SNR}}^{(k)}, p_{\text{succ,SNR}}^{(k)}), \quad (5.17)$$

where $p_{\text{succ}}^{(k)} = 1 - p_{\text{out}}^{(k)}$ and $p_{\text{succ,SNR}}^{(k)} = 1 - p_{\text{out,SNR}}^{(k)}$ are the success probabilities, with $p_{\text{out,SNR}}^{(k)}$ as defined in (5.13) is the gamma CDF given by

$$p_{\text{out,SNR}}^{(k)} = \gamma\left(v_k, \frac{\epsilon \sigma_n^2}{\theta_k P_k}\right) \quad (5.18)$$

with $\gamma(\cdot, \cdot)$ being the regularized incomplete gamma function.

Proof. There are two paths for a successful detection of the k^{th} UE: it is detected correctly in the first iteration; or, it is not, but the other UE is detected correctly, and after IC, the k^{th} UE is detected interference-free in the presence of noise only. Following these events, we can approximate the success probability under IC as $p_{\text{succ,IC}}^{(k)} \approx p_{\text{succ}}^{(k)} + p_{\text{succ}}^{(j)} p_{\text{succ,SNR}}^{(k)}$. The detection sequence just mentioned is not of fully independent events; hence, the approximation sign. To further improve the approximation, we make use of the fact that the performance cannot be better than that of the interference-free noise-only case. We get the final results by taking the minimum between that expression and the noise-only case. \square

5.2.3 Analysis of an Example Scenario

We consider a scenario where the communication between the NOMA UEs and the BS takes place through a 32-elements RIS, and evaluate the outage performance using (5.17). We assume the RIS to have a strong LOS connection to the BS, while the UEs have moderate LOS to the RIS, with UE1 having a stronger LOS than UE2. This is set by adjusting the corresponding Nakagami m parameters. The pathloss of UE1 is fixed to -70 dB, while for UE2, it varies. Without loss of generality, we assume that both UEs are transmitting with the same power, i.e., $P_1 = P_2$. In practice, the UEs might transmit with different powers; however, that does not affect the validity of our analysis here. It holds for any choice of P_1 and P_2 . The simulation parameters are summarized in Table 5.1.

Impact of Elements Splitting

We define the split factor α as the percentage of elements that are allocated for the coherent combining of the signal of UE1. Given that N_1 elements are allocated to

5.2. RIS-Assisted Two-User NOMA Uplink

Parameter	Value
RIS elements	$N = 32$
Transmit powers	$P_1 = P_2 = 30$ dBm
Nakagami parameters	$m_{\text{BS}} = 6$ $m_{h_1} = 3, m_{h_2} = 1.5$
Pathlosses	$\ell_{\text{BS}} = -65$ dB $\ell_{h_1} = -70$ dB, ℓ_{h_2} is variable
Noise power	$\sigma_n^2 = -110$ dBm

Table 5.1: Simulation parameters for the elements' split scenario.

UE1, the split factor then is defined as

$$\alpha = N_1/N, \quad (5.19)$$

and thus the number of elements allocated for UE2 is

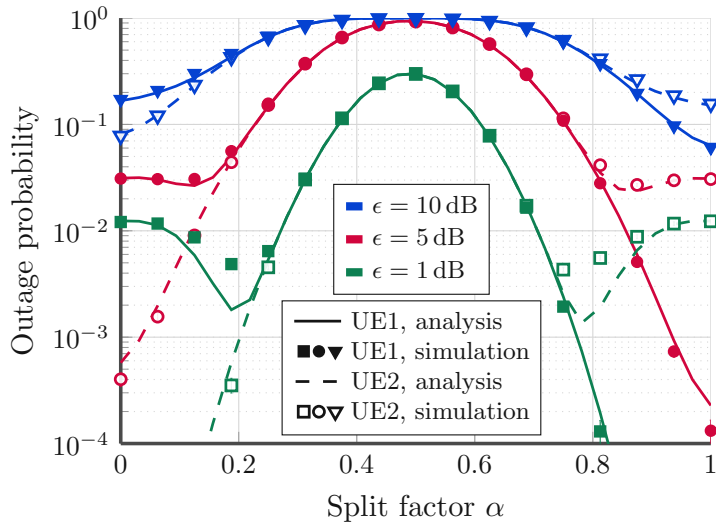
$$N_2 = N - \lceil \alpha N \rceil. \quad (5.20)$$

When $\alpha = 1$, all the elements are allocated to UE1, while for $\alpha = 0$, all the elements are allocated to UE2, etc.

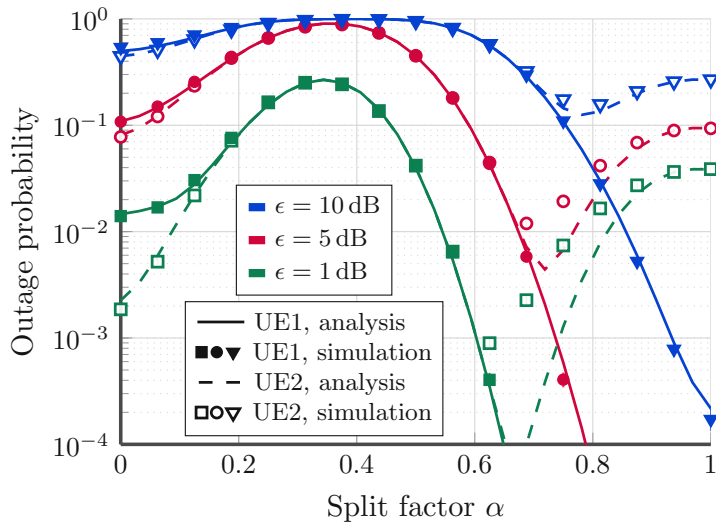
We investigate the outage performance over the split factor α for different values of the outage threshold ϵ . Figure 5.4a shows the performance when the UEs have equal pathloss. In this case, the outage probability for both UEs is minimized, if most of the elements are configured to boost one of the UEs. The reason for this is, when the split is close to 50%, then it is more likely that both UEs will be received with similar strength at the BS. This, in turn, makes IC more difficult, since the UEs would suffer strong interference from each other. Therefore, when the pathloss difference between the two UEs is small, it makes sense to focus on boosting one of the UEs, such that the power gap between them increases, allowing the stronger UE to be detected correctly with high probability at the first IC iteration. This should be done such that the weaker UE still gets assigned a sufficient number of elements, ensuring that its signal is also sufficiently boosted. For example, in Figure 5.4a for $\epsilon = 1$ dB, operating around a split of 0.2 or 0.8 offers a balanced setup for both UEs.

Figure 5.4b and Figure 5.4c show the outage performance when the pathloss of UE2 is 5 dB and 10 dB higher than UE1, respectively. We observe that as the gap increases, and at low outage thresholds, the split moves towards boosting UE2. In this case, the two UEs have a natural power gap due to the pathloss difference, and therefore the RIS can be used to enhance the performance of the weaker user (UE2). It can also be observed that as the gap increases, better performance is achieved for both UEs. This indicates that when it comes to NOMA user pairing, the BS should avoid pairing users with similar pathlosses. However, the weak UE should be strong enough such that after the combining at the surface, it is able to overcome the noise at the BS receiver.

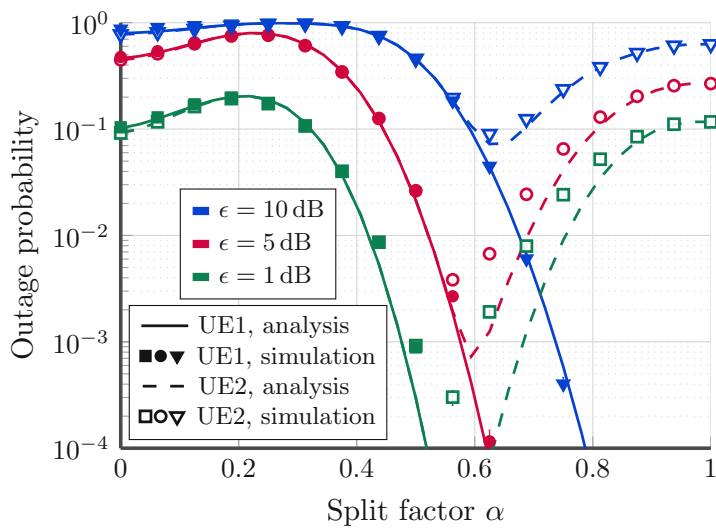
5.2. RIS-Assisted Two-User NOMA Uplink



(a) Equal pathloss case.



(b) UE2 5 dB weaker.



(c) UE2 10 dB weaker.

Figure 5.4: Outage probability vs. the split factor for different outage thresholds.

Regarding the accuracy of the analysis, we observe that the approximations hold well for the strong UE. As for the weak UE, and at low outage thresholds, a relatively large gap exists between analysis and simulation for some values of the split factor, suggesting that the gamma approximation of the power does not hold well under such splitting conditions.

Selection of the Split Factor

We consider the selection of the split factor α from a robust perspective. To ensure that boosting the performance of one UE does not come at the cost of degrading the performance of the other one, the split factor is chosen according to

$$\alpha_{\text{robust}} = \arg \min_{\alpha} \max_k p_{\text{out,IC}}^{(k)}. \quad (5.21)$$

In Figure 5.4a to 5.4c, this would correspond to the points where the UE2 outage diverges from UE1 and starts saturating (on the right side). However, at high outage thresholds, it can be observed that the outage probability of UE2 is very high, no matter what split is applied. For that reason, we introduce the notion of a limiting threshold λ . If the weak UE outage probability is higher than λ , then the entire RIS is used to boost the strong UE, as allocating elements to the weak UE would be a waste of the surface elements. Assuming UE2 is the weaker UE, (5.21) is modified as follows

$$\alpha_{\text{robust}} = \begin{cases} \arg \min_{\alpha} \max_k p_{\text{out,IC}}^{(k)}, & \text{if } p_{\text{out,IC}}^{(2)} < \lambda, \\ 1, & \text{otherwise.} \end{cases} \quad (5.22)$$

Although the analysis results shown in the previous subsection is not very accurate at low outage thresholds, it can be seen that the robust point occurs almost at the same α for both analysis and simulation. This motivates the use of the analysis as a method to determine α_{robust} . Solving (5.22) in closed-form is difficult due to the complexity of the functions involved. We thus rely on performing a search for determining the optimal point. Recall that $\alpha = N_1/N$ with $N_1 = 1, 2, \dots, N$ (i.e., the maximum number of possibilities is N), meaning that the search can be performed quickly.

Figure 5.5 shows α_{robust} obtained by search through exhaustive simulations vs. analysis for different pathloss gaps between the two UEs. We observe that at low pathloss gaps, the split is chosen to boost UE1 (the stronger UE), since it improves the performance of the NOMA IC. As the pathloss of UE2 increases, the robust RIS strategy attempts to compensate for the high pathloss by allocating more elements to UE2. The sudden jumps to 1 are due to the limiting threshold, which is set to $\lambda = 10^{-1}$ here. This indicates that at those outage thresholds, the performance of UE2 is unacceptable anyway, that it is better to use the entire surface to boost UE1. Also, at low outage thresholds, we observe that a split close to 50% seems to be the robust selection, while at high outage thresholds, the selection across the

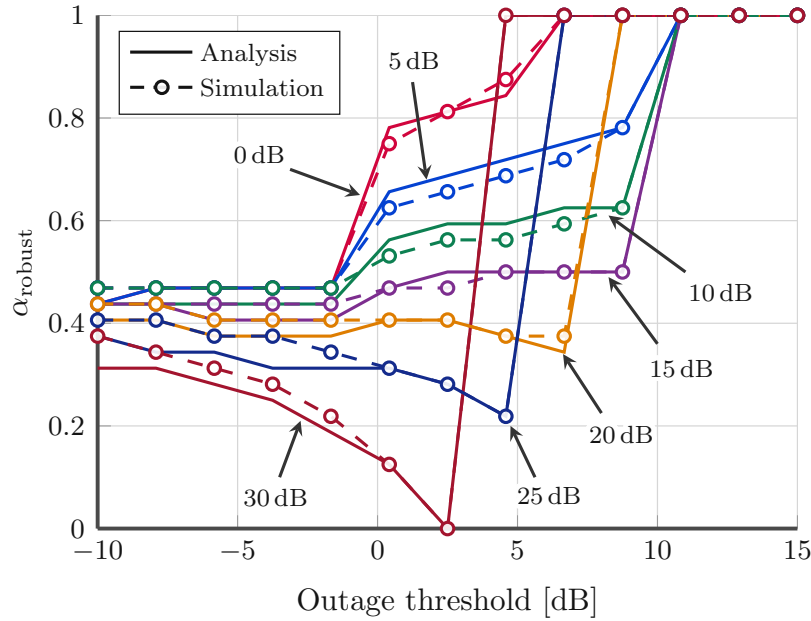


Figure 5.5: Robust selection of the split factor for different pathloss gaps (search via analysis vs. exhaustive simulations).

different gaps can vary substantially. Note that at low outage thresholds, the outage probability is very low, which makes the simulation-based selection inaccurate, since it would require a huge number of simulation samples. This is the advantage of the analytical based approach, since it can predict the performance, even at very low outage probabilities.

5.3 Combination with Code-Domain NOMA

In the previous part, we considered the combination of RISs with power-domain NOMA for two UEs. We saw how the RIS configuration can heavily impact the NOMA IC performance, as it can alter the receive powers of the UEs at the BS. We now investigate the combination of a RIS with code-domain NOMA, in the context of a cluster-based massive MIMO deployment. As it is likely that those surfaces would be deployed at rooftops, we make the assumption that each cluster is served by a RIS having a strong LOS connection to the BS. The BS forms beams towards the clusters' RISs, allowing to simultaneously boost the received power of the target cluster and suppressing inter-cluster interference, as depicted in Figure 5.6. In order to support massive connectivity, code-domain NOMA via short spreading is employed in each cluster. At the BS, and after spatial filtering, MMSE-IC detection is carried out to detect the NOMA UEs. The question then is, how to configure the RIS such that a large number of UEs is supported?

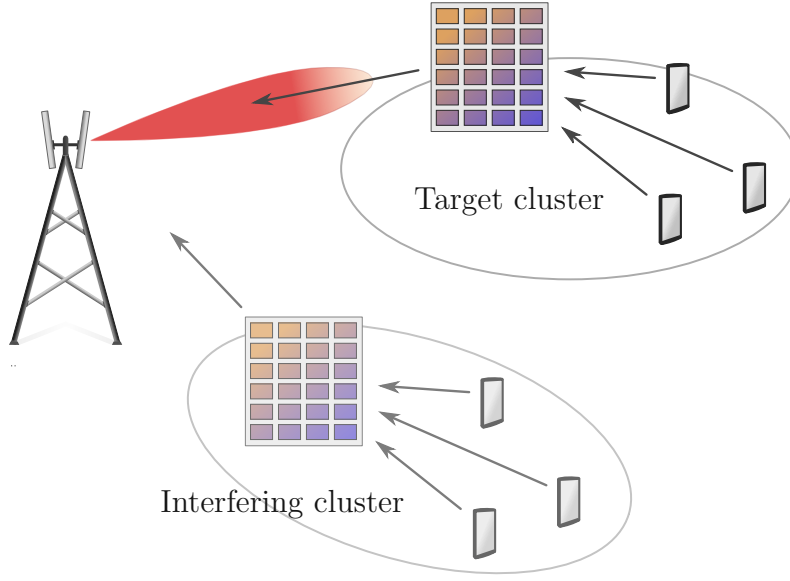


Figure 5.6: The cluster-based RIS-assisted NOMA uplink.

5.3.1 System Model Combined with Code-Domain NOMA

We consider a cluster of K single-antenna UEs communicating with an N_R -antennas BS through an N -elements RIS. Due to blockage, we assume the communication to take place primarily through the RIS, and therefore we drop the direct paths between the UEs and the BS. We will further justify this assumption later. With each UE transmitting using a short spreading signature \mathbf{s}_k , the received signal at the BS, $\mathbf{y} \in \mathbb{C}^{LN_R \times 1}$, is given by

$$\mathbf{y} = \sum_{k=1}^K \sqrt{\ell_{\text{BS}} \ell_{h_k} P_k L} (\mathbf{H}_{\text{BS}} \Phi \mathbf{h}_k \otimes \mathbf{s}_k) x_k + \mathbf{z} + \mathbf{n}, \quad (5.23)$$

where the matrix $\mathbf{H}_{\text{BS}} \in \mathbb{C}^{N_R \times N}$ represents the small-scale MIMO fading channel of the BS-RIS link. The term $\mathbf{z} \in \mathbb{C}^{LN_R \times 1}$ is the sum of all received signals from outside the intended cluster, i.e., inter-cluster interference. Since the RIS is deployed in a LOS to the BS, the BS-RIS channel is rank-1 (assuming far-field with respect to the BS), and is given by $\mathbf{H}_{\text{BS}} = \mathbf{a} \mathbf{b}^H$, where \mathbf{a} and \mathbf{b} are the array responses at the BS and RIS, respectively. The received signal can then be written as

$$\mathbf{y} = (\mathbf{a} \otimes \mathbf{I}_L) \sum_{k=1}^K \sqrt{\ell_{\text{BS}} \ell_{h_k} P_k L} (\mathbf{b}^H \Phi \mathbf{h}_k \otimes \mathbf{s}_k) x_k + \mathbf{z} + \mathbf{n}. \quad (5.24)$$

Equipped with a large number of antennas, the BS forms a beam towards the cluster's RIS, boosting the received power on the one hand, and on the other hand, suppressing inter-cluster interference. To achieve this, beam forming via maximum-

ratio combining (MRC) is performed. This further justifies the assumption of dropping the direct path between the UEs and the BS, as it would be even weaker after beamforming. The MRC spatially filtered signal $\tilde{\mathbf{y}} = (\mathbf{a}^H \otimes \mathbf{I}_L)\mathbf{y}$ is given by

$$\tilde{\mathbf{y}} = (\mathbf{a}^H \mathbf{a} \otimes \mathbf{I}_L) \sum_{k=1}^K \sqrt{\ell_{\text{BS}} \ell_{h_k} P_k L} (\mathbf{b}^H \Phi \mathbf{h}_k \otimes \mathbf{s}_k) x_k + (\mathbf{a}^H \otimes \mathbf{I}_L) \mathbf{z} + (\mathbf{a}^H \otimes \mathbf{I}_L) \mathbf{n}. \quad (5.25)$$

With the beamforming towards the target RIS, inter-cluster interference is greatly reduced, i.e., $(\mathbf{a}^H \otimes \mathbf{I}_L)\mathbf{z} \approx \mathbf{0}$. Since we have $\mathbf{a}^H \mathbf{a} = N_R$, and letting $\tilde{\mathbf{n}} = (\mathbf{a}^H \otimes \mathbf{I}_L)\mathbf{n}$ be the spatially filtered noise with $\sigma_{\tilde{\mathbf{n}}}^2 = N_R \sigma_{\mathbf{n}}^2$, (5.25) is further developed as

$$\tilde{\mathbf{y}} = \sum_{k=1}^K \sqrt{N_R^2 \ell_{\text{BS}} \ell_{h_k} P_k L} (\mathbf{b}^H \Phi \mathbf{h}_k \otimes \mathbf{s}_k) x_k + \tilde{\mathbf{n}}. \quad (5.26)$$

Notice that $\mathbf{b}^H \Phi \mathbf{h}_k$ is a scalar and therefore \otimes is no longer necessary. Let $\beta_k = \sqrt{N_R^2 \ell_{\text{BS}} \ell_{h_k} P_k L}$, $\mathbf{w} = \text{diag}(\Phi^H)$, and $\hat{\mathbf{h}}_k = \mathbf{b}^* \circ \mathbf{h}_k$, where \circ denotes the Hadamard product, the post-spatially filtered signal can finally be written as

$$\tilde{\mathbf{y}} = \sum_{k=1}^K \beta_k (\mathbf{w}^H \hat{\mathbf{h}}_k) \mathbf{s}_k x_k + \tilde{\mathbf{n}}. \quad (5.27)$$

In order to detect the UEs within the cluster, the BS performs MMSE-IC detection, with a UE being detected correctly if its SINR exceeds a certain rate threshold. Assuming, for simplicity, a successive IC in which one UE is detected per IC stage, and assuming a detection order of UE1, UE2, ..., UEK, the post-filtering SINR of the k^{th} UE is given by

$$\text{SINR}_k = \frac{|\beta_k (\mathbf{w}^H \hat{\mathbf{h}}_k) \mathbf{v}_k^H \mathbf{s}_k|^2}{\sum_{l=k+1}^K |\beta_l (\mathbf{w}^H \hat{\mathbf{h}}_l) \mathbf{v}_k^H \mathbf{s}_l|^2 + \sigma_{\tilde{\mathbf{n}}}^2 \|\mathbf{v}_k\|^2}, \quad (5.28)$$

where \mathbf{v}_k is the MMSE filter applied at the k^{th} stage. As can be seen, every time a UE is removed, the next UE in the next IC stage experiences less interference, until we reach the last UE, in which it only has to deal with noise. The goal now is to design \mathbf{w} such that

$$\text{SINR}_k \geq \epsilon_k, \quad \forall k, k = 1, 2, \dots, K, \quad (5.29)$$

where ϵ_k is the detection (or outage) threshold of the k^{th} UE. In other words, we choose the phase-shifts at the RIS such that the power gaps between the UEs combined with the MMSE filtering and IC result in SINRs exceeding the required threshold for decodability, at each of the IC stages.

We have multiple problems here; first, the detection order of the UEs, to begin

with, is unknown and it depends on the choice of \mathbf{w} . This can be clearly seen in (5.27), where the received power of the users is directly impacted by the choice of \mathbf{w} . In other words, the optimal detection order and \mathbf{w} need to be determined jointly, requiring a search over all possible detection orders, which can be of prohibitive complexity for large K ; second, the resulting SINR at each stage depends on the MMSE filter \mathbf{v}_k ; however, \mathbf{v}_k also depends on \mathbf{w} and the detection order (coupled), and therefore determining \mathbf{w} depends on the resulting \mathbf{v}_k ; third, even if everything is known, how do we find a \mathbf{w} satisfying all of the K inequalities in (5.29)?

5.3.2 Sum-Rate Optimized Phase-Shifts

It is known from the MIMO literature that MMSE-IC is a sum-rate optimal detection scheme [96]. Therefore, one way to avoid the aforementioned problems with the detection order and the choice of the MMSE filter, is to optimize \mathbf{w} such that the sum-rate of the cluster is maximized. To that end, the sum-rate is given by

$$R_{\text{sum}} = \frac{1}{L} \log_2 \det \left(\mathbf{I}_L + \frac{1}{\sigma_{\hat{\mathbf{n}}}^2} \sum_{k=1}^K \beta_k^2 (\mathbf{w}^H \hat{\mathbf{h}}_k) \mathbf{s}_k \mathbf{s}_k^H (\hat{\mathbf{h}}_k^H \mathbf{w}) \right). \quad (5.30)$$

Due to the determinant operator $\det(\cdot)$ and the $\mathbf{s}_k \mathbf{s}_k^H$ term, maximizing the above sum-rate expression is not an easy task. To manage that, we drop the spreading, and optimize the system as if no spreading is employed, i.e., we set $L = 1$ and $\mathbf{s}_k = 1, \forall k$. Such an optimization would correspond to a pure power-domain NOMA system, i.e., a worst-case scenario in which the spreading has no impact. Then, (5.30) becomes

$$R_{\text{sum}}^{(\text{no spread.})} = \log_2 \left(1 + \frac{1}{\sigma_{\hat{\mathbf{n}}}^2} \sum_{k=1}^K \beta_k^2 \mathbf{w}^H \hat{\mathbf{h}}_k \hat{\mathbf{h}}_k^H \mathbf{w} \right). \quad (5.31)$$

Let $\mathbf{H} = \sum_{k=1}^K \beta_k^2 \hat{\mathbf{h}}_k \hat{\mathbf{h}}_k^H$, the sum-rate maximizer is given by

$$\begin{aligned} \mathbf{w}_{\text{sum}} &= \arg \max_{\mathbf{w}} \mathbf{w}^H \mathbf{H} \mathbf{w} \\ \text{s.t. } & |[\mathbf{w}]_n| = 1, \quad n = 1, 2, \dots, N, \end{aligned} \quad (5.32)$$

where the condition $|[\mathbf{w}]_n| = 1$ refers to the n^{th} element of \mathbf{w} performing a phase-shift only. In order to solve (5.32), we relax it to a conventional quadratic problem. Therefore, the maximizer of $\mathbf{w}^H \mathbf{H} \mathbf{w}$ is given by the eigenvector of \mathbf{H} corresponding to its maximum eigenvalue. Let \mathbf{u}_{max} be that eigenvector, the elements of \mathbf{w}_{sum} are then set to

$$[\mathbf{w}_{\text{sum}}]_n = \exp(j\angle[\mathbf{u}_{\text{max}}]_n), \quad n = 1, 2, \dots, N, \quad (5.33)$$

i.e., \mathbf{w}_{sum} is set such that it performs the same phase-shifts as \mathbf{u}_{max} .

The issue with the sum-rate optimized shifts is that if the UEs have similar

receive powers, then the RIS would boost all of them by an equal amount, i.e., it only provides a SNR gain (the strongest eigenvector would point in the direction that favors all the UEs). This is beneficial if the system suffers from low SNR; however, our major problem here is multi-user interference, and the goal is to boost the UEs with different portions, such that sufficient power gaps are created between them, allowing the IC to operate successfully. Also, in our optimization above, spreading is not taken into account. However, if the UEs have sufficient power gaps between them (e.g., due to different pathlosses), then \mathbf{w}_{sum} can provide a good solution, as the strongest eigenvector would point in the direction of the strongest UEs, and this helps to further enlarge the gaps (the RIS would boost the stronger UEs further), resulting in better sequential SINRs under IC. We will see this effect later in Section 5.3.4.

5.3.3 Proposed Optimization Approach

Robust optimization of the phase-shifts requires solving the inequalities of (5.29). However, as we mentioned before, the optimal solution is difficult to obtain, due to the coupling between the detection order and the MMSE filter with the RIS phase-shifts. In the following, we propose a suboptimal procedure that allows us to obtain a solution to the problem.

Detection Order

The optimal solution requires an exhaustive search over all possible detection orders, consisting of $K!$ possibilities. This can be prohibitive for large K , and it is the large K that we are interested in. A suboptimal approach that can provide a good performance [27], is to order the UEs based on their received signal strength, i.e., $|\beta_k \mathbf{w}^H \hat{\mathbf{h}}_k|$. However, we can see that it depends on \mathbf{w} , which we seek to find in the first place. For that reason, we do the ordering based on the sum-rate optimized shifts, i.e., by ordering the UEs according to $|\beta_k \mathbf{w}_{\text{sum}}^H \hat{\mathbf{h}}_k|$. In other words, \mathbf{w}_{sum} is employed as the initial solution for determining the detection order. In the following, and without loss of generality, we assume the resultant UEs' ordering is

$$|\beta_1 \mathbf{w}_{\text{sum}}^H \hat{\mathbf{h}}_1| \geq |\beta_2 \mathbf{w}_{\text{sum}}^H \hat{\mathbf{h}}_2| \geq \dots \geq |\beta_K \mathbf{w}_{\text{sum}}^H \hat{\mathbf{h}}_K|, \quad (5.34)$$

that is, after ordering, UE 1 is the strongest user, while UE K is the weakest one. This assumption is only applied to simplify notation for the next parts.

MMSE Filtering

The next coupled variable is the MMSE filter. We follow a similar approach as with the detection order. We calculate the MMSE filters based on the sum-rate solution. Therefore, given our determined detection order and \mathbf{w}_{sum} , the MMSE filters applied

in (5.28) are such that

$$\mathbf{v}_k^H = \mathbf{g}_k^H \left(\sum_{l=k}^K \mathbf{g}_l \mathbf{g}_l^H + \mathbf{I}_L \sigma_{\mathbf{n}}^2 \right)^{-1}, \quad (5.35)$$

where $\mathbf{g}_k = \beta_k (\mathbf{w}_{\text{sum}}^H \hat{\mathbf{h}}_k) \mathbf{s}_k$.

Phase-Shifts Optimization

Having both the detection order and MMSE filter determined based on \mathbf{w}_{sum} , we now proceed to finding our final phase-shifts. First, we rewrite (5.28) as

$$\frac{\mathbf{w}^H \left(\beta_k^2 |\mathbf{v}_k^H \mathbf{s}_k|^2 \hat{\mathbf{h}}_k \hat{\mathbf{h}}_k^H \right) \mathbf{w}}{\mathbf{w}^H \left(\sum_{l=k+1}^K \beta_l^2 |\mathbf{v}_k^H \mathbf{s}_l|^2 \hat{\mathbf{h}}_l \hat{\mathbf{h}}_l^H + \frac{\sigma_{\mathbf{n}}^2 \|\mathbf{v}_k\|^2}{N} \mathbf{I}_N \right) \mathbf{w}}, \quad (5.36)$$

where the fact that $\mathbf{w}^H \mathbf{w} = N$ has been applied to the noise term. Let

$$\begin{aligned} \mathbf{A}_k &= \beta_k^2 |\mathbf{v}_k^H \mathbf{s}_k|^2 \hat{\mathbf{h}}_k \hat{\mathbf{h}}_k^H, \\ \mathbf{B}_k &= \sum_{l=k+1}^K \beta_l^2 |\mathbf{v}_k^H \mathbf{s}_l|^2 \hat{\mathbf{h}}_l \hat{\mathbf{h}}_l^H + \frac{\sigma_{\mathbf{n}}^2 \|\mathbf{v}_k\|^2}{N} \mathbf{I}_N, \end{aligned} \quad (5.37)$$

our optimization problem is then formulated as

$$\begin{aligned} &\text{find } \mathbf{w} \\ &\text{s.t. } \frac{\mathbf{w}^H \mathbf{A}_k \mathbf{w}}{\mathbf{w}^H \mathbf{B}_k \mathbf{w}} \geq \epsilon_k, \quad k = 1, 2, \dots, K, \\ & \quad |\mathbf{w}_n| = 1, \quad n = 1, 2, \dots, N. \end{aligned} \quad (5.38)$$

To find a solution to these series of inequalities, we relax (5.38) into a semidefinite programming (SDP) problem, which can be solved efficiently using convex optimization algorithms [97]. Let $\mathbf{W} = \mathbf{w} \mathbf{w}^H$; using the trace operator, we have $\mathbf{w}^H \mathbf{A}_k \mathbf{w} = \text{tr}(\mathbf{A}_k \mathbf{w} \mathbf{w}^H) = \text{tr}(\mathbf{A}_k \mathbf{W})$. Similarly, we have $\mathbf{w}^H \mathbf{B}_k \mathbf{w} = \text{tr}(\mathbf{B}_k \mathbf{W})$. The SINR condition is then written as

$$\begin{aligned} \text{tr}(\mathbf{A}_k \mathbf{W}) - \epsilon_k \text{tr}(\mathbf{B}_k \mathbf{W}) &\geq 0, \\ \text{tr}([\mathbf{A}_k - \epsilon_k \mathbf{B}_k] \mathbf{W}) &\geq 0. \end{aligned} \quad (5.39)$$

Finally, our SDP-relaxed problem is given by

$$\begin{aligned}
 & \text{find } \mathbf{W} \\
 & \text{s.t. } \text{tr}\left([\mathbf{A}_k - \epsilon_k \mathbf{B}_k] \mathbf{W}\right) \geq 0, \quad k = 1, 2, \dots, K, \\
 & \quad \mathbf{W} \succeq 0, [\mathbf{W}]_{n,n} = 1, \quad n = 1, 2, \dots, N,
 \end{aligned} \tag{5.40}$$

where $[\mathbf{W}]_{n,n}$ is the n^{th} diagonal element of \mathbf{W} . In this work, the optimizer used is based on CVX [98]. If a solution is found, then we set \mathbf{w}_{prop} (proposed) such that it performs the same phase-shifts as the eigenvector of \mathbf{W} corresponding to its maximum eigenvalue (best rank-1 approximation) in a similar fashion as in (5.33). If no solution is feasible, then we rely on the sum-rate solution, i.e., we set $\mathbf{w}_{\text{prop}} = \mathbf{w}_{\text{sum}}$.

5.3.4 Investigation of an Example Scenario

We consider a scenario where K active UEs in the target cluster communicate with a 32-antennas BS through a 32-elements RIS. A 4×16 Grassmannian codebook is employed for the spreading. The BS-RIS channel is LOS with pathloss $\ell_{\text{BS}} = -65$ dB, while the RIS-UE channels are modeled as Rayleigh fading with the pathloss uniformly disturbed as $\ell_{h_k} \sim \mathcal{U}(-65 - s, -65 + s)$ dB, i.e., a mean component of -65 dB plus a spread of $\pm s$. By adjusting s , we can control the pathloss differences across the UEs, and thus the average received power difference between the UEs at the BS. We assume the UEs to transmit with an equal power of $P_k = P = 30$ dBm, $\forall k$. The noise power is set to $\sigma_{\mathbf{n}}^2 = -110$ dBm. Also, we assume all the UEs to have the same outage threshold of $\epsilon_k = \epsilon$.

In Figure 5.7, we compare the detection performance using our proposed approach versus the sum-rate optimized phase-shifts and random ones. The results are shown for a pathloss spread of ± 3 dB, and over the outage thresholds of $\epsilon = 1, 4,$ and 9 dB. The desired result here is a 1:1 line, i.e., all active UEs are detected correctly. We observe that our proposed RIS adaptation allows for a substantial increase of the number of correctly detected users. This is achieved at both low and high outage thresholds. As the threshold increases, it becomes more challenging for the RIS to satisfy all the inequalities of (5.29). If the threshold is too high for the number of active UEs, then no feasible solution would be possible, and the number of correctly detected UEs begins to drop.

Next, we set $K = 12$ and $\epsilon = 4$ dB and investigate the performance over different pathloss spreads. We also compare our NOMA (Grassmannian) codebook to an OMA codebook, e.g., a 4×4 identity matrix. In the case of the OMA codebook, we only have 4 signatures, and therefore unique signature assignment to the 12 UEs is not possible, i.e., the orthogonal signatures must be reused between the users. The results are shown in Figure 5.8. We observe that, at least for the considered configuration, the NOMA codebook offers substantial improvement over the OMA codebook reuse strategy, across the different adaptation approaches. The

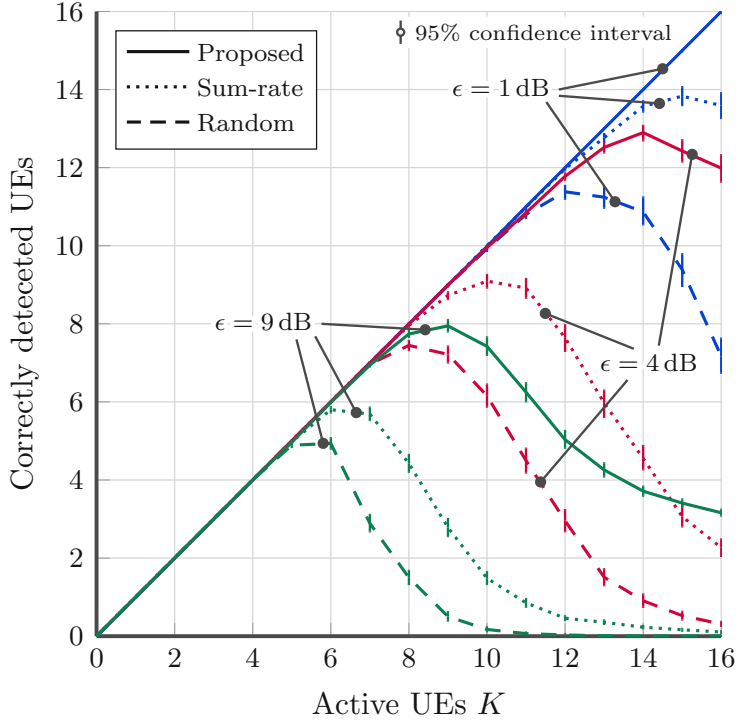


Figure 5.7: Detectability under different K . Here $N = 32$, $P = 30$ dBm, and the pathloss spread is ± 3 dB.

NOMA codebook distributes the interference across all the signatures, which leads to improved performance under IC, as canceling one UE would help in improving the SINR of all remaining UEs, and not just to a limited subset as would occur when reusing the OMA codebook. We see that our approach provides a robust adaptation of the RIS phase-shifts with respect to the pathloss spread, and is able to create the necessary power gaps that result in the required SINR levels. As for the sum-rate optimized phases, we observe that the performance improves as the pathloss spread increases. As explained in Section 5.3.2, the power gap between the UEs resulting from the larger pathloss spreads goes in the favor of the sum-rate solution. At low pathloss spreads, their user-separability performance approaches that of the random shifts. The gain at those ranges is mostly an SNR gain, which is not visible in the figure due to the relatively high transmit power.

To further investigate that, in Figure 5.9 we show the performance over the number of RIS elements N , for low and high transmit powers of $P = -5$ dBm and 30 dBm, respectively. We set the pathloss spread to ± 0 dB and $\epsilon = 3$ dB. First, we make the observation that a certain number of elements is required in order for (5.29) to be solvable; second, at low transmit powers and 0 dB pathloss spread, the sum-rate optimized phase-shifts clearly provide an SNR gain compared to the random phase-shifts, converging towards the performance of that at high transmit power as N increases. Our approach, as can be seen, is capable of providing both SNR and SINR gains under IC.

5.3. Combination with Code-Domain NOMA

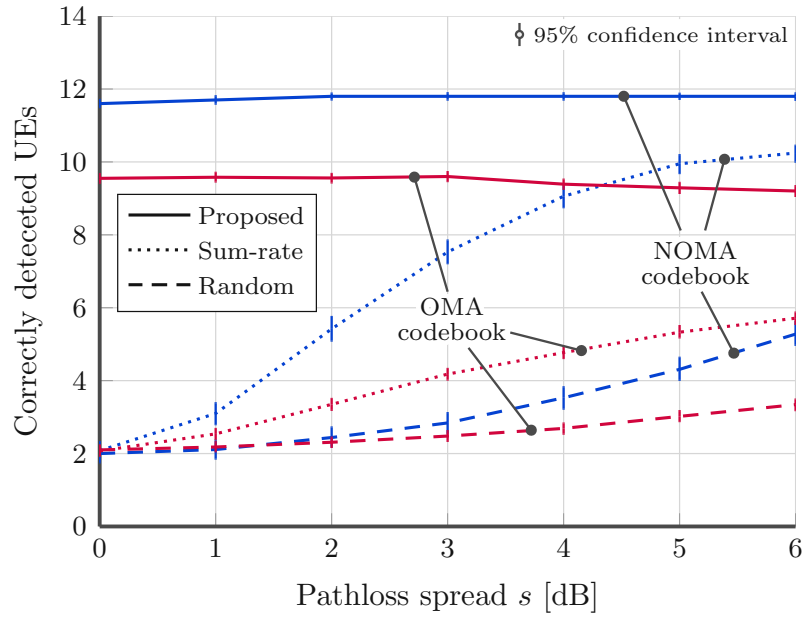


Figure 5.8: Impact of the pathloss spread and codebook design on the performance. Here $N = 32$, $K = 12$, and $\epsilon = 4$ dB.

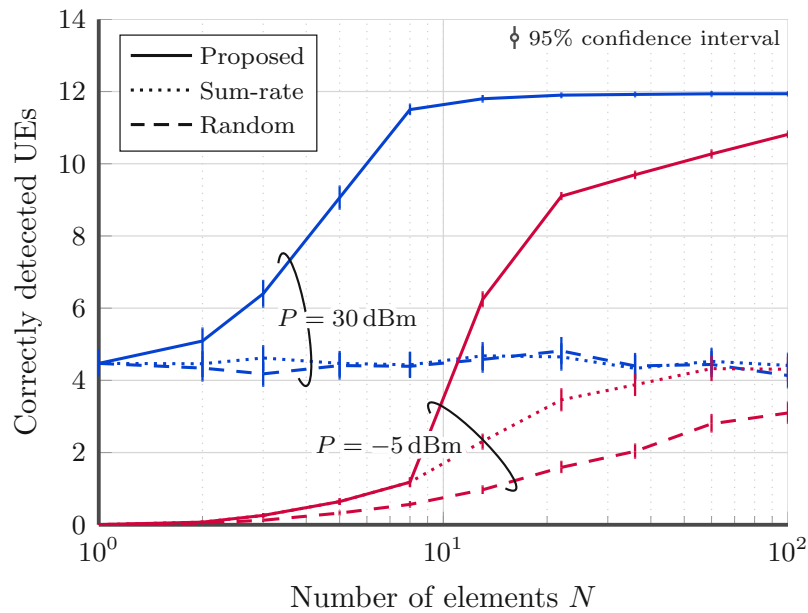


Figure 5.9: Scaling of the performance with N . Here $K = 12$, $\epsilon = 3$ dB, and the pathloss spread is ± 0 dB.

5.4 Final Remarks

This chapter considered the combination of uplink NOMA with RISs. In the first part, we characterized the outage performance of a two-UE RIS-assisted NOMA uplink, where the surface elements are split in boosting either of the UEs. Our results show the strong impact of the RIS optimization on the detection performance under IC, and how it is possible to optimize it in a way that guarantees robust transmission for both UEs. We made here the assumption that the phase-shifts applied at the RIS elements are from a continuous range. In practice, however, it is most likely that the RIS implementation will be based on a discrete set of phase-shifts. Therefore, it makes sense to extend these results taking into account the effects of discrete phase-shifting. Another aspect is the channel estimation for those surfaces; the results given assume perfect knowledge of the channel, which is not available in practice.

In the second part of the chapter, we considered the combination of RISs with code-domain NOMA under a clustered massive MIMO deployment, where each cluster is served by a RIS. We showed how the clusters' RISs can be optimized in order to improve the number of UEs supported. We utilized sum-rate optimized phase-shifts as an initial solution to determine the detection order and the applied filters, and then found the final phase-shifts via a semi-definite relaxation of the problem. Further aspects not considered are as follows. The spatial filtering applied is based on MRC, following the assumption of the BS being equipped with a large antenna array. If strong inter-cluster interference is present, then MMSE-based spatial filtering is a possibility. The model in (5.27) would still hold, except that an MMSE filter is applied instead of \mathbf{a}^H . Another solution, which is based on the code-domain, is to jointly design the spreading codebooks across the different clusters as in Section 3.3, such that the cross-correlation between the signatures of the neighbouring clusters is reduced. Although we assumed the use of successive MMSE-IC detection, several alternatives are possible. For example, parallel IC can be performed, allowing multiple UEs to be detected per IC iteration, which reduces the detection latency for large K . Also, under certain conditions, the NOMA filtering can be performed in a low-complexity manner as discussed in Section 4.3. Finally, for the extraction of the final phase-shifts in (5.40), a rank-1 approximation based on the strongest eigenvector was applied. In general, a better solution might be achieved via a Gaussian randomization procedure [97].

6

Conclusion and Outlook

Addressing the demands of future mobile networks for high data-rates and increased connectivity requires the exploration of new technologies. In the context of supporting massive connectivity, NOMA has been identified as a promising solution, relaxing the current limitations of OMA by allowing multiple UEs to access the same time-frequency resources simultaneously. Motivated by the potential gains of NOMA, we investigated in this dissertation various aspects of the NOMA communication chain, focusing on uplink code-domain transmission.

6.1 Summary of Contribution

The first part of the dissertation dealt with the optimization of the transmitter. We considered the problem of designing the spreading signatures across the different UEs, and investigated a Grassmannian-based approach to the design problem. We proposed an iterative algorithm for constructing Grassmannian codebooks, and showed that the resulting codebooks enjoy close-to-optimal correlation properties, enabling good detection performance even under suboptimal detectors. We then extended the design problem to the case where the UEs are available as groups, such as cells or spatial clusters, and proposed a joint codebook design approach. The results show that the proposed designs can improve the detection performance of the UEs, especially those suffering from high inter-group interference.

In the second part, we switched our attention to the receiver side. We considered the problem of activity detection in the context of grant-free access, and formulated a subspace detector with MUSIC based on a practical frame-structure, and realistic data and pilots' allocation. We found out that the selectivity of the channel can have a substantial impact on the performance. Accordingly, in order to address the influence of strong time-frequency correlation, we proposed to overlay the pilot-blocks with user-specific masking sequences. This allows to fully recover the rank of the signal part of the autocorrelation matrix, which then guarantees a successful operation of the subspace detector. On the other hand, in order to deal with strong time-frequency selectivity, we proposed to readjust the resources' allocation in order

to reduce the experienced channel variability along the sequences. Then, we considered reducing the detection complexity for the data part of the transmission. We showed how the time-frequency correlation of the channel can be utilized to reduce the number of calculated filters across the time-frequency grid, and allow their reuse among neighbouring spreading blocks. We also utilized the availability of multiple receive antennas at the BS to replace exact MMSE filtering with an approximate filter having a lower computational complexity.

The last part of the dissertation focused on the controllability of the channel via RISs. We first characterized the outage performance of a two-UE NOMA uplink, in which part of the RIS elements are configured to boost the signal of one of the UEs and the other part is configured to boost the second one. We proposed to apply a gamma approximation of the receive powers and derived expressions for the outage probability under IC. Our results further illustrate the impact of the RIS optimization on the NOMA performance under IC, and also allow to identify robust operating points that enable sufficient link reliability for both UEs. We finally considered the combination of RISs with a K -UE code-domain NOMA uplink, in the context of cluster-based massive MIMO deployment, where each cluster is served by a RIS. We proposed to optimize the RIS in two steps: first, using an initial solution based on the sum-rate, we determine the detection order and applied filters; then, having those determined, the final phase-shifts are found via a semi-definite relaxation of the problem. Our results show that the considered optimization approach can greatly improve the number of UEs supported by the system.

6.2 Possible Future Work Directions

Throughout this work, we assumed the spreading signatures are constructed once and then get assigned randomly to the UEs. However, in the case of grant-based access, and also in grant-free access with preconfiguration, controlled signature assignment can be beneficial. For example, if the BS knows that on a certain resources' region the number of active UEs is less than the spreading length, then it makes sense to assign them orthogonal signatures, avoiding multi-user interference. Only when the activity increases beyond the spreading length, then the BS switches to the non-orthogonal codebook. This holds true for both the data and pilot spreading. Also, the resources' allocation was assumed to be static, while in practice it might be adjusted over time depending on the load. Therefore, it can be interesting to investigate the influence of resources' allocation and signature assignment under varying user load and activity.

The considered results were mainly based on link-level simulations investigating the performance of certain signal processing stages. It would be interesting to perform system-level abstraction of the considered aspects, which then would allow us to evaluate the performance on a network-wide basis, involving the simulation of hundreds, or even thousands of UEs. This allows to properly gauge the benefit that NOMA brings to an entire network, across various supported service types,

and not just considering the performance under a specific link setup. Alternatively, evaluating the performance analytically of the entire system, including the effects of codebook design, activity detection, channel estimation, and data detection would be an interesting topic as well, albeit difficult to achieve.

In the 5th chapter, the RISs considered were based on a reflective model. That is, the incident waves upon the RISs are reflected with modified phase, amplitude, etc. However, the concept of simultaneous transmitting and reflecting reconfigurable intelligent surfaces, known as STAR-RISs, has also been considered lately [99]. In this case, the incident waves are not just reflected off the surface, but also transmitted through it. With this comes the modification of amplitude and phase not only of the reflected waves, but also of the transmitted waves as well. For example, such a surface can be deployed along windows, and therefore can serve both indoor and outdoor UEs simultaneously. It is then interesting to investigate the considered topics here with this type of surfaces, especially that they naturally provide clusters of indoor and outdoor UEs, which can be utilized by our joint codebook design.

Finally, almost all the considered topics in this dissertation can be investigated from a machine-learning perspective. On the transmitter side, deep-learning can be used for constructing the codebooks (e.g., as in [100]). From the receiver-side perspective, the detection procedure can also be implemented with the aid of deep-learning. For example, it is possible to train neural networks that are capable of performing activity detection [78]. Similarly for the optimization of the RISs, deep learning can also be utilized [101]. Therefore, investigating machine-learning models that generalize well to practical scenarios for our code-domain NOMA setup, and using them to solve problems that are otherwise difficult to address with direct processing, can be an interesting research direction.

List of Abbreviations

1G	1st generation
2G	2nd generation
3G	3rd generation
4G	4th generation
5G	5th generation
AIC	Akaike information criterion
AMP	approximate message passing
AP	alternating projection
B5G	beyond fifth-generation
BIC	Bayesian information criterion
BLER	block error ratio
BS	base station
CBGC	coherence-based Grassmannian codebook
CDF	cumulative distribution function
CDMA	code-division multiple-access
CLT	central limit theorem
CRC	cyclic-redundancy-check
CS	compressed sensing
CSI	channel state information
EM	expectation maximization
EP	expectation propagation
ETF	equiangular tight frame
FDMA	frequency-division multiple access
FFT	fast-Fourier-transform
IC	interference cancellation
ICBP	iterative collision-based packing
IoT	Internet-of-things
IRS	intelligent reflecting surface
KKT	Karush-Kuhn-Tucker
LDPC	low-density parity-check

List of Abbreviations

LOS	line-of-sight
LS	least-squares
LTE	long-term evolution
MA	multiple access
MF	matched filter
MIMO	multiple-input multiple-output
ML	maximum likelihood
MMSE	minimum mean square error
MRC	maximum-ratio combining
MTC	machine-type communication
MUD	multiuser detection
MUSIC	MUltiple SIgnal Classification
NLOS	non-line-of-sight
NOMA	non-orthogonal multiple access
OFDM	orthogonal frequency-division multiplexing
OFDMA	orthogonal frequency-division multiple access
OMA	orthogonal multiple access
OMP	orthogonal matching pursuit
PAPR	peak-to-average-power ratio
PIC	parallel interference cancellation
QAM	quadrature amplitude modulation
RB	resource-block
RE	resource-element
RIS	reconfigurable intelligent surface
RMS	root-mean-square
RV	random variable
SDP	semidefinite programming
SIC	successive interference cancellation
SINR	signal-to-interference-plus-noise ratio
SNR	signal-to-noise ratio
SCMA	sparse-code multiple access
SVD	singular value decomposition
TDL-C	tapped-delay-line-C
TDMA	time-division multiple access
UE	user equipment
WBE	Welch-bound-equality

Notation

In the following table, we describe the notation used throughout this work.

x, X	Non-boldface letters denote scalars
\mathbf{x}	Lowercase boldface letters denote vectors
\mathbf{X}	Uppercase boldface letters denote matrices
\mathcal{X}	Calligraphic letters denote sets
\mathbf{I}_L	Identity matrix of size $L \times L$
$\mathbf{0}_{L \times K}$	All-zeros matrix of size $L \times K$
$\text{diag}(\mathbf{x})$	Diagonal matrix with the elements of \mathbf{x} on its diagonal
$ x $	Magnitude of a scalar
$ \mathbf{x} $	Number of elements in a vector
$\ \mathbf{x}\ $	Euclidean norm of a vector
$\ \mathbf{X}\ _F$	Frobenius norm of a matrix
$\ \mathbf{X}\ _{\max}$	Element-wise maximum norm of a matrix
$ \mathcal{X} $	Size of a set
$[\mathbf{X}]_{i,j}$	Element at the i^{th} row and j^{th} column of a matrix
$(\cdot)^T$	Transpose operation
$(\cdot)^H$	Hermitian operation
$(\cdot)^{-1}$	Inverse operation
$(\cdot)^{1/2}$	Square-root operation
$\text{tr}(\cdot)$	Trace of a matrix
$\det(\cdot)$	Determinant of a matrix
$\rho(\cdot)$	Spectral radius of a matrix
\otimes	Kronecker product
$*$	Column-wise Khatri–Rao product
\circ	Element-wise (Hadamard) product
\oslash	Element-wise division

Notation

$\mathbb{E}\{.\}$	Expected value of a random variable
$\text{Var}\{.\}$	Variance of a random variable
$\mathcal{CN}(\mu, \sigma^2)$	Complex Gaussian distribution with mean μ and variance σ^2
$\text{Nakagami}(m, 1)$	Nakagami distribution with shape m and a spread of 1
$\Gamma(v, \theta)$	Gamma distribution with shape v and scale θ
$I(., ., .)$	Regularized incomplete beta function
$\gamma(., .)$	Regularized incomplete gamma function

References

- [1] Ericsson, “Ericsson Mobility Report,” Nov. 2021. [Online]. Available: <https://www.ericsson.com/4ad7e9/assets/local/reports-papers/mobility-report/documents/2021/ericsson-mobility-report-november-2021.pdf>
- [2] Y. Saito, Y. Kishiyama, A. Benjebbour, T. Nakamura, A. Li, and K. Higuchi, “Non-Orthogonal Multiple Access (NOMA) for Cellular Future Radio Access,” in *2013 IEEE 77th Vehicular Technology Conference (VTC Spring)*, 2013, pp. 1–5.
- [3] K. Higuchi and A. Benjebbour, “Non-orthogonal Multiple Access (NOMA) with Successive Interference Cancellation for Future Radio Access,” *IEICE Transactions on Communications*, vol. E98.B, no. 3, pp. 403–414, 2015.
- [4] Z. Ding, Y. Liu, J. Choi, Q. Sun, M. Elkashlan, I. Chih-Lin, and H. V. Poor, “Application of Non-Orthogonal Multiple Access in LTE and 5G Networks,” *IEEE Communications Magazine*, vol. 55, no. 2, pp. 185–191, 2017.
- [5] K. Higuchi and Y. Kishiyama, “Non-Orthogonal Access with Random Beamforming and Intra-Beam SIC for Cellular MIMO Downlink,” in *2013 IEEE 78th Vehicular Technology Conference (VTC Fall)*, 2013, pp. 1–5.
- [6] B. Kimy, S. Lim, H. Kim, S. Suh, J. Kwun, S. Choi, C. Lee, S. Lee, and D. Hong, “Non-orthogonal Multiple Access in a Downlink Multiuser Beamforming System,” in *MILCOM 2013 - 2013 IEEE Military Communications Conference*, 2013, pp. 1278–1283.
- [7] Z. Ding, F. Adachi, and H. V. Poor, “The Application of MIMO to Non-Orthogonal Multiple Access,” *IEEE Transactions on Wireless Communications*, vol. 15, no. 1, pp. 537–552, 2016.
- [8] Y. Liu, Z. Qin, M. Elkashlan, Z. Ding, A. Nallanathan, and L. Hanzo, “Nonorthogonal Multiple Access for 5G and Beyond,” *Proceedings of the IEEE*, vol. 105, no. 12, pp. 2347–2381, 2017.
- [9] L. Dai, B. Wang, Y. Yuan, S. Han, I. Chih-lin, and Z. Wang, “Non-orthogonal multiple access for 5G: solutions, challenges, opportunities, and future research trends,” *IEEE Communications Magazine*, vol. 53, no. 9, pp. 74–81, 2015.

- [10] M. Vaezi, R. Schober, Z. Ding, and H. V. Poor, “Non-Orthogonal Multiple Access: Common Myths and Critical Questions,” *IEEE Wireless Communications*, vol. 26, no. 5, pp. 174–180, 2019.
- [11] N. Abu-Ali, A. M. Taha, M. Salah, and H. Hassanein, “Uplink Scheduling in LTE and LTE-Advanced: Tutorial, Survey and Evaluation Framework,” *IEEE Communications Surveys Tutorials*, vol. 16, no. 3, pp. 1239–1265, Third 2014.
- [12] S. Ali, N. Rajatheva, and W. Saad, “Fast Uplink Grant for Machine Type Communications: Challenges and Opportunities,” *IEEE Communications Magazine*, vol. 57, no. 3, pp. 97–103, 2019.
- [13] L. Tian, C. Yan, W. Li, Z. Yuan, W. Cao, and Y. Yuan, “On uplink non-orthogonal multiple access for 5g: opportunities and challenges,” *China Communications*, vol. 14, no. 12, pp. 142–152, December 2017.
- [14] L. Liu, E. G. Larsson, W. Yu, P. Popovski, C. Stefanovic, and E. de Carvalho, “Sparse Signal Processing for Grant-Free Massive Connectivity: A Future Paradigm for Random Access Protocols in the Internet of Things,” *IEEE Signal Processing Magazine*, vol. 35, no. 5, pp. 88–99, 2018.
- [15] G. Wunder, P. Jung, and C. Wang, “Compressive random access for post-LTE systems,” in *2014 IEEE International Conference on Communications Workshops (ICC)*, 2014, pp. 539–544.
- [16] G. Hannak, M. Mayer, A. Jung, G. Matz, and N. Goertz, “Joint channel estimation and activity detection for multiuser communication systems,” in *2015 IEEE International Conference on Communication Workshop (ICCW)*, 2015, pp. 2086–2091.
- [17] B. Wang, L. Dai, Y. Zhang, T. Mir, and J. Li, “Dynamic Compressive Sensing-Based Multi-User Detection for Uplink Grant-Free NOMA,” *IEEE Communications Letters*, vol. 20, no. 11, pp. 2320–2323, 2016.
- [18] M. B. Shahab, R. Abbas, M. Shirvanimoghaddam, and S. J. Johnson, “Grant-Free Non-Orthogonal Multiple Access for IoT: A Survey,” *IEEE Communications Surveys Tutorials*, vol. 22, no. 3, pp. 1805–1838, 2020.
- [19] L. Dai, B. Wang, Z. Ding, Z. Wang, S. Chen, and L. Hanzo, “A Survey of Non-Orthogonal Multiple Access for 5G,” *IEEE Communications Surveys Tutorials*, vol. 20, no. 3, pp. 2294–2323, 2018.
- [20] M. D. Renzo, M. Debbah, D.-T. Phan-Huy, A. Zappone, M.-S. Alouini, C. Yuen, V. Sciancalepore, G. C. Alexandropoulos, J. Hoydis, H. Gacanin, J. d. Rosny, A. Bounceur, G. Lerosey, and M. Fink, “Smart radio environments empowered by reconfigurable AI meta-surfaces: an idea whose time has come,” *EURASIP Journal on Wireless Communications and Networking*, vol. 2019, no. 1, p. 129, May 2019.

- [21] Q. Wu and R. Zhang, "Intelligent Reflecting Surface Enhanced Wireless Network via Joint Active and Passive Beamforming," *IEEE Transactions on Wireless Communications*, vol. 18, no. 11, pp. 5394–5409, 2019.
- [22] E. Basar, M. Di Renzo, J. De Rosny, M. Debbah, M.-S. Alouini, and R. Zhang, "Wireless Communications Through Reconfigurable Intelligent Surfaces," *IEEE Access*, vol. 7, pp. 116 753–116 773, 2019.
- [23] Q. Wu and R. Zhang, "Towards Smart and Reconfigurable Environment: Intelligent Reflecting Surface Aided Wireless Network," *IEEE Communications Magazine*, vol. 58, no. 1, pp. 106–112, 2020.
- [24] M. Di Renzo, A. Zappone, M. Debbah, M.-S. Alouini, C. Yuen, J. de Rosny, and S. Tretyakov, "Smart Radio Environments Empowered by Reconfigurable Intelligent Surfaces: How It Works, State of Research, and The Road Ahead," *IEEE Journal on Selected Areas in Communications*, vol. 38, no. 11, pp. 2450–2525, 2020.
- [25] Z. Ding and H. Vincent Poor, "A Simple Design of IRS-NOMA Transmission," *IEEE Communications Letters*, vol. 24, no. 5, pp. 1119–1123, 2020.
- [26] M. Fu, Y. Zhou, and Y. Shi, "Intelligent Reflecting Surface for Downlink Non-Orthogonal Multiple Access Networks," in *2019 IEEE Globecom Workshops (GC Wkshps)*, 2019, pp. 1–6.
- [27] G. Yang, X. Xu, and Y. Liang, "Intelligent Reflecting Surface Assisted Non-Orthogonal Multiple Access," in *2020 IEEE Wireless Communications and Networking Conference (WCNC)*, 2020, pp. 1–6.
- [28] Y. Cheng, K. H. Li, Y. Liu, K. C. Teh, and H. Vincent Poor, "Downlink and Uplink Intelligent Reflecting Surface Aided Networks: NOMA and OMA," *IEEE Transactions on Wireless Communications*, vol. 20, no. 6, pp. 3988–4000, 2021.
- [29] Z. Ding, R. Schober, and H. V. Poor, "On the Impact of Phase Shifting Designs on IRS-NOMA," *IEEE Wireless Communications Letters*, vol. 9, no. 10, pp. 1596–1600, 2020.
- [30] A. S. d. Sena, D. Carrillo, F. Fang, P. H. J. Nardelli, D. B. d. Costa, U. S. Dias, Z. Ding, C. B. Papadias, and W. Saad, "What Role Do Intelligent Reflecting Surfaces Play in Multi-Antenna Non-Orthogonal Multiple Access?" *IEEE Wireless Communications*, vol. 27, no. 5, pp. 24–31, 2020.
- [31] T. Hou, Y. Liu, Z. Song, X. Sun, Y. Chen, and L. Hanzo, "Reconfigurable Intelligent Surface Aided NOMA Networks," *IEEE Journal on Selected Areas in Communications*, vol. 38, no. 11, pp. 2575–2588, 2020.

- [32] 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,” *EURASIP Journal on Wireless Communications and Networking*, vol. 2018, no. 1, p. 226, Sep. 2018.
- [33] B. Tahir, S. Schwarz, and M. Rupp, “Constructing Grassmannian Frames by an Iterative Collision-Based Packing,” *IEEE Signal Processing Letters*, vol. 26, no. 7, pp. 1056–1060, 2019.
- [34] —, “Joint Codebook Design for Multi-Cell NOMA,” in *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 4814–4818.
- [35] —, “Low-Complexity Detection of Uplink NOMA by Exploiting Properties of the Propagation Channel,” in *ICC 2020 - 2020 IEEE International Conference on Communications (ICC)*, 2020, pp. 1–6.
- [36] —, “Collision Resilient V2X Communication via Grant-Free NOMA,” in *2020 28th European Signal Processing Conference (EUSIPCO)*, 2021, pp. 1732–1736.
- [37] —, “Impact of Channel Correlation on Subspace-Based Activity Detection in Grant-Free NOMA,” *submitted to the 2022 IEEE 95th Vehicular Technology Conference (VTC2022-Spring)*, 2022. [Online]. Available: <https://arxiv.org/abs/2202.11161>
- [38] —, “Analysis of Uplink IRS-Assisted NOMA Under Nakagami-m Fading via Moments Matching,” *IEEE Wireless Communications Letters*, vol. 10, no. 3, pp. 624–628, 2021.
- [39] —, “Outage Analysis of Uplink IRS-Assisted NOMA under Elements Splitting,” in *2021 IEEE 93rd Vehicular Technology Conference (VTC2021-Spring)*, 2021, pp. 1–5.
- [40] —, “RIS-Assisted Code-Domain MIMO-NOMA,” in *2021 29th European Signal Processing Conference (EUSIPCO)*, 2021, pp. 821–825.
- [41] H. Nikopour and H. Baligh, “Sparse code multiple access,” in *2013 IEEE 24th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC)*, 2013, pp. 332–336.
- [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,” *EURASIP Journal on Wireless Communications and Networking*, vol. 2018, no. 1, p. 17, Sep. 2018.

- [43] J. H. Conway, R. H. Hardin, and N. J. A. Sloane, "Packing Lines, Planes, etc.: Packings in Grassmannian Space," *ArXiv Mathematics e-prints*, Jul. 2002.
- [44] T. Strohmer and R. W. Heath, "Grassmannian frames with applications to coding and communication," *Applied and Computational Harmonic Analysis*, vol. 14, no. 3, pp. 257 – 275, 2003.
- [45] L. Welch, "Lower bounds on the maximum cross correlation of signals (Corresp.)," *IEEE Transactions on Information Theory*, vol. 20, no. 3, pp. 397–399, May 1974.
- [46] P. Viswanath, V. Anantharam, and D. Tse, "Optimal sequences, power control, and user capacity of synchronous CDMA systems with linear MMSE multiuser receivers," *IEEE Transactions on Information Theory*, vol. 45, no. 6, pp. 1968–1983, 1999.
- [47] J. Zhang and E. Chong, "CDMA systems in fading channels: admissibility, network capacity, and power control," *IEEE Transactions on Information Theory*, vol. 46, no. 3, pp. 962–981, 2000.
- [48] R. W. Heath, T. Strohmer, and A. J. Paulraj, "Grassmannian signatures for CDMA systems," in *GLOBECOM '03. IEEE Global Telecommunications Conference (IEEE Cat. No.03CH37489)*, vol. 3, Dec 2003, pp. 1553–1557 vol.3.
- [49] M. Fickus, D. G. Mixon, and J. C. Tremain, "Steiner equiangular tight frames," *ArXiv e-prints*, Sep. 2010.
- [50] B. G. Bodmann and J. Haas, "Achieving the orthoplex bound and constructing weighted complex projective 2-designs with Singer sets," *Linear Algebra and its Applications*, vol. 511, pp. 54 – 71, 2016.
- [51] P. Xia, S. Zhou, and G. B. Giannakis, "Achieving the Welch bound with difference sets," *IEEE Transactions on Information Theory*, vol. 51, no. 5, pp. 1900–1907, May 2005.
- [52] B. M. Hochwald, T. L. Marzetta, T. J. Richardson, W. Sweldens, and R. Urbanke, "Systematic design of unitary space-time constellations," *IEEE Transactions on Information Theory*, vol. 46, no. 6, pp. 1962–1973, Sep 2000.
- [53] J. A. Tropp, I. S. Dhillon, R. W. Heath, and T. Strohmer, "Designing structured tight frames via an alternating projection method," *IEEE Transactions on Information Theory*, vol. 51, no. 1, pp. 188–209, Jan 2005.
- [54] I. S. Dhillon, J. R. W. Heath, T. Strohmer, and J. A. Tropp, "Constructing Packings in Grassmannian Manifolds via Alternating Projection," *Experimental Mathematics*, vol. 17, no. 1, pp. 9–35, 2008. [Online]. Available: <https://doi.org/10.1080/10586458.2008.10129018>

- [55] H. Zörlein and M. Bossert, “Coherence Optimization and Best Complex Antipodal Spherical Codes,” *IEEE Transactions on Signal Processing*, vol. 63, no. 24, pp. 6606–6615, Dec 2015.
- [56] H. E. A. Laue and W. P. du Plessis, “A Coherence-Based Algorithm for Optimizing Rank-1 Grassmannian Codebooks,” *IEEE Signal Processing Letters*, vol. 24, no. 6, pp. 823–827, June 2017.
- [57] R. A. Rankin, “The Closest Packing of Spherical Caps in n Dimensions,” *Proceedings of the Glasgow Mathematical Association*, vol. 2, no. 3, p. 139–144, 1955.
- [58] G. A. Kabatiansky and V. I. Levenshtein, “Bounds for packings on the sphere and in space,” *Problemy Peredači Informacii*, vol. 14, no. 1, pp. 3–25, 1978.
- [59] K. K. Mukkavilli, A. Sabharwal, E. Erkip, and B. Aazhang, “On beamforming with finite rate feedback in multiple-antenna systems,” *IEEE Transactions on Information Theory*, vol. 49, no. 10, pp. 2562–2579, Oct 2003.
- [60] D. J. Love, R. W. Heath, and T. Strohmer, “Grassmannian beamforming for multiple-input multiple-output wireless systems,” *IEEE Transactions on Information Theory*, vol. 49, no. 10, pp. 2735–2747, Oct 2003.
- [61] M. Fickus and D. G. Mixon, “Tables of the existence of equiangular tight frames,” *ArXiv e-prints*, Apr. 2015.
- [62] J. von Neumann, *Functional Operators (AM-22), Volume 2: The Geometry of Orthogonal Spaces. (AM-22)*. Princeton University Press, 1950.
- [63] W. Cheney and A. A. Goldstein, “Proximity Maps for Convex Sets,” *Proceedings of the American Mathematical Society*, vol. 10, no. 3, pp. 448–450, 1959.
- [64] I. S. Dhillon, R. W. Heath, Jr, T. Strohmer, and J. A. Tropp, “Constructing packings in Grassmannian manifolds via alternating projection,” *ArXiv e-prints*, Sep. 2007.
- [65] R. T. Rockafellar, *Convex analysis*, ser. Princeton Mathematical Series. Princeton, N. J.: Princeton University Press, 1970.
- [66] S. Kim, S. Kim, J. Kim, K. Lee, S. Choi, and B. Shim, “Low Latency Random Access for Small Cell Toward Future Cellular Networks,” *IEEE Access*, vol. 7, pp. 178 563–178 576, 2019.
- [67] M. Mohammadkarimi, M. A. Raza, and O. A. Dobre, “Signature-Based Nonorthogonal Massive Multiple Access for Future Wireless Networks: Uplink Massive Connectivity for Machine-Type Communications,” *IEEE Vehicular Technology Magazine*, vol. 13, no. 4, pp. 40–50, 2018.

- [68] L. Marijanović, S. Schwarz, and M. Rupp, “Multiplexing Services in 5G and Beyond: Optimal Resource Allocation Based on Mixed Numerology and Mini-Slots,” *IEEE Access*, vol. 8, pp. 209 537–209 555, 2020.
- [69] J. Kim, G. Lee, S. Kim, T. Taleb, S. Choi, and S. Bahk, “Two-Step Random Access for 5G System: Latest Trends and Challenges,” *IEEE Network*, vol. 35, no. 1, pp. 273–279, 2021.
- [70] B. Wang, L. Dai, T. Mir, and Z. Wang, “Joint User Activity and Data Detection Based on Structured Compressive Sensing for NOMA,” *IEEE Communications Letters*, vol. 20, no. 7, pp. 1473–1476, 2016.
- [71] Y. Du, B. Dong, W. Zhu, P. Gao, Z. Chen, X. Wang, and J. Fang, “Joint Channel Estimation and Multiuser Detection for Uplink Grant-Free NOMA,” *IEEE Wireless Communications Letters*, vol. 7, no. 4, pp. 682–685, 2018.
- [72] Y. Du, C. Cheng, B. Dong, Z. Chen, X. Wang, J. Fang, and S. Li, “Block-Sparsity-Based Multiuser Detection for Uplink Grant-Free NOMA,” *IEEE Transactions on Wireless Communications*, vol. 17, no. 12, pp. 7894–7909, 2018.
- [73] J. Zhang, Y. Pan, and J. Xu, “Compressive Sensing for Joint User Activity and Data Detection in Grant-Free NOMA,” *IEEE Wireless Communications Letters*, vol. 8, no. 3, pp. 857–860, 2019.
- [74] S. M. Hasan, K. Mahata, and M. M. Hyder, “Uplink Grant-Free NOMA With Sinusoidal Spreading Sequences,” *IEEE Transactions on Communications*, vol. 69, no. 6, pp. 3757–3770, 2021.
- [75] L. Cheng, L. Liu, and S. Cui, “A Covariance-based User Activity Detection and Channel Estimation Approach with Novel Pilot Design,” in *2020 IEEE 21st International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, 2020, pp. 1–5.
- [76] C. Wei, H. Liu, Z. Zhang, J. Dang, and L. Wu, “Approximate Message Passing-Based Joint User Activity and Data Detection for NOMA,” *IEEE Communications Letters*, vol. 21, no. 3, pp. 640–643, 2017.
- [77] J. Ahn, B. Shim, and K. B. Lee, “EP-Based Joint Active User Detection and Channel Estimation for Massive Machine-Type Communications,” *IEEE Transactions on Communications*, vol. 67, no. 7, pp. 5178–5189, 2019.
- [78] W. Kim, Y. Ahn, and B. Shim, “Deep Neural Network-Based Active User Detection for Grant-Free NOMA Systems,” *IEEE Transactions on Communications*, vol. 68, no. 4, pp. 2143–2155, 2020.

- [79] G. Schwarz, “Estimating the Dimension of a Model,” *The Annals of Statistics*, vol. 6, no. 2, pp. 461 – 464, 1978. [Online]. Available: <https://doi.org/10.1214/aos/1176344136>
- [80] P. Stoica and Y. Selen, “Model-order selection: a review of information criterion rules,” *IEEE Signal Processing Magazine*, vol. 21, no. 4, pp. 36–47, 2004.
- [81] H. Akaike, “A new look at the statistical model identification,” *IEEE Transactions on Automatic Control*, vol. 19, no. 6, pp. 716–723, 1974.
- [82] J. J. Dziak, D. L. Coffman, S. T. Lanza, R. Li, and L. S. Jermin, “Sensitivity and specificity of information criteria,” *Briefings in Bioinformatics*, vol. 21, no. 2, pp. 553–565, 03 2019. [Online]. Available: <https://doi.org/10.1093/bib/bbz016>
- [83] V. T. Ermolaev, A. A. Mal’tsev, and K. V. Rodyushkin, “Statistical Characteristics of the AIC and MDL Criteria in the Problem of Estimating the Number of Sources of Multivariate Signals in the Case of a Short Sample,” *Radiophysics and Quantum Electronics*, vol. 44, no. 12, pp. 977–983, Dec 2001. [Online]. Available: <https://doi.org/10.1023/A:1014882129741>
- [84] N. Y. Yu, “Binary Golay Spreading Sequences and Reed-Muller Codes for Uplink Grant-Free NOMA,” *IEEE Transactions on Communications*, vol. 69, no. 1, pp. 276–290, 2021.
- [85] 3rd Generation Partnership Project (3GPP), “Technical Specification Group Radio Access Network; Study on channel model for frequency spectrum above 6 GHz,” 3rd Generation Partnership Project (3GPP), TR 38.900, Jun. 2018.
- [86] E. Zöchmann, S. Schwarz, S. Pratschner, L. Nagel, M. Lerch, and M. Rupp, “Exploring the physical layer frontiers of cellular uplink,” *EURASIP Journal on Wireless Communications and Networking*, vol. 2016, no. 1, pp. 1–18, 2016. [Online]. Available: <http://dx.doi.org/10.1186/s13638-016-0609-1>
- [87] F. Rusek, D. Persson, B. K. Lau, E. G. Larsson, T. L. Marzetta, O. Edfors, and F. Tufvesson, “Scaling Up MIMO: Opportunities and Challenges with Very Large Arrays,” *IEEE Signal Processing Magazine*, vol. 30, no. 1, pp. 40–60, Jan 2013.
- [88] D. Zhu, B. Li, and P. Liang, “On the matrix inversion approximation based on neumann series in massive MIMO systems,” in *2015 IEEE International Conference on Communications (ICC)*, June 2015, pp. 1763–1769.
- [89] J. Minango and C. de Almeida, “Low-complexity MMSE detector based on the first-order Neumann series expansion for massive MIMO systems,” in *2017 IEEE 9th Latin-American Conference on Communications (LATINCOM)*, Nov 2017, pp. 1–5.

References

- [90] M. Zeng, X. Li, G. Li, W. Hao, and O. A. Dobre, “Sum Rate Maximization for IRS-Assisted Uplink NOMA,” *IEEE Communications Letters*, vol. 25, no. 1, pp. 234–238, 2021.
- [91] I. Atzeni, J. Arnau, and M. Kountouris, “Downlink Cellular Network Analysis With LOS/NLOS Propagation and Elevated Base Stations,” *IEEE Transactions on Wireless Communications*, vol. 17, no. 1, pp. 142–156, 2018.
- [92] J. Lyu and R. Zhang, “Spatial Throughput Characterization for Intelligent Reflecting Surface Aided Multiuser System,” *IEEE Wireless Communications Letters*, vol. 9, no. 6, pp. 834–838, 2020.
- [93] I. Florescu, *Probability and Stochastic Processes*. Wiley, 2014.
- [94] A. Papoulis, *The Fourier integral and its applications*. New York: McGraw-Hill, 1962.
- [95] S. Atapattu, R. Fan, P. Dharmawansa, G. Wang, J. Evans, and T. A. Tsiftsis, “Reconfigurable Intelligent Surface Assisted Two-Way Communications: Performance Analysis and Optimization,” *IEEE Transactions on Communications*, vol. 68, no. 10, pp. 6552–6567, 2020.
- [96] D. Tse and P. Viswanath, *Fundamentals of Wireless Communication*. Cambridge University Press, 2005.
- [97] Z.-q. Luo, W.-k. Ma, A. M.-c. So, Y. Ye, and S. Zhang, “Semidefinite Relaxation of Quadratic Optimization Problems,” *IEEE Signal Processing Magazine*, vol. 27, no. 3, pp. 20–34, 2010.
- [98] M. Grant and S. Boyd, “CVX: Matlab Software for Disciplined Convex Programming, version 2.1,” <http://cvxr.com/cvx>, Mar. 2014.
- [99] J. Xu, Y. Liu, X. Mu, and O. A. Dobre, “STAR-RISs: Simultaneous Transmitting and Reflecting Reconfigurable Intelligent Surfaces,” *IEEE Communications Letters*, vol. 25, no. 9, pp. 3134–3138, 2021.
- [100] M. Han, H. Seo, A. T. Abebe, and C. G. Kang, “Deep Learning-Based Multi-User Multi-Dimensional Constellation Design in Code Domain Non-Orthogonal Multiple Access,” in *2020 IEEE International Conference on Communications Workshops (ICC Workshops)*, 2020, pp. 1–6.
- [101] C. Huang, G. C. Alexandropoulos, C. Yuen, and M. Debbah, “Indoor Signal Focusing with Deep Learning Designed Reconfigurable Intelligent Surfaces,” in *2019 IEEE 20th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, 2019, pp. 1–5.