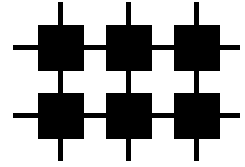




Lehrstuhl für Netzwerktheorie
und Signalverarbeitung



Forschungsberichte

Herausgeber: Prof. Dr. techn. Dr. h. c. Josef A. Nossek

Hela Jedda

**Quantized Constant Envelope Transmit
Signal processing**

Quantized Constant Envelope Transmit Signal Processing

Hela Jedda

Vollständiger Abdruck der von der Fakultät für Elektrotechnik und Informationstechnik der Technischen Universität München zur Erlangung des akademischen Grades eines
Doktor-Ingenieurs
genehmigten Dissertation.

Vorsitzender: Prof. Dr. Holger Boche

Prüfer der Dissertation:

1. Prof. Dr. techn. Dr. h. c. Josef A. Nossek
2. Prof. A. Lee Swindlehurst, Ph.D.,
University of California, Irvine/ USA

Die Dissertation wurde am 18.09.2018 bei der Technischen Universität München eingereicht und durch die Fakultät für Elektrotechnik und Informationstechnik am 22.11.2018 angenommen.

Acknowledgments

I would like to express my sincere gratitude to Professor Josef A. Nossek for giving me the opportunity to pursue a doctoral degree at the Institute of Circuit Theory and Signal Processing and supervising me with his steady support and valuable discussions.

I thank Professor A. Lee Swindlehurst for being my mentor and my second supervisor during the Hans-Fischer fellowship at the Institute for Advanced Study. I additionally would like to thank him for giving me the opportunity to spend a wonderful research stay at the University of California in Irvine during the fall term of 2016. I also thank Professor Holger Boche for serving as a chairman of the examination committee.

I thank Professor Wolfgang Utschick for welcoming me in his research group since April 2016 and providing an enjoyable atmosphere.

Many thanks to all my present and former colleagues for the pleasant and friendly work atmosphere. I would like to mention: Jawad Munir, Israa Slim, Leonardo Baltar, Dr. Michael Joham, Hans Brunner, Manuel Stein, Kilian Roth, Qing Bai, and Bernhard Lehmeier. A special thanks goes to Amine Mezghani for his continuous help and the long constructive discussions we had together.

I also would like to thank my students for all the questions and fruitful discussions. Let me mention Oliver De Candido, Daniel Plabst, Benedikt Fesl, Ferhad Askerbeyli, Andreas Noll, Donia Ben Amor, and Ovais Bin Usman.

A very special gratitude goes to my parents and my family for their everlasting love and all the positive energy they offered me.

Finally, I thank my husband for his love and encouragement. He has been my model in diligence and discipline.

Munich, December 2018

Hela Jedda

Abstract

In this doctoral thesis, we investigate the effect of Quantized Constant Envelope (QCE) signaling at the transmitter on the performance of communication systems. Although massive Multiple-Input Multiple-Output (MIMO) systems, which are characterized by the large number of Base Station (BS) antennas, offer a significant increase in power and spectral efficiency, the hardware power consumption remains a crucial concern especially at Millimeter-Wave (mmW) frequencies. The QCE signaling at the transmitter is a promising approach to enhance the power efficiency. Therefore, we opt for this solution and develop digital signal processing techniques to mitigate the resulting performance loss. Our contribution is two-fold. First, we derive the signal statistical properties of the Constant Envelope Quantizer (CEQ). To this end, we extend Price's theorem to the complex-valued case. Second, we develop new precoding techniques, linear and non-linear, to counteract additionally the CEQ distortions. Finally, we discuss the benefits of QCE systems in terms of power efficiency and their challenges in terms of spectral shaping.

Contents

1. Introduction	1
1.1 Related Works	2
1.2 Outline and Main Contributions	3
1.3 Notation	5
2. Motivation	9
2.1 Power Amplifiers	9
2.2 DACs	9
2.3 Polar Transmitters	10
3. On Constant Envelope Quantizers	11
3.1 One-Dimensional CEQ	11
3.2 Multi-Dimensional CEQ	12
3.3 Statistical Properties	12
3.4 Optimal CEQ	23
3.5 Statistical Properties of the Optimal CEQ	25
3.6 Bussgang Linearization of the CEQ	29
4. System Model	31
4.1 Input-Output Relationship	31
4.2 Compact Input-Output Relationship	32
4.3 Input Set	33
4.4 Transmit Processing: Precoding	33
4.5 Transmit Power Constraint	33
4.6 Channel Model	35
4.7 Receive Processing	36
4.8 Potential Dual Uplink System Model with the Optimal CEQ	37
I Linear Transmit Signal Processing	39
5. Flat-Fading Channels	41
5.1 Input-Output Relationship	41
5.2 Optimization Problem	41
5.3 Precoder Designs in the Primal Domain	42

5.4	Dual Optimization Problem	46
5.5	Simulation Results	51
6.	Frequency-Selective Channels	55
6.1	Input-Output Relationship	55
6.2	Optimization Problem	56
6.3	Precoder Designs in the Primal Domain	58
6.4	Dual Optimization Problem	62
6.5	Simulation Results	66
II	Non-linear Transmit Signal Processing	71
7.	Flat-Fading Channels	73
7.1	Transmit Power Constraint	73
7.2	Optimization Problem	73
7.3	Problem Formulation for PSK Signaling	76
7.4	Problem Formulation for QAM Signaling	79
7.5	Computational Complexity of MSM	87
7.6	Simulation Results	89
8.	Frequency-Selective Channels	95
8.1	Input-Output Relationship	95
8.2	Optimization Problem	96
8.3	Computational Complexity	99
8.4	Simulation Results	100
9.	Discussion: Benefits and Challenges	103
9.1	Power Efficiency	103
9.2	Strong Non-linearities vs. Spectral Shaping	104
10.	Conclusion	119
	Appendix	121
A1	Complex-Valued Gaussian Joint PDF	121
A2	Proof of Theorem 1	122
A3	Second Order Derivative in the Polar Representation	123
A4	Special Integrals	124
A5	Derivative of arcsin(A)	125
	Bibliography	133

1. Introduction

The remarkable increase in the number of personal mobile devices with internet access combined with the excessive emerging of bandwidth-hungry mobile applications such as multimedia communications, e.g. video streaming with High Definition Television (HDTV) and Ultra-High Definition Video (UHDV), and cloud computing, leads to an ever-increasing demand for higher data rates. However, the existing commercial standards cannot meet the ever increasing need for high-speed wireless connectivity anymore. Therefore, the next generation of mobile communication aims at increasing 1000-fold the network capacity, 10-100-fold the number of connected devices and decreasing 5-fold the latency time and the power consumption compared to 4G networks [1]. To achieve these challenging requirements, the following technologies are the subject of current research:

- massive Multiple-Input Multiple-Output (MIMO) systems, where the Base Stations (BSs) are equipped with a very large number of antennas (100 or more) that can simultaneously serve many users [2–6],
- Millimeter-Wave (mmW) communication, i.e. frequencies ranging between 30 GHz and 300 GHz, where the spectrum is less crowded and greater bandwidth is available [7–9] and
- smaller cells with ranges on the order of 10-200 m, i.e. pico- and femtocells.

First, massive MIMO systems lead to a drastic increase in the number of Radio Frequency (RF) chains at the BS and hence in the number of the wireless front-end hardware components. Second, mmW communication implies that the wireless front-end hardware components are operated at much higher frequencies. Third, reducing the cell size means that the number of cells per unit area is increased and thus results in a much more dense wireless network. Combining the three technologies means a dramatic increase in the number of RF hardware elements operating at very high frequencies per unit area. Hence, the RF power consumption per unit area alarmingly increases. While the above technologies are foreseen as key technologies for future communication systems, the increase in power consumption represents a crucial concern.

The most critical front-end elements in terms of power consumption, depending on whether the large number of antennas is situated at the receiver or at the transmitter, are the Analog-to-Digital Converters (ADCs) in the uplink scenario, and mainly the Power Amplifiers (PAs) and secondarily the Digital-to-Analog Converters (DACs) in the downlink scenario, which is the focus of this contribution. According to [10, 11], the PA is considered as the most power hungry device at the transmitter side. When the PA is run in the saturation region, i.e. the highly non-linear region, high power efficiency is achieved and hence less power is consumed [12]. However, the saturation region implies strong non-linear signal

distortions. To omit the signal distortions, while keeping the PA operate in the saturation region, the input signals should fulfill the Constant Envelope (CE) property, which leads to a unit Peak-to-Average-Power Ratio (PAPR).

To this end, polar (phase-based) DACs at the transmitter are designed to convert the discrete-time and discrete-value base-band signals into continuous-time but discrete-value (at the sample rate), i.e. discrete-phase, CE signals. The number of possible discrete phases is determined by the resolution of the DAC. Thus, the polar DAC can be considered as a Constant Envelope Quantizer (CEQ). The larger the resolution is, the more accurate the phase information at the DAC's output is, but the larger its power consumption is [13]. To further reduce the hardware power consumption, the DAC's resolution can be reduced. The use of coarsely quantized DACs is also beneficial in terms of reduced cost and circuit area and can further simplify the surrounding RF circuitry due to the relaxed linearity constraint, leading to very efficient hardware implementations. In this way, the power consumption is reduced twofold: power efficient PAs due to the CE signals and less power consuming polar DACs due to the low resolution. However, this approach can lead to non-linear distortions that degrade the system performance and have to be mitigated by the precoder design in massive Multi-User (MU) MIMO downlink systems.

1.1 Related Works

Many works have addressed the precoding problem in the context of CE transmit signals for massive MIMO systems [14–18], where the Multi-User Interference (MUI) is minimized subject to the CE constraint. Another work [19] opts for minimizing an upper bound of the Symbol Error Ratio (SER) in the case of single-user Multiple-Input Single-Output (MISO) systems for two strategies: antenna-subset selection, where a subset of the antennas is selected for transmission, and unequal power allocation among the antennas, where the magnitude of the transmit signal at each antenna is kept constant over a transmission period but the signal magnitudes at distinct transmit antennas are not necessarily equal. The authors of [20] jointly optimize the transmit CE precoding and the constellation in order to minimize the SER in a MISO multicast system. Recent works in [21] and [22] exploit the constructive part of the MUI to design the CE precoder. The authors in [23] design a CE precoder to maximize the Signal-to-Leakage-plus-Noise Ratio (SLNR). In [24], a CE precoder is jointly designed with the receive beamforming to minimize the SER for point-to-point MIMO systems, while adopting antenna grouping for multi-stream transmission. In the above contributions, the DACs are assumed to have infinite resolution.

The contribution in [25] is an early work that addressed the precoding task with low resolution DACs at the transmitter. A linear Minimum Mean Squared Error (MMSE) precoder is designed, while quantization distortion is taken into account. This precoding design is not given in the context of coarsely Quantized Constant Envelope (QCE) signals since the DACs are not polar but cartesian (in-phase- and quadrature). However, the extreme case of 1-bit DACs in [25] represents a special case of coarsely QCE signals. Many contributions in the literature have studied this special case. They can be categorized in two groups: linear and non-linear precoders. In addition to the linear precoder in [25], we introduced in [26] another linear precoder, where the second-order statistics of the 1-bit DAC signals are computed based on Price's theorem [27]. Non-linear precoding techniques in this context were introduced in [28–34]. The non-linear methods can be classified with respect to two

design criteria: the symbol-wise Mean Squared Error (MSE) and the symbol-wise SER. In the context of the symbol-wise MSE, the authors in [29] presented a convex formulation of the problem and applied it to higher-order modulations [30]. The problem formulation is based on semidefinite relaxation and squared ℓ_∞ -norm relaxation. The same optimization problem was solved more efficiently in [32] and [35].

In the context of the symbol-wise SER, we presented in [28] a precoding technique based on a minimum Bit Error Ratio (BER) criterion and made use of the box norm (ℓ_∞) to relax the 1-bit constraint. Recently, the work in [33] proposed a method to significantly improve linear precoding solutions in conjunction with 1-bit quantization by properly perturbing the linearly precoded signal for each given input signal to favorably impact the probability of correct detection. In [31] the safety margin to the decision thresholds of the received Phase-Shift Keying (PSK) symbols is maximized subject to a relaxed 1-bit constraint using linear programming for flat-fading channels and extended in [36] for frequency-selective channels. This approach is called the Maximum Safety Margin (MSM) method. The same optimization problem was solved by the Branch-and Bound algorithm in [37] for the special case of 4-PSK. The approach in [34] is based on minimizing an upper bound of the SER under the 1-bit constraint. By the use of a powerful Lemma, the problem was reformulated as a convex optimization problem, of which the solution is discrete.

To the best of our knowledge, the only works that have considered the case of coarsely QCE transmit signals are [38–41]. In [38], we proposed a symbol-wise MSE precoder based on the gradient-descent method under a strict CE constraint or a relaxed polygon constraint. In [39], the authors extended the method in [29] to fit the context of QCE transmit signals. In [40], the authors use a greedy approach for the precoder design while using the symbol-wise MSE as the design criterion. The contribution in [41] addressed the task of QCE precoding in the context of using a single common PA and separate digital phase shifters for the antenna front-ends. The optimization problem consists of designing the QCE precoder while minimizing the MUI, and the idea of constructive interference, [42, 43], is not exploited. In [44], the MSM method was extended to QCE precoding for general constellations and for flat-fading channels. The extension of the MSM method for frequency-selective channels was studied in [45] for Quadrature Amplitude Modulation (QAM) signaling only.

It is worth mentioning that the QCE precoding can be combined with appropriate pulse shaping strategies as in [46, 47] to ensure an efficient spectral confinement. In [48], it was shown that CE precoding is still power efficient even when considering time-based processing. The same investigation can be conducted for the case of QCE precoding. Here, we focus rather on the spatial design problem.

1.2 Outline and Main Contributions

The thesis is organized as follows.

- Chapter 3 is devoted to introduce the mathematical model of the CEQ. Moreover, Price's theorem is extended and applied for the CEQ to obtain its signal statistical properties under the assumption of Gaussian distributed input signals. The main results are summarized in [49]. To ensure minimal distortions, the optimal CEQ is introduced, whose statistical properties can be approximated by the Linear Covariance Approximation (LCA). Finally, a linearized model of the CEQ is derived using Bussgang's theorem.

- Chapter 4 presents the system model of the QCE MU massive MIMO system that is considered throughout the thesis. Two channel models, i.e. the independent and identically distributed (i.i.d.) and the mmW channel models, are also described. Moreover, a potential uplink system model is illustrated.
- Part I, which includes Chapter 5 and Chapter 6, addresses the linear precoding task for flat-fading and frequency-selective channels, respectively. For the precoder design, we choose the MMSE criterion. With the help of Bussgang's theorem, the CEQ can be linearized, which leads to a linear system model with an additional quantization noise. The Wiener Filter (WF) precoder that considers the quantization noise is designed under the assumption of equal and unequal power allocation at the BS antennas. The statistics of the quantization noise are computed by applying Price's theorem or the LCA. Furthermore, the MSE duality between the uplink and the downlink QCE systems is investigated to show that no MSE duality holds and only virtual or approximate duality can be achieved. Linear precoders for the virtual and approximate dual uplink systems are derived.
- Part II, which contains Chapter 7 and Chapter 8, is concerned with the non-linear precoder design for flat-fading and frequency-selective channels, respectively. A significant part of Chapter 7 is published in [44], whereas Chapter 8 is partially published in [45]. The design criterion consists of the safety margin to the decision thresholds. Maximizing the safety margin leads to decreased SER. In contrast to the linear precoding scheme, no precoding matrix is designed but the transmit vector for each given input signal and at a given channel realization is optimized. To take the QCE constraint into account, the entries of the optimized vector should belong to a relaxed convex version of the QCE constraint. The obtained optimization problem is a linear programming problem, for which there exist very efficient methods to solve.
- Chapter 9 discusses the benefits and challenges of QCE systems in terms of power efficiency and spectral shaping. We show the potential power savings of QCE systems compared to the ideal linear systems. Afterwards, we investigate the spectral regrowth in the presence of coarse quantization.

1.3 Notation

Signals, Channels, Filters

$\boldsymbol{\eta}$	Additive White Gaussian Noise (AWGN) vector
\mathbf{F}	Linear equalizer in the dual domain
\mathbf{G}	Receive matrix
\mathbf{H}	Channel matrix
\mathbf{P}	Precoding matrix
\mathbf{r}	Received signal vector
\mathbf{s}	input signal vector
$\hat{\mathbf{s}}$	Estimate of the input signal vector
\mathbf{t}	CEQ output signal vector
\mathbf{T}	Linear transmit matrix in the dual domain
\mathbf{u}	Received filtered signal vector
\mathbf{x}	CEQ input signal vector
\mathbf{y}	Noiseless received signal vector

Numbers and Quantities

B	Block length for the non-linear processing
δ	Safety margin to the decision thresholds at the receiver
L	Number of taps of the channel impulse response matrix
L_p	Number of taps of the linear precoder impulse response matrix
M	Number of single-antenna users
N	Number of transmit antennas at the BS
P_{tx}	Available transmit power
S	Cardinality of the input set
θ	$\frac{\pi}{S}$
T	Block length for the blind estimation at the receiver
T_s	Sample period
T_{sym}	Symbol period

Matrix, vector and complex number operations

$(\bullet)^T$	Transpose of a matrix
$(\bullet)^H$	Hermitian of a matrix
$\text{tr}(\bullet)$	Trace of a matrix
$\text{null}(\bullet)$	Null space of a matrix
$\text{range}(\bullet)$	Range of a matrix
$(\bullet)^*$	Complex conjugate
$\ \bullet\ _p$	p -norm
$\Re\{\bullet\}$	Real part of a complex-valued number
$\Im\{\bullet\}$	Imaginary part of a complex-valued number
$*$	Convolution
\circ	Hadamard product
\otimes	Kronecker product
$\mathbf{0}_N$	N -dimensional zero vector
$\mathbf{1}_N$	N -dimensional all ones vector
\mathbf{I}_N	N -dimensional identity matrix
\mathbf{e}_n	n -th column of the identity matrix
$\mathbf{0}_{N,M}$	$N \times M$ -dimensional zero matrix
\mathbf{X}	Matrix
\mathbf{x}	Vector
x	Scalar
x_i	i -th element of the vector \mathbf{x}
$\text{diag}(\mathbf{X})$	Diagonal matrix containing the diagonal elements of \mathbf{X}
$\text{nondiag}(\mathbf{X})$	$\mathbf{X} - \text{diag}(\mathbf{X})$
$f[t]$	Function f evaluated at discrete instants t
$f(t)$	Function f evaluated at continuous instants t

Quantization

α_q	$1 - \beta_q$
β_q	Distortion factor of the CEQ
\mathbf{d}_Q	Quantization noise after the Bussgang linearization of the CEQ
\mathbf{L}_Q	Linear matrix between the input and output after the Bussgang linearization of the CEQ
q	Resolution in bits
Q	Number of possible discrete phases at the output of the CEQ
Q_{CE}	Element-wise CEQ
Q_ϕ	Element-wise phase quantization
ψ	$\frac{\pi}{Q}$
Ξ_n	Envelope magnitude at the output of the n -th CEQ
Ξ	Diagonal matrix $\sum_{n=1}^N \Xi_n \mathbf{e}_n \mathbf{e}_n^T$

Probability and Stochastics

$p_X(x)$	Probability Density Function (PDF) of the random variable X
$\Pr(\text{statement})$	Probability that the statement holds true
\wedge	Logical and
\vee	Logical or
$E[\bullet]$	Expectation operator
σ_x^2	Variance of the signal x
$\rho_{x,y}$	Correlation factor between x and y
\mathbf{C}_{xy}	Covariance matrix between \mathbf{x} and \mathbf{y} , i.e. $\mathbf{C}_{xy} = E[\mathbf{xy}^H] - E[\mathbf{x}]E[\mathbf{y}^H]$
\mathbf{R}_{xy}	Correlation matrix, i.e. $\mathbf{R}_{xy} = \text{diag}(\mathbf{C}_{xx})^{-1/2}\mathbf{C}_{xy}\text{diag}(\mathbf{C}_{yy})^{-1/2}$
$\mathcal{CN}_{\mathbb{C}}(\mathbf{0}, \mathbf{C})$	Circularly-symmetric complex-valued Gaussian distribution with zero-mean and covariance matrix \mathbf{C}

Sets

\mathbb{C}	Set of complex-valued numbers
\mathbb{S}	Input set, i.e. QAM or PSK
\mathbb{T}	CEQ output set
\mathbb{X}	CEQ input set

2. Motivation

2.1 Power Amplifiers

The PA in wireless communication systems is a device at the end of the transmit front-end to amplify the RF signal and drive it into the antenna. The amplification is in the ideal case without distortion and with 100% efficiency. However, this is not possible with real devices. Therefore, each PA is characterized by two main features, i.e. linearity and efficiency.

Linear amplification avoids introducing strong distortions to the signal and hence avoiding the higher out-of band radiations. To run the PA in the linear region, the input signal has to have a peak magnitude well below the PA output peak. This implies, however, that the supply power is not used fully and power dissipation arises.

Power efficient amplification ensures efficient power usage and hence less requirements on cooling systems at the BS. The power efficiency of the PA is defined as

$$\eta = \frac{P_{\text{out}}}{P_{\text{DC}}}, \quad (2.1)$$

where P_{out} is the output power and P_{DC} is the supply power. It describes how much percentage out of the supply power is converted to the RF power (output power). If the power efficiency is equal to 1, it means that the supply power is totally transmitted to the output. Hence, no power is consumed at the PA in the ideal case.

Both fundamental characteristics depend on the class of the PA. There are different classes of amplifiers of different theoretical efficiency values [50]; that is

- class A ($\eta = 50\%$), class AB ($\eta = 68\%$), class B ($\eta = 78.5\%$), class C ($\eta = 87\%$): behave like linear transconductors
- class D ($\eta = 100\%$), class E ($\eta = 100\%$), class S ($\eta = 100\%$), and M ($\eta = 100\%$) [51]: behave like non-linear switches with high efficiency levels.

There is always a trade-off between the linearity and the power efficiency of the PA. The linear PA is less power efficient than the non-linear PA. However, having a PA input signal of constant magnitude, the PA can be operated in the saturation region while achieving the maximal possible power efficiency and without introducing any distortion to the signal. The PA of class M shows a higher efficiency in real systems compared to other highly efficient PAs. Therefore, we suggest its use in this thesis.

2.2 DACs

The conventional DACs in wireless communication systems are the current-steering DACs thanks to their high-speed characteristic [52]. A current steering DAC consists of a number of

parallel switched current sources. Depending on the digital input signal, the switchers steer the currents to either the output or dump them to ground. The resulting weighted currents are summed at the output and converted into a voltage by a transimpedance stage.

The number of switched current sources grows exponentially with the resolution of the DAC. Therefore, a high resolution leads to a huge number of devices, increased complexity of the corresponding wiring and certainly higher hardware power consumption. Thus, low resolution ensures less power consumption and simplified circuitry.

Since we are interested in having CE signals at the PA input for power efficiency enhancement, the DAC should be designed accordingly. In CE signaling, the information is conveyed by the signal phase. Hence, instead of having two DACs for the inphase and quadrature signal parts each, the DAC should process both signals jointly to recover the phase information. Therefore, we do not talk about cartesian DACs anymore but about polar DACs that are applied in polar transmitters. The DAC can be modeled as a two-step operation. The first step consists of the quantization of the continuous-value input signal. The second step is the conversion from the discrete-time to the continuous-time representation. The quantization operation of the polar DAC can be then modeled by the CEQ, which is described in Chapter 3. Note that the cartesian DAC with resolution of one bit at each dimension is equivalent to the polar DAC with a resolution of two bits; that is only four discrete phase values can be generated at the output.

2.3 Polar Transmitters

Polar transmission can enhance simultaneously the linearity and the efficiency of wireless communication systems [53,54]. The information is conveyed by the magnitude and the phase in contrast to cartesian systems, where the information lies in the inphase and quadrature signal parts. For our study, the magnitude is constant due to the CE property. Thus, the information lies only in the phase and only phase modulation is required.

For the conversion from the digital to the analog world, Digital-to-Time Converter (DTC) are deployed in polar transmitters. The DTC is composed of an array of switchers with cascaded delay elements. The DTC input signal is an instantaneous period that is converted into a certain time delay. The time delay is obtained by turning on the corresponding switch. The number of delay elements defines the number of possible discrete phase values at the quantizer output. In other words, the resolution determines the number of delay elements. The higher the resolution is, the more delay elements we have in the converter. For less power consumption and reduced chip area we restrict our work to low resolution and therefore small number of delay elements in the DTC. Note that our derivations assume that the number of delay elements is a power of two. However, the work can be generalized for all possible numbers of delay elements.

3. On Constant Envelope Quantizers

3.1 One-Dimensional CEQ

The input of the CEQ, denoted by x , is a complex-valued signal in contrast to the quantizers with real-valued inputs that are largely studied in the literature. In the polar representation, the output signal t has a constant magnitude and a quantized version of the input signal phase as

$$\begin{aligned}
 t &= \mathcal{Q}_{\text{CE}}(x) \\
 &= \Xi e^{j \mathcal{Q}_\phi(\arg(x))} \\
 &= \Xi \sum_{k=1}^Q (u(\arg(x) - (2k-2)\psi) - u(\arg(x) - 2k\psi)) e^{j(2k-1)\psi}, \tag{3.1}
 \end{aligned}$$

where Ξ is constant and denotes the magnitude of the CE, $u(\bullet)$ denotes the unit step function as

$$u(\phi) = \begin{cases} 0 & \phi < 0 \\ 1 & \phi \geq 0, \end{cases} \tag{3.2}$$

ψ is defined as

$$\psi = \frac{\pi}{Q}, \tag{3.3}$$

and $\arg(x)$ gives the phase of the complex-valued signal x . So, the information after the CEQ lies only in the phase. The phase quantizer $\mathcal{Q}_\phi(\bullet)$ is a symmetric uniform real-valued quantizer. It is characterized by its resolution q that defines the number of the discrete output phases, i.e Q , that is

$$Q = 2^q. \tag{3.4}$$

In other words, the 2π phase range is divided into Q , $\frac{2\pi}{Q}$ -rotationally symmetric, sectors. The input signal that belongs to the k -th sector gets quantized (mapped) to $\Xi e^{j(2k-1)\psi}$ as shown in Fig. 3.1.

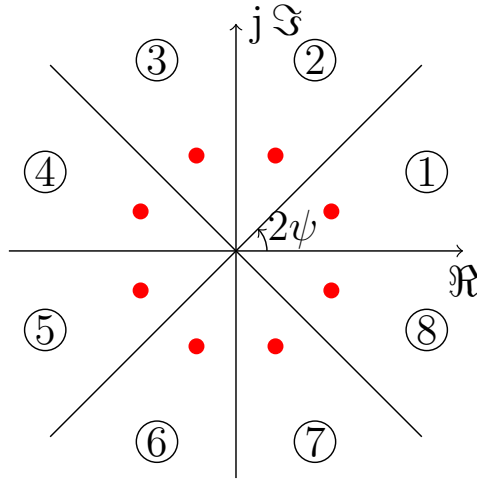


Fig. 3.1: Illustration of the CEQ output set for $q = 3$, i.e. $Q = 8$.

3.2 Multi-Dimensional CEQ

Having an N -dimensional vector \mathbf{x} at the input of an N -dimensional CEQ, the output vector \mathbf{t} is given by

$$\begin{aligned}
 \mathbf{t} &= \mathcal{Q}_{\text{CE}}(\mathbf{x}) \\
 &= \sum_{n=1}^N \mathcal{Q}_{\text{CE}}(x_n) \mathbf{e}_n \\
 &= \sum_{n=1}^N \Xi_n e^{j \mathcal{Q}_\phi(\arg(x_n))} \mathbf{e}_n,
 \end{aligned} \tag{3.5}$$

which corresponds to an element-wise quantization. The output vector can be compactly expressed as

$$\mathbf{t} = \Xi e^{j \mathcal{Q}_\phi(\arg(\mathbf{x}))}, \tag{3.6}$$

where Ξ is a diagonal matrix containing the magnitudes for each output; that is

$$\Xi = \sum_{n=1}^N \Xi_n \mathbf{e}_n \mathbf{e}_n^T. \tag{3.7}$$

3.3 Statistical Properties

In this section, we aim to find the impact of the CEQ on the statistical properties of Gaussian input signals. To this end, we consider N input signals x_n , $n = 1, \dots, N$, that build the N -dimensional vector \mathbf{x} . The N -dimensional CEQ output vector is denoted by \mathbf{t} . We assume that the input signals are joint complex-valued Gaussian distributed with zero mean and variance $\sigma_{x_n}^2$, that is $x_n \sim \mathcal{N}_{\mathbb{C}}(0, \sigma_{x_n}^2)$, $n = 1, \dots, N$. Thus, the Probability Density Function (PDF) of the complex-valued Gaussian distributed random variable $x_n = r_n e^{j\phi_n}$, $n = 1, \dots, N$,

reads as, [55, Theorem 2.4],

$$p_{X_n}(r_n, \phi_n) = \frac{1}{\pi \sigma_{x_n}^2} e^{-\frac{r_n^2}{\sigma_{x_n}^2}}. \quad (3.8)$$

Additionally, the complex-valued correlation coefficient ρ_{x_i, x_j} between the input signals x_i and x_j , $i, j = 1, \dots, N$, is defined as

$$\rho_{x_i, x_j} = \frac{\mathbb{E}[x_i x_j^*]}{\sigma_{x_i} \sigma_{x_j}} \quad (3.9)$$

$$= \frac{1}{\sigma_{x_i} \sigma_{x_j}} \int_0^{2\pi} \int_0^{2\pi} \int_0^\infty \int_0^\infty r_i^2 r_j^2 e^{j(\phi_i - \phi_j)} p_{X_i, X_j}(r_i, r_j, \phi_i, \phi_j) dr_i dr_j d\phi_j d\phi_i, \quad (3.10)$$

where the joint PDF for $x_i \neq x_j$ is expressed as

$$p_{X_i, X_j, i \neq j}(r_i, r_j, \phi_i, \phi_j) = \frac{1}{\pi^2 (1 - |\rho_{x_i, x_j}|^2) \sigma_{x_i}^2 \sigma_{x_j}^2} e^{-\frac{1}{1 - |\rho_{x_i, x_j}|^2} \left(\frac{r_i^2}{\sigma_{x_i}^2} + \frac{r_j^2}{\sigma_{x_j}^2} - \frac{2r_i r_j \Re\{\rho_{x_i, x_j} e^{-j(\phi_i - \phi_j)}\}}{\sigma_{x_i} \sigma_{x_j}} \right)}, \quad (3.11)$$

as explained in A1. Moreover, we define the input covariance matrix as

$$\begin{aligned} \mathbf{C}_{\mathbf{xx}} &= \mathbb{E}[\mathbf{xx}^H] \\ &= \begin{bmatrix} \sigma_{x_1}^2 & \cdots & \mathbb{E}[x_1 x_N^*] \\ \vdots & \ddots & \\ \mathbb{E}[x_N x_1^*] & \cdots & \sigma_{x_N}^2 \end{bmatrix}. \end{aligned} \quad (3.12)$$

Consequently, the input correlation matrix reads as

$$\begin{aligned} \mathbf{R}_{\mathbf{xx}} &= \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \mathbf{C}_{\mathbf{xx}} \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \\ &= \begin{bmatrix} 1 & \cdots & \rho_{x_1, x_N} \\ \vdots & \ddots & \\ \rho_{x_N, x_1} & \cdots & 1 \end{bmatrix}. \end{aligned} \quad (3.13)$$

Assuming Gaussian input signals, our target consists of finding the expressions of

- the covariance matrix $\mathbf{C}_{\mathbf{tx}} = \mathbb{E}[\mathbf{tx}^H]$ between the QCE signal vector and the unquantized signal vector, i.e. $\mathbb{E}[t_i x_j^*]$, $i, j = 1, \dots, N$, and
- the covariance matrix $\mathbf{C}_{\mathbf{tt}} = \mathbb{E}[\mathbf{tt}^H]$ of the QCE signal vector for a given input correlation matrix $\mathbf{R}_{\mathbf{xx}}$, i.e. $\mathbb{E}[t_i t_j^*]$, $i, j = 1, \dots, N$.

To this end, we introduce Price's theorem.

3.3.1 Price's Theorem

Price's theorem, first introduced in [27], consists in ordinary or partial differential equations of different orders to describe the expected value of the product of n non-linear functions of n real-valued jointly distributed Gaussian random variables. Each non-linear function depends only on one random variable. The ℓ -th derivative of the expected value with respect

to (w.r.t.) the covariances between the random variables equates the expected value of the product of the ℓ -th functions' derivatives w.r.t. the corresponding random variables. The theorem was reformulated in [56] to obtain a partial differential equation relating the derivative of the expected value of a non-linear function of two jointly distributed Gaussian random variables w.r.t. their covariance and the expected value of the second order derivative of the non-linear function w.r.t. the random variables. This relationship holds under a certain boundedness constraint given in [57]. The number of the real-valued joint Gaussian variables n determines the number of the differential equations to solve, that is $\binom{n}{2} = n(n-1)/2$. The modified version of Price's theorem in [58] extends the expression for non-linear functions of more than two real-valued variables to end up with a single ordinary differential equation involving all derivatives' orders. In [59], Price's theorem was extended for two circularly symmetric complex-valued Gaussian random variables and generalized in [60] under no circular symmetry assumption. The two latter extensions describe the complex-valued version of Price's theorem for the derivative of first order only. We could stick to the version of Price's theorem for complex-valued Gaussian random variables for the first order only given in [60]. However, we extend the theorem in [58] to obtain Price's theorem for two circularly symmetric complex-valued Gaussian random variables that involve all derivatives' orders.

Theorem 1. *Let $f(\mathbf{x}) = f(x_1, x_2)$ be a function of two joint complex-valued circularly symmetric Gaussian variables. The covariance matrix and the correlation matrix of the vector \mathbf{x} are denoted by $\mathbf{C}_{\mathbf{xx}}$ and $\mathbf{R}_{\mathbf{xx}}$ of elements ρ_{x_i, x_j} , $i, j = 1, 2$. The off-diagonal entries of $\mathbf{C}_{\mathbf{xx}}$ are multiplied by a perturbing term v . Then, it holds that*

$$\begin{aligned} \frac{d^\ell \mathbb{E}[f(x_i, x_j)]}{dv^\ell} &= \sigma_{x_i}^\ell \sigma_{x_j}^\ell \mathbb{E}_v \left[\Re\{\rho_{x_i, x_j}\}^\ell \left(\frac{\partial^{2\ell} f(x_i, x_j)}{\partial x_i^\ell \partial x_j^{*\ell}} + \frac{\partial^{2\ell} f(x_i, x_j)}{\partial x_i^{*\ell} \partial x_j^\ell} \right) \right. \\ &\quad \left. + \mathfrak{J}^\ell \{\rho_{x_i, x_j}\}^\ell \left(\frac{\partial^{2\ell} f(x_i, x_j)}{\partial x_i^\ell \partial x_j^{*\ell}} - \frac{\partial^{2\ell} f(x_i, x_j)}{\partial x_i^{*\ell} \partial x_j^\ell} \right) \right], \end{aligned} \quad (3.14)$$

where \mathbb{E}_v denotes the expectation operation based on the resulting perturbed PDF.

The proof is provided in A2. ■

Price's theorem is usually an alternative way to find the analytical expression for $\mathbb{E}[f(x_i, x_j)]$, which can be obtained by integrating the first-order derivative w.r.t. v . Therefore, we state Theorem 3.14 for the special case $\ell = 1$

$$\frac{d \mathbb{E}[f(x_i, x_j)]}{dv} = \sigma_{x_i} \sigma_{x_j} \mathbb{E}_v \left[\rho_{x_i, x_j} \frac{\partial^2 f(x_i, x_j)}{\partial x_i \partial x_j^*} + \rho_{x_i, x_j}^* \frac{\partial^2 f(x_i, x_j)}{\partial x_i^* \partial x_j} \right]. \quad (3.15)$$

The latter expression can be reformulated as

$$\begin{aligned} d \mathbb{E}[f(x_i, x_j)] &= \sigma_{x_i} \sigma_{x_j} \mathbb{E}_v \left[\frac{\partial^2 f(x_i, x_j)}{\partial x_i \partial x_j^*} \right] \rho_{x_i, x_j} dv + \sigma_{x_i} \sigma_{x_j} \mathbb{E}_v \left[\frac{\partial^2 f(x_i, x_j)}{\partial x_i^* \partial x_j} \right] \rho_{x_i, x_j}^* dv \\ &= \sigma_{x_i} \sigma_{x_j} \mathbb{E}_v \left[\frac{\partial^2 f(x_i, x_j)}{\partial x_i \partial x_j^*} \right] d\rho'_{x_i, x_j} + \sigma_{x_i} \sigma_{x_j} \mathbb{E}_v \left[\frac{\partial^2 f(x_i, x_j)}{\partial x_i^* \partial x_j} \right] d\rho'^*_{x_i, x_j}, \end{aligned} \quad (3.16)$$

where $\rho'_{x_i, x_j} = \nu \rho_{x_i, x_j}$. Therefore, we obtain

$$\frac{\partial \mathbf{E}[f(x_i, x_j)]}{\partial \rho'_{x_i, x_j}} = \sigma_{x_i} \sigma_{x_j} \mathbf{E}_\nu \left[\frac{\partial^2 f(x_i, x_j)}{\partial x_i \partial x_j^*} \right], \quad (3.17)$$

$$\frac{\partial \mathbf{E}[f(x_i, x_j)]}{\partial \rho'^*_{x_i, x_j}} = \sigma_{x_i} \sigma_{x_j} \mathbf{E}_\nu \left[\frac{\partial^2 f(x_i, x_j)}{\partial x_i^* \partial x_j} \right], \quad (3.18)$$

which simplify for the case of $\nu = 1$ to

$$\frac{\partial \mathbf{E}[f(x_i, x_j)]}{\partial \rho_{x_i, x_j}} = \sigma_{x_i} \sigma_{x_j} \mathbf{E} \left[\frac{\partial^2 f(x_i, x_j)}{\partial x_i \partial x_j^*} \right], \quad (3.19)$$

$$\frac{\partial \mathbf{E}[f(x_i, x_j)]}{\partial \rho^*_{x_i, x_j}} = \sigma_{x_i} \sigma_{x_j} \mathbf{E} \left[\frac{\partial^2 f(x_i, x_j)}{\partial x_i^* \partial x_j} \right]. \quad (3.20)$$

The expression in (3.19) was also obtained in [59] under the circular symmetry assumption, while (3.19) and (3.20) were obtained in [60] without any circular symmetry assumption. This implies that Theorem 1 holds true for the first-order derivative for general complex-valued signals, too.

In the polar representation, (3.19) and (3.20) read as

$$\frac{\partial \mathbf{E}[f(x_i, x_j)]}{\partial \rho_{x_i, x_j}} = \frac{1}{4} \sigma_{x_i} \sigma_{x_j} \cdot \mathbf{E} \left[e^{-j(\phi_i - \phi_j)} \left(\frac{\partial^2 f(x_i, x_j)}{\partial r_i \partial r_j} + \frac{1}{r_i r_j} \frac{\partial^2 f(x_i, x_j)}{\partial \phi_i \partial \phi_j} + j \frac{1}{r_j} \frac{\partial^2 f(x_i, x_j)}{\partial r_i \partial \phi_j} - j \frac{1}{r_i} \frac{\partial^2 f(x_i, x_j)}{\partial \phi_i \partial r_j} \right) \right] \quad (3.21)$$

$$\frac{\partial \mathbf{E}[f(x_i, x_j)]}{\partial \rho^*_{x_i, x_j}} = \frac{1}{4} \sigma_{x_i} \sigma_{x_j} \cdot \mathbf{E} \left[e^{j(\phi_i - \phi_j)} \left(\frac{\partial^2 f(x_i, x_j)}{\partial r_i \partial r_j} + \frac{1}{r_i r_j} \frac{\partial^2 f(x_i, x_j)}{\partial \phi_i \partial \phi_j} - j \frac{1}{r_j} \frac{\partial^2 f(x_i, x_j)}{\partial r_i \partial \phi_j} + j \frac{1}{r_i} \frac{\partial^2 f(x_i, x_j)}{\partial \phi_i \partial r_j} \right) \right]. \quad (3.22)$$

The conversion of the second order derivative of a function w.r.t. two complex-valued variables to the second order derivatives w.r.t. to their corresponding polar arguments is detailed in A3.

3.3.2 Cross-correlation between QCE Signal and Unquantized Signal

The cross-correlation factor between x_i and t_i is expressed as

$$\begin{aligned}
\mathbb{E}[t_i x_i^*] &= \int_0^{2\pi} \int_0^\infty \Xi_i e^{j\mathcal{Q}_\phi(\phi_i)r_i} e^{-j\phi_i} p_{X_i}(r_i, \phi_i) r_i dr_i d\phi_i \\
&= \frac{1}{\pi\sigma_{x_i}^2} \int_0^\infty \sum_{k=1}^Q \int_{(2k-2)\psi}^{2k\psi} \Xi_i e^{j(2k-1)\psi} r_i^2 e^{-j\phi_i} e^{-\frac{r_i^2}{\sigma_{x_i}^2}} d\phi_i dr_i \\
&= \frac{1}{\pi\sigma_{x_i}^2} \int_0^\infty r_i^2 e^{-\frac{r_i^2}{\sigma_{x_i}^2}} dr_i \sum_{k=1}^Q \Xi_i e^{j(2k-1)\psi} \int_{(2k-2)\psi}^{2k\psi} e^{-j\phi_i} d\phi_i \\
&= \frac{1}{\pi\sigma_{x_i}^2} \frac{\sigma_{x_i}^2}{4} \sqrt{\pi\sigma_{x_i}^2} \sum_{k=1}^Q \Xi_i e^{j(2k-1)\psi} \left(e^{j(\frac{\pi}{2}-2k\psi)} - e^{j(\frac{\pi}{2}-(2k-2)\psi)} \right) \\
&= \frac{\sigma_{x_i}}{4\sqrt{\pi}} \sum_{k=1}^Q \Xi_i \left(e^{j(\frac{\pi}{2}-\psi)} - e^{j(\frac{\pi}{2}+\psi)} \right) \\
&= \frac{\Xi_i \sigma_{x_i}}{2\sqrt{\pi}} Q \sin(\psi).
\end{aligned} \tag{3.23}$$

The computation of the cross-correlation factor between t_i and x_j with $i \neq j$ implies the evaluation of the following quadruple integral

$$\begin{aligned}
\mathbb{E}[t_i x_j^*] &= \int_0^{2\pi} \int_0^{2\pi} \int_0^\infty \int_0^\infty \Xi_i e^{j\mathcal{Q}_\phi(\phi_i)} r_j e^{-j\phi_j} p_{X_i, X_j, i \neq j}(r_i, r_j, \phi_i, \phi_j) r_i r_j dr_i dr_j d\phi_i d\phi_j \\
&= \sum_{k=1}^Q \int_{(2k-2)\psi}^{2k\psi} \Xi_i e^{j(2k-1)\psi} \int_0^{2\pi} e^{-j\phi_j} \int_0^\infty \int_0^\infty r_i r_j^2 p_{X_i, X_j, i \neq j}(r_i, r_j, \phi_i, \phi_j) dr_i dr_j d\phi_j d\phi_i,
\end{aligned} \tag{3.24}$$

which is not straightforward to calculate analytically. Therefore, we make use of Price's theorem. To this end, we compute the following derivatives

$$\frac{\partial^2 t_i x_j^*}{\partial r_i \partial r_j} = 0, \tag{3.25}$$

$$\frac{\partial^2 t_i x_j^*}{\partial r_i \partial \phi_j} = 0, \tag{3.26}$$

$$\frac{\partial^2 t_i x_j^*}{\partial \phi_i \partial r_j} = \Xi_i \sum_{k=1}^Q (\delta(\phi_i - (2k-2)\psi) - \delta(\phi_i - 2k\psi)) e^{j((2k-1)\psi - \phi_j)}, \tag{3.27}$$

$$\begin{aligned}
\frac{\partial^2 t_i x_j^*}{\partial \phi_i \partial \phi_j} &= -j \Xi_i \sum_{k=1}^Q (\delta(\phi_i - (2k-2)\psi) - \delta(\phi_i - 2k\psi)) r_j e^{j((2k-1)\psi - \phi_j)} \\
&= -j r_j \frac{\partial^2 t_i x_j^*}{\partial \phi_i \partial r_j}.
\end{aligned} \tag{3.28}$$

Hence, we apply Price's theorem in the polar representation from (3.21) and (3.22). We start with the first property and we get

$$\begin{aligned}
 \frac{\partial \mathbb{E} [t_i x_j^*]}{\partial \rho_{x_i, x_j}} &= \sigma_{x_i} \sigma_{x_j} \mathbb{E} \left[\frac{\partial^2 t_i x_j^*}{\partial x_i \partial x_j^*} \right] \\
 &= \frac{1}{4} \sigma_{x_i} \sigma_{x_j} \mathbb{E} \left[e^{-j(\phi_i - \phi_j)} \left(\frac{-j}{r_i} \frac{\partial^2 t_i x_j^*}{\partial \phi_i \partial r_j} + \frac{1}{r_i r_j} \frac{\partial^2 t_i x_j^*}{\partial \phi_i \partial \phi_j} \right) \right] \\
 &= -j \frac{1}{2} \sigma_{x_i} \sigma_{x_j} \mathbb{E} \left[e^{-j(\phi_i - \phi_j)} \frac{1}{r_i} \frac{\partial^2 t_i x_j^*}{\partial \phi_i \partial r_j} \right] \\
 &= -j \frac{1}{2} \sigma_{x_i} \sigma_{x_j} \int_0^{2\pi} \int_0^{2\pi} \int_0^\infty \int_0^\infty e^{-j(\phi_i - \phi_j)} \frac{\partial^2 t_i x_j^*}{\partial \phi_i \partial r_j} p_{X_i, X_j, i \neq j}(r_i, r_j, \phi_i, \phi_j) r_j \, dr_i \, dr_j \, d\phi_j \, d\phi_i.
 \end{aligned} \tag{3.29}$$

Plugging (3.27) in (3.29) results in

$$\begin{aligned}
 \frac{\partial \mathbb{E} [t_i x_j^*]}{\partial \rho_{x_i, x_j}} &= -j \frac{1}{2} \sigma_{x_i} \sigma_{x_j} \Xi_i \int_0^{2\pi} \sum_{k=1}^Q e^{j((2k-1)\psi - \phi_j)} \cdot \\
 &\quad \left(e^{-j((2k-2)\psi - \phi_j)} \int_0^\infty \int_0^\infty p_{X_i, X_j, i \neq j}(r_i, r_j, (2k-2)\psi, \phi_j) r_j \, dr_i \, dr_j \right. \\
 &\quad \left. - e^{-j(2k\psi - \phi_j)} \int_0^\infty \int_0^\infty p_{X_i, X_j, i \neq j}(r_i, r_j, 2k\psi, \phi_j) r_j \, dr_i \, dr_j \right) \\
 &= -j \frac{1}{2} \sigma_{x_i} \sigma_{x_j} \Xi_i \int_0^{2\pi} \sum_{k=1}^Q (e^{j\psi} w_1((2k-2)\psi - \phi_j) - e^{-j\psi} w_1(2k\psi - \phi_j)) \, d\phi_j,
 \end{aligned} \tag{3.30}$$

where

$$w_1(\phi_i - \phi_j) = \int_0^\infty \int_0^\infty p_{X_i, X_j, i \neq j}(r_i, r_j, \phi_i, \phi_j) r_j \, dr_i \, dr_j. \tag{3.31}$$

The latter double integral is computed analytically in (A2) and the obtained expression is

$$w_1(\phi_i - \phi_j) = \frac{\sqrt{1 - |\rho_{x_i, x_j}|^2}}{4\pi\sqrt{\pi}\sigma_{x_i}(1 - \kappa(\phi_i - \phi_j))}, \tag{3.32}$$

where

$$\kappa(\phi) = \Re\{\rho_{x_i, x_j} e^{-j\phi}\}. \tag{3.33}$$

Hence, (3.30) reduces to

$$\begin{aligned}
\frac{\partial \mathbb{E} [t_i x_j^*]}{\partial \rho_{x_i, x_j}} &= -j \frac{1}{2} \sigma_{x_i} \sigma_{x_j} \Xi_i \int_0^{2\pi} \sum_{k=1}^Q w_1(2k\psi - \phi_j) (e^{j\psi} - e^{-j\psi}) d\phi_j \\
&= \sigma_{x_i} \sigma_{x_j} \Xi_i \sin(\psi) \int_0^{2\pi} \sum_{k=1}^Q w_1(2k\psi - \phi_j) d\phi_j \\
&\stackrel{(a)}{=} \sigma_{x_i} \sigma_{x_j} \Xi_i \sin(\psi) Q \int_0^{2\pi} w_1(\phi) d\phi \\
&= \frac{\sqrt{1 - |\rho_{x_i, x_j}|^2} \sigma_{x_j} \Xi_i}{4\pi\sqrt{\pi}} \sin(\psi) Q \int_0^{2\pi} \frac{1}{1 - \Re\{\rho_{x_i, x_j} e^{-j\phi}\}} d\phi \\
&\stackrel{(b)}{=} \frac{\sqrt{1 - |\rho_{x_i, x_j}|^2} \sigma_{x_j} \Xi_i}{4\pi\sqrt{\pi}} \sin(\psi) Q \int_0^{2\pi} \frac{1}{1 - |\rho_{x_i, x_j}| \cos(\phi')} d\phi' \\
&= \frac{\sqrt{1 - |\rho_{x_i, x_j}|^2} \sigma_{x_j} \Xi_i}{4\pi\sqrt{\pi}} \sin(\psi) Q \int_0^{\pi} \frac{2}{1 - |\rho_{x_i, x_j}|^2 \cos^2(\phi')} d\phi' \\
&\stackrel{(c)}{=} \frac{Q}{2\sqrt{\pi}} \sin(\psi) \Xi_i \sigma_{x_j}, \tag{3.34}
\end{aligned}$$

where in (a) we introduce the variable $\phi = 2k\psi - \phi_j$. Since $w_1(\phi)$ is 2π -periodic, the integral operation from $(2k\psi - 2\pi)$ to $2k\psi$ is equivalent to the integration from 0 to 2π . In (b) we define $\phi' = \arg(\rho_{x_i, x_j}) - \phi$. Again the integral boundaries do not change due to the 2π -periodicity of the integrand. In (c) we use the equality

$$\int_{-1}^1 \frac{2}{1 - a^2 x^2} \frac{1}{\sqrt{1 - x^2}} dx = \frac{2\pi}{1 - a^2}, \quad \forall a \in \mathbb{R} \quad \text{with } a \leq 1. \tag{3.35}$$

Additionally, we obtain

$$\begin{aligned}
\frac{\partial \mathbb{E} [t_i x_j^*]}{\partial \rho_{x_i, x_j}^*} &= \sigma_{x_i} \sigma_{x_j} \mathbb{E} \left[\frac{\partial^2 t_i x_j^*}{\partial x_i^* \partial x_j} \right] \\
&= \frac{1}{4} \sigma_{x_i} \sigma_{x_j} \mathbb{E} \left[e^{j(\phi_i - \phi_j)} \left(\frac{j}{r_i} \frac{\partial^2 t_i x_j^*}{\partial \phi_i \partial r_j} + \frac{1}{r_i r_j} \frac{\partial^2 t_i x_j^*}{\partial \phi_i \partial \phi_j} \right) \right] \\
&= 0. \tag{3.36}
\end{aligned}$$

Since $\mathbb{E} [t_i x_j^*]$ depends only on ρ_{x_i, x_j} and not on ρ_{x_i, x_j}^* , we can get the covariance between t_i and x_j by integrating the expression in (3.34) w.r.t. ρ_{x_i, x_j} as

$$\begin{aligned}
\mathbb{E} [t_i x_j^*] &= \int_0^{\rho_{x_i, x_j}} \frac{\partial \mathbb{E} [t_i x_j^*]}{\partial \rho_{x_i, x_j}'} d\rho_{x_i, x_j}' \\
&= \frac{Q}{2\sqrt{\pi}} \sin(\psi) \Xi_i \sigma_{x_j} \rho_{x_i, x_j}. \tag{3.37}
\end{aligned}$$

3.3.3 Cross-correlation between Two QCE Signals

The auto-correlation of the CEQ output t_i is given by

$$\begin{aligned} \mathbb{E} [|t_i|^2] &= \int_0^{2\pi} \int_0^\infty |t_i|^2 p_{X_i}(r_i, \phi_i) r_i \, dr_i \, d\phi_i \\ &= |t_i|^2 \int_0^{2\pi} \int_0^\infty p_{X_i}(r_i, \phi_i) r_i \, dr_i \, d\phi_i \\ &= |t_i|^2 = \Xi_i^2. \end{aligned} \quad (3.38)$$

However, the computation of the cross-correlation between t_i and t_j with $i \neq j$ turns to be more complicated

$$\begin{aligned} \mathbb{E} [t_i t_j^*] &= \int_0^{2\pi} \int_0^{2\pi} \int_0^\infty \int_0^\infty \Xi_i \Xi_j e^{j(\mathcal{Q}_\phi(\phi_i) - \mathcal{Q}_\phi(\phi_j))} p_{X_i, X_j, i \neq j}(r_i, r_j, \phi_i, \phi_j) r_i r_j \, dr_i \, dr_j \, d\phi_i \, d\phi_j \\ &= \sum_{k=1}^Q \sum_{k'=1}^Q \int_{(2k-2)\psi}^{2k\psi} \int_{(2k'-2)\psi}^{2k'\psi} \Xi_i \Xi_j e^{j(2k-2k')\psi} \int_0^\infty \int_0^\infty r_i r_j p_{X_i, X_j, i \neq j}(r_i, r_j, \phi_i, \phi_j) \, dr_i \, dr_j \, d\phi_j \, d\phi_i. \end{aligned} \quad (3.39)$$

Therefore, we make use of Price's theorem to compute the cross-correlation between t_i and t_j with $i \neq j$. Since the information after the CEQ lies only in the phase, all derivatives w.r.t. the radii vanish and only the derivative w.r.t. to the phases is different from zero, i.e.

$$\frac{\partial^2 t_i t_j^*}{\partial r_i \partial r_j} = 0, \quad (3.40)$$

$$\frac{\partial^2 t_i t_j^*}{\partial r_i \partial \phi_j} = 0, \quad (3.41)$$

$$\frac{\partial^2 t_i t_j^*}{\partial \phi_i \partial r_j} = 0, \quad (3.42)$$

$$\begin{aligned} \frac{\partial^2 t_i t_j^*}{\partial \phi_i \partial \phi_j} &= \sum_{k=1}^Q \sum_{k'=1}^Q (\delta(\phi_i - (2k-2)\psi) - \delta(\phi_i - 2k\psi)) \cdot \\ &\quad (\delta(\phi_j - (2k'-2)\psi) - \delta(\phi_j - 2k'\psi)) \Xi_i \Xi_j e^{j(2(k-k')\psi)}. \end{aligned} \quad (3.43)$$

Hence, we apply Price's theorem in (3.21) and (3.22). We get

$$\begin{aligned} \frac{\partial \mathbb{E} [t_i t_j^*]}{\partial \rho_{x_i, x_j}} &= \sigma_{x_i} \sigma_{x_j} \mathbb{E} \left[\frac{\partial^2 t_i t_j^*}{\partial x_i \partial x_j^*} \right] \\ &= \frac{1}{4} \sigma_{x_i} \sigma_{x_j} \mathbb{E} \left[e^{-j(\phi_i - \phi_j)} \frac{1}{r_i r_j} \frac{\partial^2 t_i t_j^*}{\partial \phi_i \partial \phi_j} \right] \\ &= \frac{1}{4} \sigma_{x_i} \sigma_{x_j} \int_0^{2\pi} \int_0^{2\pi} \int_0^\infty \int_0^\infty e^{-j(\phi_i - \phi_j)} \frac{\partial^2 t_i t_j^*}{\partial \phi_i \partial \phi_j} p_{X_i, X_j, i \neq j}(r_i, r_j, \phi_i, \phi_j) \, dr_i \, dr_j \, d\phi_i \, d\phi_j \end{aligned} \quad (3.44)$$

Plugging (3.43) in (3.44) leads to

$$\begin{aligned}
\frac{\partial E [t_i t_j^*]}{\partial \rho_{x_i, x_j}} &= \frac{1}{4} \sigma_{x_i} \sigma_{x_j} \Xi_i \Xi_j \sum_{k=1}^Q \sum_{k'=1}^Q e^{j(2(k-k')\psi)} \\
&\quad \left(e^{-j(2(k-k')\psi)} \int_0^\infty \int_0^\infty p_{X_i, X_j, i \neq j}(r_i, r_j, (2k-2)\psi, (2k'-2)\psi) dr_i dr_j \right. \\
&\quad + e^{-j(2(k-k')\psi)} \int_0^\infty \int_0^\infty p_{X_i, X_j, i \neq j}(r_i, r_j, 2k\psi, 2k'\psi) dr_i dr_j \\
&\quad - e^{-j(2(k-k'-1)\psi)} \int_0^\infty \int_0^\infty p_{X_i, X_j, i \neq j}(r_i, r_j, (2k-2)\psi, 2k'\psi) dr_i dr_j \\
&\quad \left. - e^{-j(2(k-k'+1)\psi)} \int_0^\infty \int_0^\infty p_{X_i, X_j, i \neq j}(r_i, r_j, 2k\psi, (2k'-2)\psi) dr_i dr_j \right) \\
&= \frac{1}{4} \sigma_{x_i} \sigma_{x_j} \Xi_i \Xi_j \sum_{k=1}^Q \sum_{k'=1}^Q \left(2w_2(2(k-k')\psi) \right. \\
&\quad \left. - e^{j2\psi} w_2(2(k-k'-1)\psi) - e^{-j2\psi} w_2(2(k-k'+1)\psi) \right), \quad (3.45)
\end{aligned}$$

where

$$w_2(\phi_i - \phi_j) = \int_0^\infty \int_0^\infty p_{X_i, X_j, i \neq j}(r_i, r_j, \phi_i, \phi_j) dr_i dr_j. \quad (3.46)$$

The latter double integral is computed analytically in (A3) and the obtained expression is

$$w_2(\phi_i - \phi_j) = \frac{1}{\pi \sigma_{x_i} \sigma_{x_j} \sqrt{1 - \kappa(\phi_i - \phi_j)^2}} \left(\frac{1}{4} + \frac{1}{2\pi} \arcsin(\kappa(\phi_i - \phi_j)) \right), \quad (3.47)$$

where $\kappa(\phi)$ is defined in (3.33). Thus, (3.45) can be rewritten as

$$\begin{aligned}
\frac{\partial E [t_i t_j^*]}{\partial \rho_{x_i, x_j}} &= \frac{1}{4} \sigma_{x_i} \sigma_{x_j} \Xi_i \Xi_j Q \sum_{\Delta k=0}^{Q-1} (2w_2(2\Delta k\psi) - e^{j2\psi} w_2(2(\Delta k-1)\psi) - e^{-j2\psi} w_2(2(\Delta k+1)\psi)) \\
&= \frac{\Xi_i \Xi_j Q}{4\pi} \sum_{\Delta k=0}^{Q/2-1} \left(\frac{1}{\sqrt{1 - \kappa(2\Delta k\psi)^2}} - \frac{e^{j2\psi}}{2\sqrt{1 - \kappa(2(\Delta k-1)\psi)^2}} - \frac{e^{-j2\psi}}{2\sqrt{1 - \kappa(2(\Delta k+1)\psi)^2}} \right) \\
&= \frac{\Xi_i \Xi_j Q}{4\pi} \sum_{\Delta k=0}^{Q/2-1} \frac{1}{\sqrt{1 - \kappa(2\Delta k\psi)^2}} \left(1 - \frac{1}{2} (e^{-j(2\psi)} - e^{j(2\psi)}) \right) \\
&= \frac{\Xi_i \Xi_j Q}{4\pi} (1 - \cos(2\psi)) \sum_{\Delta k=0}^{Q/2-1} \frac{1}{\sqrt{1 - \kappa(2\Delta k\psi)^2}} \\
&= \frac{\Xi_i \Xi_j Q}{2\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} \frac{1}{\sqrt{1 - \kappa(2\Delta k\psi)^2}}. \quad (3.48)
\end{aligned}$$

Similarly, we compute $\frac{\partial \mathbb{E}[t_i t_j^*]}{\partial \rho_{x_i, x_j}^*}$ and get

$$\frac{\partial \mathbb{E}[t_i t_j^*]}{\partial \rho_{x_i, x_j}^*} = \frac{\Xi_i \Xi_j Q}{2\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} \frac{e^{j4\Delta k\psi}}{\sqrt{1 - \kappa(2\Delta k\psi)^2}}. \quad (3.49)$$

Since $\frac{\partial^2 \mathbb{E}[t_i t_j^*]}{\partial \rho_{x_i, x_j} \partial \rho_{x_i, x_j}^*} = \frac{\partial^2 \mathbb{E}[t_i t_j^*]}{\partial \rho_{x_i, x_j}^* \partial \rho_{x_i, x_j}}$, there exists a potential for the vector field built by the partial derivatives. Hence, the expression of the covariance between t_i and t_j is given by

$$\begin{aligned} \mathbb{E}[t_i t_j^*] &= \int_0^{\rho_{x_i, x_j}} \frac{\partial \mathbb{E}[t_i t_j^*]}{\partial \rho'_{x_i, x_j}} d\rho'_{x_i, x_j} = \int_0^{\rho_{x_i, x_j}^*} \frac{\partial \mathbb{E}[t_i t_j^*]}{\partial \rho_{x_i, x_j}^*} d\rho_{x_i, x_j}^* \\ &= \int_0^{\rho_{x_i, x_j}} \frac{\Xi_i \Xi_j Q}{2\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} \frac{1}{\sqrt{1 - \kappa(2\Delta k\psi)^2}} d\rho'_{x_i, x_j} \\ &= \frac{\Xi_i \Xi_j Q}{2\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} \int_0^{\kappa(2\Delta k\psi)} \frac{2 e^{j2\Delta k\psi}}{\sqrt{1 - \kappa'(2\Delta k\psi)^2}} d\kappa'(2\Delta k\psi) \\ &= \frac{\Xi_i \Xi_j Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \arcsin(\kappa(2\Delta k\psi)) \\ &\stackrel{(3.33)}{=} \frac{\Xi_i \Xi_j Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \arcsin(\Re\{\rho_{x_i, x_j} e^{-j2\Delta k\psi}\}). \end{aligned} \quad (3.50)$$

3.3.4 Numerical Validation

In this section, we aim to check the correctness of the expressions in (3.37) and (3.50) by comparing them with numerical results. To this end, we take two sequences of length $N_x = 10^6$ of two arbitrary complex-valued Gaussian signals x_1 and x_2 of variances $\sigma_{x_1}^2$ and $\sigma_{x_2}^2$. We pass them through a two-dimensional CEQ to obtain t_1 and t_2 . The input signals are correlated with a given input correlation coefficient $\rho_{\text{in}} = \rho_{x_1, x_2}$. With Monte-Carlo simulations we compute numerically the cross-correlation coefficient $\rho_{\text{cross}} = \rho_{t_1, x_2} = \mathbb{E}[t_1 x_2^*] / \Xi_1 \sigma_{x_2}$ and the output correlation coefficient, i.e. $\rho_{\text{out}} = \rho_{t_1, t_2} = \mathbb{E}[t_1 t_2^*] / \Xi_1 \Xi_2$. The obtained results (dashed lines with markers) are compared to the derived closed-form expressions (solid lines) in Fig. 3.2 for $Q = 4, 8, \infty$. Almost the same distortion behavior is observed between $Q = 16$ and $Q = \infty$. For this reason the results related to $Q = 16$ are dropped out. Since we are dealing with complex-valued correlation coefficients, we need 4 plots to study the relationship between ρ_{in} and ρ_{cross} and another 4 plots for ρ_{in} and ρ_{out} . As can be seen, the numerical simulations show the correctness of the derived expressions. Moreover, it can be deduced that the phases of the output correlation coefficient and the cross-correlation coefficient remain almost unchanged compared to the phase of the input correlation coefficient. However, most of the distortion is concerning the magnitude depending on Q .

3.3.5 In a Nutshell

In summary, applying (3.37) and (3.50) for the multi-dimensional case, the signal statistical properties of an N -dimensional CEQ of a Gaussian input signal \mathbf{x} and an output signal \mathbf{t}

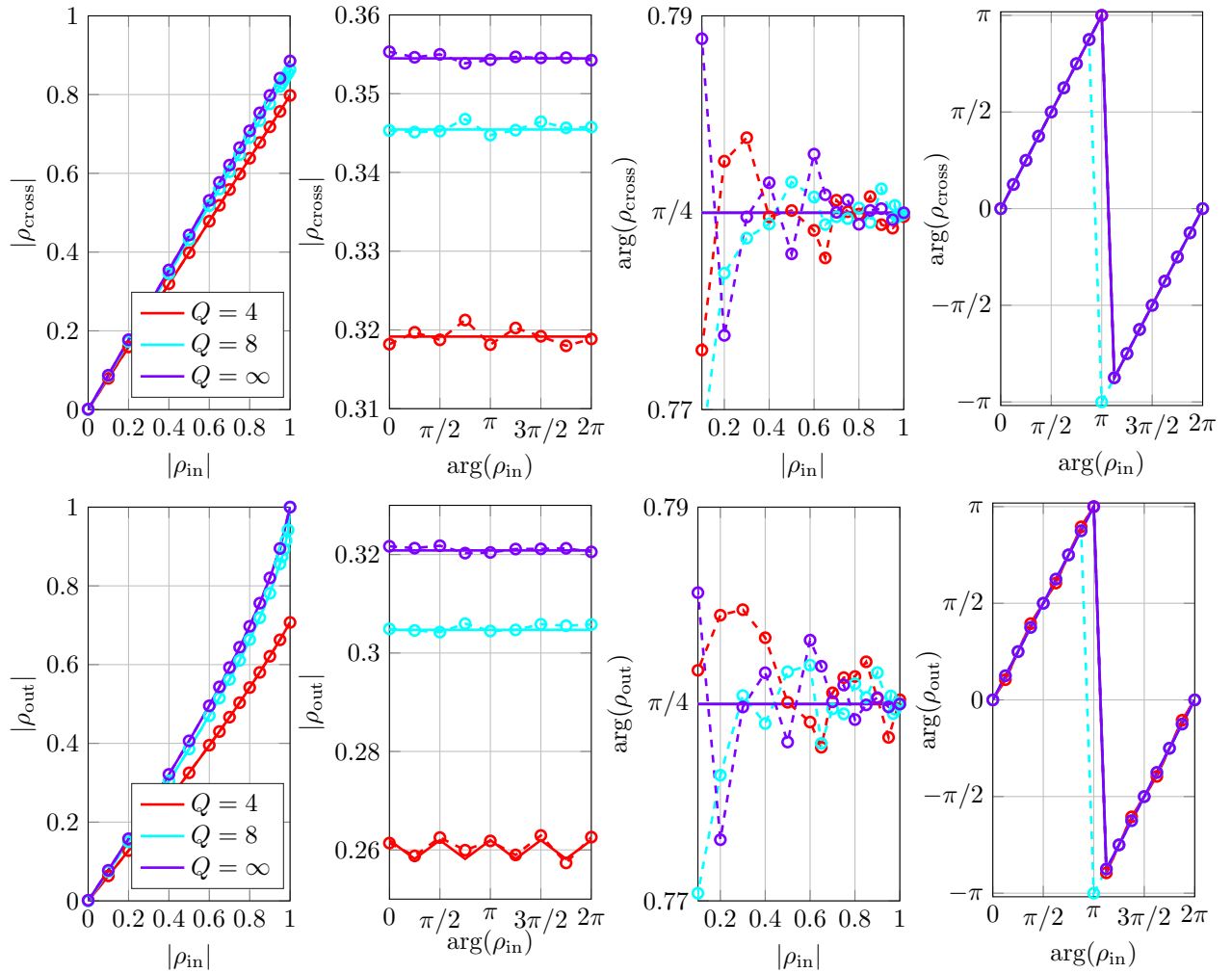


Fig. 3.2: Comparison of the derived closed-form expressions in (3.37) and (3.50) (solid lines) with the numerical results obtained by Monte-Carlo simulations (dashed lines with markers). In each figure, one dimension of ρ_{in} is fixed either $|\rho_{\text{in}}| = 0.4$ or $\arg(\rho_{\text{in}}) = \frac{\pi}{4}$, © 2018 IEEE.

read as

$$\mathbf{C}_{\mathbf{t}\mathbf{x}} = \frac{Q}{2\sqrt{\pi}} \sin(\psi) \Xi \operatorname{diag}(\mathbf{C}_{\mathbf{x}\mathbf{x}})^{-1/2} \mathbf{C}_{\mathbf{x}\mathbf{x}}, \quad (3.51)$$

and

$$\mathbf{C}_{\mathbf{t}\mathbf{t}} = \frac{Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \Xi \arcsin \left(\Re\{\mathbf{R}_{\mathbf{x}\mathbf{x}} e^{-j2\Delta k\psi}\} \right) \Xi. \quad (3.52)$$

3.4 Optimal CEQ

3.4.1 One-Dimensional Optimal CEQ

In this section, the magnitude Ξ is optimized to minimize the quantization distortions introduced to the signal. In general, the one-dimensional CEQ output can be expressed as

$$t = x + \eta_q, \quad (3.53)$$

where η_q denotes the complex-valued quantization distortion term. Since we assume that the input signal $x = r e^{j\phi}$ is a complex-valued Gaussian distributed signal with zero mean and variance σ_x^2 , we get

$$\mathbb{E}[\eta_q] = \mathbb{E}[t] - \mathbb{E}[x] = 0, \quad (3.54)$$

and

$$\mathbb{E}[|\eta_q|_2^2] = \sigma_{\eta_q}^2 = \beta_q \sigma_x^2, \quad (3.55)$$

where β_q denotes the distortion factor of the CEQ. The task is to find the optimal Ξ that minimizes the power of the quantization distortion term; that is

$$\Xi_{\text{opt}} = \arg \min_{\Xi} \sigma_{\eta_q}^2. \quad (3.56)$$

To this end, we find the expression of $\sigma_{\eta_q}^2$.

$$\begin{aligned}
\sigma_{\eta_q}^2 &= \text{E} \left[|t - x|_2^2 \right] \\
&= \int_0^{2\pi} \int_0^\infty |t - x|_2^2 p_X(r, \phi) r \, dr \, d\phi \\
&= \sum_{k=1}^Q \int_{(2k-2)\psi}^{2k\psi} \int_0^\infty \left| \Xi e^{j(2k-1)\psi} - r e^{j\phi} \right|_2^2 \frac{1}{\pi\sigma_x^2} e^{-\frac{r^2}{\sigma_x^2}} r \, dr \, d\phi \\
&= \frac{1}{\pi\sigma_x^2} \sum_{k=1}^Q \int_{(2k-2)\psi}^{2k\psi} \int_0^\infty \left(\Xi^2 + r^2 - \Xi r e^{j((2k-1)\psi - \phi)} - \Xi r e^{-j((2k-1)\psi - \phi)} \right) e^{-\frac{r^2}{\sigma_x^2}} r \, dr \, d\phi \\
&= \frac{1}{\pi\sigma_x^2} \left(2\pi\Xi^2 \int_0^\infty e^{-\frac{r^2}{\sigma_x^2}} r \, dr + 2\pi \int_0^\infty e^{-\frac{r^2}{\sigma_x^2}} r^3 \, dr \right. \\
&\quad - \Xi \sum_{k=1}^Q e^{j((2k-1)\psi)} \int_{(2k-2)\psi}^{2k\psi} e^{-j\phi} \int_0^\infty e^{-\frac{r^2}{\sigma_x^2}} r^2 \, dr \, d\phi \\
&\quad \left. - \Xi \sum_{k=1}^Q e^{-j((2k-1)\psi)} \int_{(2k-2)\psi}^{2k\psi} e^{j\phi} \int_0^\infty e^{-\frac{r^2}{\sigma_x^2}} r^2 \, dr \, d\phi \right) \\
&= \frac{1}{\pi\sigma_x^2} \left(2\pi\Xi^2 \frac{\sigma_x^2}{2} + 2\pi \frac{\sigma_x^4}{2} \right. \\
&\quad \left. - \Xi \frac{\sigma_x^2}{4} \left(\sum_{k=1}^Q e^{j((2k-1)\psi)} \int_{(2k-2)\psi}^{2k\psi} e^{-j\phi} \sqrt{\pi\sigma_x^2} \, d\phi + \sum_{k=1}^Q e^{-j((2k-1)\psi)} \int_{(2k-2)\psi}^{2k\psi} e^{j\phi} \sqrt{\pi\sigma_x^2} \, d\phi \right) \right) \\
&= \Xi^2 + \sigma_x^2 - Q \frac{\sigma_x}{4\sqrt{\pi}} \Xi (2j(e^{-j\psi} - e^{j\psi})) \\
&= \Xi^2 + \sigma_x^2 - Q \frac{\sigma_x}{\sqrt{\pi}} \Xi \sin(\psi). \tag{3.57}
\end{aligned}$$

The optimal magnitude Ξ_{opt} is obtained by setting the derivative of $\sigma_{\eta_q}^2$ w.r.t. Ξ equal to 0; that is

$$\frac{d\sigma_{\eta_q}^2}{d\Xi} = 2\Xi - \frac{\sigma_x}{\sqrt{\pi}} Q \sin(\psi) = 0. \tag{3.58}$$

Thus, the optimal magnitude is given by

$$\Xi_{\text{opt}} = \frac{Q}{2\sqrt{\pi}} \sin(\psi) \sigma_x. \tag{3.59}$$

The optimal magnitude of the CE signal is dependent on the standard deviation of the input signal x . Thus, a different optimal magnitude is obtained for a different input signal variance σ_x^2 . With the use of the optimal CEQ, the following covariance factors simplify to:

$$\begin{aligned}
\text{E} [t\eta_q^*] &= \text{E} [t(t^* - x^*)] \\
&= \text{E} [tt^*] - \text{E} [tx^*] \\
&= \Xi_{\text{opt}}^2 - \frac{Q}{2\sqrt{\pi}} \sin(\psi) \Xi_{\text{opt}} \sigma_x \stackrel{(3.59)}{=} 0, \tag{3.60}
\end{aligned}$$

q	Q	Ξ_{opt}	β_q
2	4	$\sqrt{2/\pi}$	$1 - 2/\pi$
3	8	0.86362402	0.2541536
4	16	0.88054342	0.2246433
5	32	0.88480399	0.2171219
∞	∞	$\sqrt{\pi/4} = 0.88622693$	$1 - \pi/4$

Table 3.1: Optimal step size for the CEQ and the corresponding distortion factor for unit variance inputs.

and

$$\begin{aligned}
 \mathbb{E}[x\eta_q^*] &= \mathbb{E}[(t - \eta_q)\eta_q^*] \\
 &= \mathbb{E}[t\eta_q^*] - \sigma_{\eta_q}^2 \\
 &\stackrel{(3.60)}{=} -\sigma_{\eta_q}^2 \\
 &\stackrel{(3.55)}{=} -\beta_q\sigma_x^2.
 \end{aligned} \tag{3.61}$$

After plugging (3.59) in (3.57), the distortion factor β_q for the optimal CEQ reduces to

$$\beta_q = 1 - \frac{Q^2 \sin^2(\psi)}{4\pi}. \tag{3.62}$$

And we introduce the new variable α_q as

$$\alpha_q = 1 - \beta_q = \frac{Q^2 \sin^2(\psi)}{4\pi}. \tag{3.63}$$

The optimal values Xi_{opt} and the resulting distortion factors of the CEQ having unit variance Gaussian complex-valued inputs, i.e. $\sigma_x^2 = 1$, are summarized in Table 3.1 for different quantization resolutions. We notice that the values of Ξ_{opt} and β_q , when rounded up to the second digit after the decimal point, do not vary for $q \geq 4$ and hence for $Q \geq 16$. This explains why we observed the same behavior between $Q = 16$ and $Q = \infty$ in Fig. 3.2.

3.4.2 Multi-Dimensional Optimal CEQ

From (3.59) and the definition in (3.63), the optimal matrix Ξ for an N -dimensional CEQ with input vector \mathbf{x} and output vector \mathbf{t} is given by

$$\Xi_{\text{opt}} = \sqrt{\alpha_q} \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2}. \tag{3.64}$$

3.5 Statistical Properties of the Optimal CEQ

The statistical properties of the optimal CEQ can be first obtained from Price's theorem introduced in Section 3.3 by just plugging (3.64) in (3.51) and (3.52). A second way to obtain the approximate statistical properties of the CEQ is the LCA that was applied in [25] and will be recalled below. One might ask why the LCA is only considered in the case of the optimal CEQ. The answer to this question is that the covariance factor between the output

signal t and the quantization distortion term η_q vanishes under the optimality condition of the CEQ as shown in (3.60). This fact simplifies the derivations and allows us to get closed-form approximations of the statistical properties of the CEQ.

To this end, we again assume an N -dimensional CEQ with joint Gaussian distributed input signals x_n , $n = 1, \dots, N$.

3.5.1 Price's Theorem

When plugging (3.64) in (3.51) and (3.52), we get

$$\mathbf{C}_{\mathbf{tx}} = \alpha_q \mathbf{C}_{\mathbf{xx}}, \quad (3.65)$$

and

$$\mathbf{C}_{\mathbf{tt}} = \frac{Q}{\pi} \sin^2(\psi) \alpha_q \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2} \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \arcsin \left(\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\} \right) \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2}. \quad (3.66)$$

3.5.2 LCA

To find the corresponding expressions of the covariance matrices when using the LCA, we write first the general expressions

$$\mathbf{C}_{\mathbf{tx}} = \mathbf{C}_{\mathbf{xx}} + \mathbf{C}_{\eta_q \mathbf{x}}, \quad (3.67)$$

and

$$\mathbf{C}_{\mathbf{tt}} = \mathbf{C}_{\mathbf{xx}} + \mathbf{C}_{\mathbf{x}\eta_q} + \mathbf{C}_{\eta_q \mathbf{x}} + \mathbf{C}_{\eta_q \eta_q}. \quad (3.68)$$

To this end, we have to find the entries of $\mathbf{C}_{\mathbf{x}\eta_q}$ and $\mathbf{C}_{\eta_q \eta_q}$. First, the correlation factor between the input signal x_i and the quantization distortion term η_{q_j} reads as

$$\begin{aligned} \mathbb{E} \left[x_i \eta_{q_j}^* \right] &= \mathbb{E}_{x_j} \left[\mathbb{E} \left[x_i \eta_{q_j}^* | x_j \right] \right] \\ &\stackrel{(a)}{=} \mathbb{E}_{x_j} \left[\mathbb{E} \left[x_i | x_j \right] \mathbb{E} \left[\eta_{q_j}^* | x_j \right] \right] \\ &\stackrel{(b)}{=} \mathbb{E}_{x_j} \left[\mathbb{E} \left[x_i x_j^* \right] \mathbb{E} \left[x_j x_j^* \right]^{-1} x_j \mathbb{E} \left[\eta_{q_j}^* | x_j \right] \right] \\ &= \mathbb{E} \left[x_i x_j^* \right] \mathbb{E} \left[x_j x_j^* \right]^{-1} \mathbb{E} \left[x_j \eta_{q_j}^* \right] \\ &\stackrel{(c)}{=} -\beta_q \mathbb{E} \left[x_i x_j^* \right], \end{aligned} \quad (3.69)$$

where in (a) we use the fact that the quantization distortion term η_{q_j} does not statistically depend on the other random variables when conditioned on x_j , in (b) the Bayesian estimator is equal to the linear estimator for jointly Gaussian distributed signals x_i and x_j and (c) results from (3.61). In summary, it holds that

$$\mathbf{C}_{\mathbf{x}\eta_q} = -\beta_q \mathbf{C}_{\mathbf{xx}}. \quad (3.70)$$

Second, the correlation factor between two quantization distortion terms η_i and η_j , with $i \neq j$, reads as

$$\begin{aligned}
 \mathbb{E} \left[\eta_{q_i} \eta_{q_j}^* \right] &= \mathbb{E}_{x_j} \left[\mathbb{E} \left[\eta_{q_i} \eta_{q_j}^* | x_j \right] \right] \\
 &= \mathbb{E}_{x_j} \left[\mathbb{E} \left[\eta_{q_i} | x_j \right] \mathbb{E} \left[\eta_{q_j}^* | x_j \right] \right] \\
 &\stackrel{(d)}{\approx} \mathbb{E}_{x_j} \left[\mathbb{E} \left[\eta_{q_i} x_j^* \right] \mathbb{E} \left[x_j x_j^* \right]^{-1} x_j \mathbb{E} \left[\eta_{q_j}^* | x_j \right] \right] \\
 &= \mathbb{E} \left[\eta_{q_i} x_j^* \right] \mathbb{E} \left[x_j x_j^* \right]^{-1} \mathbb{E} \left[x_j \eta_{q_j}^* \right] \\
 &= -\beta_q \mathbb{E} \left[\eta_{q_i} x_j^* \right] \\
 &\stackrel{(3.61)}{=} -\beta_q \left(-\beta_q \mathbb{E} \left[x_i x_j^* \right] \right) \\
 &= \beta_q^2 \mathbb{E} \left[x_i x_j^* \right], \tag{3.71}
 \end{aligned}$$

where in (d) the Bayesian estimator is approximated with the linear estimator by assuming that η_{q_i} is Gaussian distributed. The power of the quantization distortion term is related to the input signal variance as follows

$$\mathbb{E} \left[\eta_{q_i} \eta_{q_i}^* \right] = \sigma_{\eta_{q_i}}^2 = \beta_q \sigma_{x_i}^2. \tag{3.72}$$

Thus, we obtain

$$\mathbf{C}_{\eta_q \eta_q} = \beta_q^2 \mathbf{C}_{\mathbf{xx}} + \alpha_q \beta_q \text{diag} \left(\mathbf{C}_{\mathbf{xx}} \right). \tag{3.73}$$

Plugging (3.70) and (3.73) in (3.67) and (3.68), we get

$$\mathbf{C}_{\mathbf{tx}} = \alpha_q \mathbf{C}_{\mathbf{xx}}, \tag{3.74}$$

and

$$\begin{aligned}
 \mathbf{C}_{\mathbf{tt}} &= \alpha_q^2 \mathbf{C}_{\mathbf{xx}} + \alpha_q \beta_q \text{diag} \left(\mathbf{C}_{\mathbf{xx}} \right) \\
 &= \alpha_q \text{diag} \left(\mathbf{C}_{\mathbf{xx}} \right) + \alpha_q^2 \text{nondiag} \left(\mathbf{C}_{\mathbf{xx}} \right). \tag{3.75}
 \end{aligned}$$

3.5.3 Price's Theorem vs. LCA

When we compare the expressions obtained by Price's Theorem with the expressions obtained by the LCA, we observe that the only difference is in the computation of $\mathbf{C}_{\mathbf{tt}}$. To discover the relationship between both methods, we first split the expression in (3.66) into diagonal

and non-diagonal parts

$$\begin{aligned}
\mathbf{C}_{\mathbf{tt}} &= \frac{Q}{\pi} \sin^2(\psi) \alpha_q \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2} \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \arcsin \left(\Re\{\text{diag}(\mathbf{R}_{\mathbf{xx}}) e^{-j2\Delta k\psi}\} \right) \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2} \\
&+ \frac{Q}{\pi} \sin^2(\psi) \alpha_q \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2} \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \arcsin \left(\Re\{\text{nondiag}(\mathbf{R}_{\mathbf{xx}}) e^{-j2\Delta k\psi}\} \right) \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2} \\
&= \frac{Q}{\pi} \sin^2(\psi) \alpha_q \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2} \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \arcsin(\cos(2\Delta k\psi) \mathbf{I}_N) \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2} \\
&+ \frac{Q}{\pi} \sin^2(\psi) \alpha_q \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2} \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \arcsin \left(\Re\{\text{nondiag}(\mathbf{R}_{\mathbf{xx}}) e^{-j2\Delta k\psi}\} \right) \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2} \\
&= \frac{Q}{\pi} \sin^2(\psi) \alpha_q \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2} \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \left(\frac{\pi}{2} - 2\Delta k\psi \right) \mathbf{I}_N \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2} \\
&+ \frac{Q}{\pi} \sin^2(\psi) \alpha_q \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2} \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \arcsin \left(\Re\{\text{nondiag}(\mathbf{R}_{\mathbf{xx}}) e^{-j2\Delta k\psi}\} \right) \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2} \\
&\stackrel{(e)}{=} \alpha_q \text{diag}(\mathbf{C}_{\mathbf{xx}}) \\
&+ \frac{Q}{\pi} \sin^2(\psi) \alpha_q \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2} \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \arcsin \left(\Re\{\text{nondiag}(\mathbf{R}_{\mathbf{xx}}) e^{-j2\Delta k\psi}\} \right) \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2},
\end{aligned} \tag{3.76}$$

where in (e) we made use of the property

$$\frac{Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \left(\frac{\pi}{2} - 2\Delta k\psi \right) = 1. \tag{3.77}$$

Second we approximate the $\arcsin(\bullet)$ function by its first order Taylor expansion and obtain

$$\begin{aligned}
\mathbf{C}_{\mathbf{tt}} &\approx \alpha_q \text{diag}(\mathbf{C}_{\mathbf{xx}}) + \frac{Q}{\pi} \sin^2(\psi) \alpha_q \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \Re\{\text{nondiag}(\mathbf{C}_{\mathbf{xx}}) e^{-j2\Delta k\psi}\} \\
&\stackrel{(g)}{=} \alpha_q \text{diag}(\mathbf{C}_{\mathbf{xx}}) + \alpha_q^2 \text{nondiag}(\mathbf{C}_{\mathbf{xx}}),
\end{aligned} \tag{3.78}$$

where in (g) it holds that

$$\sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \Re\{x e^{-j2\Delta k\psi}\} = \frac{Q}{4} x. \tag{3.79}$$

Hence, we get the same expression as in (3.75). This implies that the approximation of the non-diagonal entries in $\mathbf{C}_{\mathbf{tt}}$ from (3.66) based on the first order Taylor expansion of $\arcsin(\bullet)$ leads to the same expression as in (3.75) for the LCA. In other words, the LCA represents the first order Taylor expansion of Price's Theorem.

3.6 Bussgang Linearization of the CEQ

Theorem 2 (Bussgang's Theorem [61]). *For two Gaussian distributed signals, the crosscorrelation function taken after one of them has undergone non-linear amplitude distortion is identical, except for a factor of proportionality, to the crosscorrelation function taken before the distortion.*

According to Bussgang's Theorem [61], it holds that

$$\mathbf{C}_{\mathbf{t}\mathbf{x}} = \mathbf{L}_Q \mathbf{C}_{\mathbf{x}\mathbf{x}}. \quad (3.80)$$

Consequently, every non-linear function with Gaussian distributed input signal can be modeled by the sum of a linear function and a distortion term that is uncorrelated with the input as

$$\mathbf{t} = \mathbf{L}_Q \mathbf{x} + \mathbf{d}_Q, \quad (3.81)$$

such that

$$\mathbb{E} [\mathbf{x} \mathbf{d}_Q^H] = \mathbf{0}. \quad (3.82)$$

This leads us to the result

$$\begin{aligned} \mathbf{L}_Q &= \mathbf{C}_{\mathbf{t}\mathbf{x}} \mathbf{C}_{\mathbf{x}\mathbf{x}}^{-1} \\ &= \sqrt{\alpha_q} \mathbf{\Xi} \text{diag}(\mathbf{C}_{\mathbf{x}\mathbf{x}})^{-\frac{1}{2}} \mathbf{C}_{\mathbf{x}\mathbf{x}} \mathbf{C}_{\mathbf{x}\mathbf{x}}^{-1} \\ &= \sqrt{\alpha_q} \mathbf{\Xi} \text{diag}(\mathbf{C}_{\mathbf{x}\mathbf{x}})^{-\frac{1}{2}}. \end{aligned} \quad (3.83)$$

For the optimal CEQ, (3.64) can be plugged in (3.83) and the expression in (3.83) reduces to

$$\mathbf{L}_Q = \alpha_q \mathbf{I}_N. \quad (3.84)$$

Indeed, when we plug (3.83) in (3.80) and compare the obtained expression with (3.51), we notice that Bussgang's theorem is a special case of Price' Theorem. Note that \mathbf{d}_Q is not Gaussian distributed but has an unknown distribution. However, the covariance matrix of \mathbf{d}_Q can be computed whether by applying Price's theorem or the LCA. First, with Price's theorem, we get

$$\begin{aligned} \mathbf{C}_{\mathbf{d}_Q \mathbf{d}_Q} &= \mathbb{E} \left[(\mathbf{t} - \mathbf{L}_Q \mathbf{x}) (\mathbf{t} - \mathbf{L}_Q \mathbf{x})^H \right] \\ &= \mathbf{C}_{\mathbf{t}\mathbf{t}} - \mathbf{L}_Q \mathbf{C}_{\mathbf{x}\mathbf{t}} - \mathbf{C}_{\mathbf{t}\mathbf{x}} \mathbf{L}_Q^H + \mathbf{L}_Q \mathbf{C}_{\mathbf{x}\mathbf{x}} \mathbf{L}_Q^H \\ &= \mathbf{C}_{\mathbf{t}\mathbf{t}} - \mathbf{L}_Q \mathbf{C}_{\mathbf{x}\mathbf{x}} \mathbf{L}_Q^H \\ &= \mathbf{C}_{\mathbf{t}\mathbf{t}} - \alpha_q \mathbf{\Xi} \mathbf{R}_{\mathbf{x}\mathbf{x}} \mathbf{\Xi} \\ &= \frac{Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \mathbf{\Xi} \arcsin \left(\Re \{ \mathbf{R}_{\mathbf{x}\mathbf{x}} e^{-j(2\Delta k\psi)} \} \right) \mathbf{\Xi} - \alpha_q \mathbf{\Xi} \mathbf{R}_{\mathbf{x}\mathbf{x}} \mathbf{\Xi} \\ &= \beta_q \mathbf{\Xi}^2 + \frac{Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \mathbf{\Xi} \arcsin \left(\Re \{ \text{nondiag}(\mathbf{R}_{\mathbf{x}\mathbf{x}}) e^{-j(2\Delta k\psi)} \} \right) \mathbf{\Xi} \\ &\quad - \alpha_q \mathbf{\Xi} \text{nondiag}(\mathbf{R}_{\mathbf{x}\mathbf{x}}) \mathbf{\Xi} \\ &\approx \beta_q \mathbf{\Xi}^2, \end{aligned} \quad (3.85)$$

where we use the first order Taylor expansion of the arcsin function. For the optimal CEQ, we get

$$\mathbf{C}_{\mathbf{d}_Q \mathbf{d}_Q} \approx \alpha_q \beta_q \text{diag}(\mathbf{C}_{\mathbf{xx}}), \quad (3.86)$$

which matches with the result when using the LCA.

In the following, we aim at getting the covariance matrix for different time instants. When applying Price's theorem, we get

$$\begin{aligned} \mathbb{E}[\mathbf{d}_Q[t-t_1] \mathbf{d}_Q^H[t-t_2]] &= \mathbb{E}\left[(\mathbf{t}[t-t_1] - \mathbf{L}_Q \mathbf{x}[t-t_1]) (\mathbf{t}[t-t_2] - \mathbf{L}_Q \mathbf{x}[t-t_2])^H\right] \\ &= \mathbb{E}[\mathbf{t}[t-t_1] \mathbf{t}^H[t-t_2]] - \mathbb{E}[\mathbf{t}[t-t_1] \mathbf{x}^H[t-t_2]] \mathbf{L}_Q^H \\ &\quad - \mathbf{L}_Q \mathbb{E}[\mathbf{x}[t-t_1] \mathbf{t}^H[t-t_2]] + \mathbf{L}_Q \mathbb{E}[\mathbf{x}[t-t_1] \mathbf{x}^H[t-t_2]] \mathbf{L}_Q^H \\ &\stackrel{(3.80)}{=} \mathbb{E}[\mathbf{t}[t-t_1] \mathbf{t}^H[t-t_2]] - \mathbf{L}_Q \mathbb{E}[\mathbf{x}[t-t_1] \mathbf{x}^H[t-t_2]] \mathbf{L}_Q^H \\ &\quad - \mathbf{L}_Q \mathbb{E}[\mathbf{x}[t-t_1] \mathbf{x}^H[t-t_2]] \mathbf{L}_Q^H + \mathbf{L}_Q \mathbb{E}[\mathbf{x}[t-t_1] \mathbf{x}^H[t-t_2]] \mathbf{L}_Q^H \\ &= \mathbb{E}[\mathbf{t}[t-t_1] \mathbf{t}^H[t-t_2]] - \mathbf{L}_Q \mathbb{E}[\mathbf{x}[t-t_1] \mathbf{x}^H[t-t_2]] \mathbf{L}_Q^H \\ &= \frac{Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} e^{j2\Delta k\psi} \Xi \\ &\quad \arcsin(\Re\{\text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \mathbb{E}[\mathbf{x}[t-t_1] \mathbf{x}^H[t-t_2]] \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} e^{-j2\Delta k\psi}\}) \\ &\Xi - \alpha_q \Xi \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \mathbb{E}[\mathbf{x}[t-t_1] \mathbf{x}^H[t-t_2]] \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \Xi. \end{aligned} \quad (3.87)$$

Again, with the first order Taylor expansion of the arcsin function, we get

$$\mathbb{E}[\mathbf{d}_Q[t-t_1] \mathbf{d}_Q^H[t-t_2]] \approx \begin{cases} \beta_q \Xi^2 & \text{if } t_1 = t_2, \\ \mathbf{0}_{N,N} & \text{otherwise.} \end{cases} \quad (3.88)$$

For the optimal CEQ, it holds that

$$\mathbb{E}[\mathbf{d}_Q[t-t_1] \mathbf{d}_Q^H[t-t_2]] \approx \begin{cases} \alpha_q \beta_q \text{diag}(\mathbf{C}_{\mathbf{xx}}) & \text{if } t_1 = t_2, \\ \mathbf{0}_{N,N} & \text{otherwise,} \end{cases} \quad (3.89)$$

where this approximation can be also obtained by the LCA.

4. System Model

4.1 Input-Output Relationship

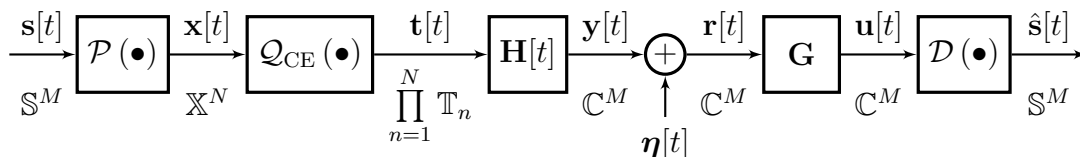


Fig. 4.1: Downlink MU-MIMO system model.

The system model shown in Fig.4.1 consists of a single-cell massive MU-MIMO downlink scenario with coarsely QCE signals at the transmitter. The BS is equipped with N antennas and serves M single-antenna users simultaneously, where $N \gg M$.

The input signal $\mathbf{s}[t] \in \mathbb{S}^M$ contains the symbols to be transmitted to each of the M users at time instant t . Each user's signal is drawn from the set \mathbb{S} that represents either an S -PSK or S -QAM constellation, where S denotes the number of constellation points. The input set is detailed Section 4.3. We assume that $\mathbb{E}[\mathbf{s}[t]] = \mathbf{0}_M$, $\mathbb{E}[\mathbf{s}[t]\mathbf{s}[t]^H] = \sigma_s^2 \mathbf{I}_M$, and $\mathbb{E}[\mathbf{s}[t]\mathbf{s}[t-\tau]^H] = \mathbf{0}_M$, $\forall t$ and $\tau \neq t$.

The signal vector $\mathbf{s}[t]$ is precoded into the vector $\mathbf{x}[t] \in \mathbb{X}^N$ prior to the polar DACs. The set \mathbb{X} can represent the set of complex numbers, i.e. \mathbb{C} or a subset of it depending on the precoder design. The choice of \mathbb{X} will be discussed in the following sections. The function $\mathcal{P}(\bullet)$ represents the precoding operation to reduce the distortions caused by the coarse quantization and the MUI. A brief overview about the precoding function $\mathcal{P}(\bullet)$ is given in Section 4.4. The whole thesis is devoted to explain in detail the precoder design.

The operator $\mathcal{Q}_{\text{CE}}(\bullet)$, defined in Chapter 3, models the non-linear behavior of the polar DACs combined with the power allocation at the PAs as

$$\mathbf{t}[t] = \mathcal{Q}_{\text{CE}}(\mathbf{x}[t]). \quad (4.1)$$

After the CEQ each entry of the transmit vector t_n , $n = 1, \dots, N$, belongs to the set \mathbb{T}_n , which is defined as

$$\mathbb{T}_n = \left\{ \Xi_n \exp \left(j(2i-1) \frac{\pi}{Q} \right) : i = 1, \dots, Q \right\}, \quad (4.2)$$

where Ξ_n^2 denotes the coefficient for the power allocation at the n -th antenna that is chosen according to the transmit power constraint detailed in Section 4.5.

The signal $\mathbf{t}[t]$ is transmitted through the channel modeled by

$$\mathbf{H}[t] = \sum_{\ell=0}^{L-1} \mathbf{H}_\ell \delta[t - \ell], \quad (4.3)$$

where δ represents the Dirac function and with the (m, n) -th element of the ℓ -th channel matrix, i.e. $[\mathbf{H}_\ell]_{m,n}$, being the ℓ -th channel tap among L taps between the n -th transmit antenna and the m -th user. At the M receive antennas, an AWGN, which is denoted by the vector $\boldsymbol{\eta} \sim \mathcal{CN}_{\mathbb{C}}(\mathbf{0}_M, \mathbf{C}_\boldsymbol{\eta} = \mathbf{I}_M)$, perturbs the received signals

$$\mathbf{r}[t] = \mathbf{H}[t] * \mathbf{t}[t] + \boldsymbol{\eta}[t]. \quad (4.4)$$

The precoder is designed such that, without any noise, the received signals would lie in their intended decision regions and no joint receive processing is necessary. Additionally, coherent data transmission with multiple BS antennas leads to an antenna gain, which depends on the channel realization. The entries of the received signal vector \mathbf{r} do not belong to the nominal decision regions of \mathbb{S} but to a scaled version of them. Therefore, rescaling the received signal at each receive antenna is required to make the signal belong to the nominal decision region. The rescaling operation is modeled by the diagonal real-valued matrix \mathbf{G} , as follows

$$\mathbf{u}[t] = \mathbf{G} (\mathbf{H}[t] * \mathbf{t}[t] + \boldsymbol{\eta}[t]), \quad (4.5)$$

where

$$\mathbf{G} = \sum_{m=1}^M g_m \mathbf{e}_m \mathbf{e}_m^T, \quad (4.6)$$

with $g_m > 0$, $m = 0, \dots, M$. The coefficients g_m are blindly estimated at the receiver over a block of T received symbols as explained in Section 4.7. Note that no receive processing \mathbf{G} is required if \mathbb{S} represents the PSK constellation. Finally, based on the decision regions to which the entries of the signal \mathbf{u} belong, the decision operation $\mathcal{D}(\bullet)$ produces the detected symbols $\hat{\mathbf{s}}$ at the users

$$\hat{\mathbf{s}}[t] = \mathcal{D}(\mathbf{G} (\mathbf{H}[t] * \mathbf{t}[t] + \boldsymbol{\eta}[t])). \quad (4.7)$$

4.2 Compact Input-Output Relationship

We drop out the time index and express $\mathbf{u} = \mathbf{u}[t]$ compactly as

$$\mathbf{u} = \mathbf{G} (\mathbf{H}\mathbf{t} + \boldsymbol{\eta}), \quad (4.8)$$

where

$$\mathbf{H} = [\mathbf{H}_0 \quad \mathbf{H}_1 \quad \dots \quad \mathbf{H}_{L-1}], \quad (4.9)$$

$$\mathbf{t} = [\mathbf{t}[t]^T \quad \mathbf{t}[t-1]^T \quad \dots \quad \mathbf{t}[t-L+1]^T]^T = \mathcal{Q}_{\text{CE}}(\mathbf{x}), \quad (4.10)$$

$$\mathbf{x} = [\mathbf{x}[t]^T \quad \mathbf{x}[t-1]^T \quad \dots \quad \mathbf{x}[t-L+1]^T]^T, \quad (4.11)$$

and

$$\boldsymbol{\eta} = \boldsymbol{\eta}[t]. \quad (4.12)$$

4.3 Input Set

Within the scope of the thesis, the input signals can be drawn from an S -PSK constellation or an S -QAM constellation. To this end, we define both constellations. First, the S -PSK constellation is given by

$$\mathbb{S} := \{\exp(j(2i-1)\theta) : i = 1, \dots, S\}, \text{ where } \theta = \frac{\pi}{S}. \quad (4.13)$$

Second, the S -QAM constellation, where S is assumed to be a power of 4, is defined as

$$\mathbb{S} := \{\pm(2i-1) \pm j(2i-1) : i = 1, \dots, \log_4(S)\}. \quad (4.14)$$

4.4 Transmit Processing: Precoding

4.4.1 Linear Precoding

For linear precoding techniques, the precoding function reduces to a Finite Impulse Response (FIR) matrix; that is

$$\mathbf{P}[t] = \sum_{\ell'=0}^{L_p-1} \mathbf{P}_{\ell'} \delta[t - \ell'], \quad (4.15)$$

where L_p represents the number of taps of the precoding FIR filter between every antenna element and every user.

4.4.2 Non-linear Precoding

In the case of non-linear precoding, no explicit expression for $\mathcal{P}(\bullet)$ is provided. However, for a given channel realization \mathbf{H} , every input vector \mathbf{s} is mapped to a precoded vector \mathbf{x} .

$$\mathbf{x} = \mathcal{P}(\mathbf{s}, \mathbf{H}). \quad (4.16)$$

4.5 Transmit Power Constraint

For an available transmit power P_{tx} , the transmit power constraint is given by

$$\text{tr}(\mathbf{C}_{\text{tt}}) \leq LP_{\text{tx}}. \quad (4.17)$$

Therefore, we should compute $\text{tr}(\mathbf{C}_{\text{tt}})$. Applying (3.52), we obtain

$$\begin{aligned}
\text{tr}(\mathbf{C}_{\text{tt}}) &= \frac{Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \text{tr}(\mathbf{\Xi} \arcsin(\Re\{\mathbf{R}_{\text{xx}} e^{-j2\Delta k\psi}\}) \mathbf{\Xi}) \\
&= \frac{Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \text{tr}(\text{diag}(\arcsin(\Re\{\mathbf{R}_{\text{xx}} e^{-j2\Delta k\psi}\})) \mathbf{\Xi}^2) \\
&= \frac{Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \text{tr}(\arcsin(\cos(2\Delta k\psi)) \mathbf{I}_N \mathbf{\Xi}^2) \\
&= \frac{Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \text{tr}\left(\left(\frac{\pi}{2} - 2\Delta k\psi\right) \mathbf{I}_N \mathbf{\Xi}^2\right) \\
&= \frac{Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} \left(\frac{\pi}{2} - 2\Delta k\psi\right) e^{j(2\Delta k\psi)} \text{tr}(\mathbf{\Xi}^2) \\
&\stackrel{(3.77)}{=} \text{tr}(\mathbf{\Xi}^2).
\end{aligned} \tag{4.18}$$

Hence, it must hold that

$$\text{tr}(\mathbf{\Xi}^2) \leq LP_{\text{tx}}. \tag{4.19}$$

We choose $\mathbf{\Xi} = \sqrt{\alpha_q} \text{diag}(\mathbf{C}_{\text{xx}})^{1/2}$ to reduce the quantization distortions, as stated in (3.64). Consequently, the transmit power constraint can be rewritten as

$$\alpha_q \text{tr}(\mathbf{C}_{\text{xx}}) \leq LP_{\text{tx}}. \tag{4.20}$$

4.5.1 Equal Power Allocation

In the case of equal power allocation, the matrix $\mathbf{\Xi}$ is just a scaled identity, i.e. $\mathbf{\Xi} = \Xi \mathbf{I}_{NL}$. This leads to

$$\text{diag}(\mathbf{C}_{\text{xx}}) = \frac{\Xi^2}{\alpha_q} \mathbf{I}_{NL} \tag{4.21}$$

due to the optimality condition of the CEQ in (3.64). The value of Ξ is found by fulfilling the power constraint in (4.19)

$$\Xi \leq \sqrt{\frac{P_{\text{tx}}}{N}}. \tag{4.22}$$

To exploit the total available transmit power we choose

$$\Xi = \sqrt{\frac{P_{\text{tx}}}{N}}. \tag{4.23}$$

Plugging (4.23) in (4.21), it follows that

$$\text{diag}(\mathbf{C}_{\text{xx}}) = \frac{P_{\text{tx}}}{\alpha_q N} \mathbf{I}_{NL}. \tag{4.24}$$

4.5.2 Unequal Power Allocation

In the case of unequal power allocation, i.e. $\Xi \neq \Xi \mathbf{I}_N$, and for maximal exploitation of the total available transmit power, it must hold that

$$\text{tr}(\mathbf{C}_{\mathbf{xx}}) = \frac{LP_{\text{tx}}}{\alpha_q}. \quad (4.25)$$

4.6 Channel Model

4.6.1 The i.i.d. Channel

The first channel model that we consider throughout this thesis is the simplistic i.i.d. channel, where the entries of the matrices \mathbf{H}_ℓ , $\ell = 0, \dots, L-1$, are of zero means and constant variances; that is

$$[\mathbf{H}_\ell]_{m,n} = \mathcal{CN}_{\mathbb{C}}\left(0, \sigma_h^{(\ell)2}\right), \quad m = 1, \dots, M, \quad n = 1, \dots, N, \quad (4.26)$$

where it must hold that

$$\sum_{\ell=0}^{L-1} \sigma_h^{(\ell)2} = 1. \quad (4.27)$$

In the case of $L = 1$, we obtain that $\sigma_h^2 = 1$. However, in the case of $L \neq 1$, the variances are defined by the exponential power delay profiles that are given in Table 4.1 and Table 4.2 for $L = 3$ and $L = 6$, respectively.

ℓ	1	2
$10 \log_{10} \left(\frac{\mathbb{E}[\ \mathbf{H}_\ell\ _{\mathbb{F}}^2]}{\mathbb{E}[\ \mathbf{H}_0\ _{\mathbb{F}}^2]} \right) / \text{dB}$	-3	-6

Table 4.1: Exponential power delay profile with $L = 3$.

4.6.2 The mmW Sparse Channel

A more realistic channel model that takes into account the use of mmW frequencies is the mmW sparse channel, which is a clustered channel based on the extended Saleh-Valenzuela model [62]. Experiments and measurements have shown that the mmW channel can be modeled by a number of clusters N_{cl} for each user that group rays of number N_{ray} departing from the BS in similar directions [63–66]. Mathematically, we can write that

$$\mathbf{e}_m^T \mathbf{H}_\ell = \sqrt{\frac{N}{N_{\text{cl}} N_{\text{ray}}}} \sum_{i=1}^{N_{\text{cl}}} \sum_{j=1}^{N_{\text{ray}}} \gamma_{i,j}^{(m,\ell)} \mathbf{a} \left(\phi_{i,j}^{(\ell)} \right)^H, \quad (4.28)$$

ℓ	1	2	3	4	5
$10 \log_{10} \left(\frac{\mathbb{E}[\ \mathbf{H}_\ell\ _{\mathbb{F}}^2]}{\mathbb{E}[\ \mathbf{H}_0\ _{\mathbb{F}}^2]} \right) / \text{dB}$	-1	-9	-10	-15	-20

Table 4.2: Vehicular A (Case II) power delay profile with $L = 6$.

where $\gamma_{i,j}^{(m,\ell)}$ denotes the complex small-scale fading gain on the j -th subpath of the i -th cluster at the ℓ -th main path for the m -th user and $\mathbf{a}(\phi)$ is the vector response function of the BS antenna arrays to the angular departures; that is, when we assume a uniform linear array, we get

$$\mathbf{a}(\phi) = \frac{1}{\sqrt{N}} [1 \quad e^{j\pi \sin(\phi)} \quad \dots \quad e^{j(N-1)\pi \sin(\phi)}]^T. \quad (4.29)$$

Each i -th cluster is defined by a certain mean Angle of Departure (AOD) that is drawn from a uniform distribution. The AODs that correspond to the rays within the same cluster, i.e. $\phi_{i,j}^{(\ell)}$, $j = 1, \dots, N_{\text{ray}}$, are drawn around the mean AOD from a truncated Laplacian PDF given by [67]

$$p_\phi(\phi) = \begin{cases} \frac{1}{\sqrt{2}\sigma_\phi(1-e^{-\sqrt{2}\pi/\sigma_\phi})} e^{-|\sqrt{2}\phi/\sigma_\phi|} & \text{if } \phi \in [-\pi, \pi), \\ 0 & \text{else.} \end{cases} \quad (4.30)$$

In this work, we assume that

$$\sigma_\phi = 1^\circ. \quad (4.31)$$

The small-scale fading factors fulfill the following property

$$\gamma_{i,j}^{(m,\ell)} \sim \mathcal{CN}_{\mathbb{C}}(0, \sigma_\gamma^{(\ell)2}), \quad (4.32)$$

where $\sigma_\gamma^{(\ell)2}$ are chosen according to the exponential power delay profile either in Table 4.1 or in Table 4.2. Additionally, it must hold that

$$\sum_{\ell=0}^{L-1} \sigma_\gamma^{(\ell)2} = 1. \quad (4.33)$$

4.7 Receive Processing

After multiplication with the receiver coefficient g_m , the m -th scaled received signal is

$$u_m[t] = g_m r_m[t] = g_m \mathbf{e}_m^T \mathbf{H}[t] * \mathbf{t}[t] + g_m \eta_m[t] = s_m[t - \tau] + \eta'_m[t], \quad (4.34)$$

where τ represents the time delay and $\eta'_m[t]$ denotes the deviation of $u_m[t]$ from the nominal point $s_m[t - \tau]$ due to the precoder design, the AWGN $\eta_m[t]$ and the quantization applied on the precoded vector \mathbf{x} . Then, we can write

$$|\Re\{r_m[t]\}| + |\Im\{r_m[t]\}| = g_m^{-1} (|\Re\{s_m[t - \tau] + \eta'_m[t]\}| + |\Im\{s_m[t - \tau] + \eta'_m[t]\}|) \quad (4.35)$$

$$\approx g_m^{-1} (|\Re\{s_m[t - \tau]\}| + |\Im\{s_m[t - \tau]\}| + \Re\{\eta'_m[t]\} + \Im\{\eta'_m[t]\}), \quad (4.36)$$

where the approximation is very accurate when the receiver Signal-to-Interference-Noise Ratio (SINR) is much larger than 1, which is the case for massive MIMO systems. Having zero-mean noise plus interference $\eta'_m[t]$, we get

$$\begin{aligned} \mathbb{E}[|\Re\{r_m[t]\}| + |\Im\{r_m[t]\}|] &\approx g_m^{-1} \mathbb{E}[|\Re\{s_m[t - \tau]\}| + |\Im\{s_m[t - \tau]\}|] \\ &\approx g_m^{-1} \sqrt{S}, \end{aligned} \quad (4.37)$$

when using the definition of the QAM constellation in (4.14). We recall that the receive processing is required only in the case of QAM signaling. Based on (4.37), we propose a blind estimation method to obtain the scaling factor g_m for each user prior to the decision operation; that is

$$g_m = T \cdot \frac{\sqrt{S}}{\sum_{t=1}^T |\Re\{r_m[t]\}| + |\Im\{r_m[t]\}|}, \quad (4.38)$$

where T is the length the received sequence. The method does not require any feedback or training from the BS nor any knowledge of the noise plus interference power at the user terminal.

4.8 Potential Dual Uplink System Model with the Optimal CEQ

As explained in Section 3.6, the CEQ can be linearized and the resulting linear system model with the optimal CEQ is shown in Fig. 4.2. The potential dual uplink system model is provided in Fig. 4.3. The MSE duality between both system models will be investigated throughout this thesis.

We have the following assumptions for the potential dual system model

$$\mathbf{F}_{\ell'} = \frac{1}{\beta} \mathbf{P}_{\ell'}^H, \quad \ell' = 0, \dots, L_p - 1, \quad (4.39)$$

$$\mathbf{T} = \beta \mathbf{C}_{\eta\eta}^{1/2, H} \mathbf{G}^H, \quad (4.40)$$

$$\mathbf{C}_{ss} = \mathbf{C}_{ss}^{\text{UL}} = \sigma_s^2 \mathbf{I}_M, \quad (4.41)$$

and

$$\mathbf{C}_{\eta\eta}^{\text{UL}} = \mathbf{I}_N. \quad (4.42)$$

The transmit power constraint in (4.20) for maximal power exploitation can be rewritten as

$$\alpha_q \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \mathbf{P}_{\ell'} \mathbf{C}_{ss} \mathbf{P}_{\ell'}^H \right) = P_{\text{tx}}. \quad (4.43)$$

Hence, the corresponding scaling factor β can be obtained by plugging (4.39) in the transmit power constraint in (4.43) as

$$\beta = \sqrt{\frac{P_{\text{tx}}}{\alpha_q \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \mathbf{F}_{\ell'}^H \mathbf{C}_{ss} \mathbf{F}_{\ell'} \right)}}. \quad (4.44)$$

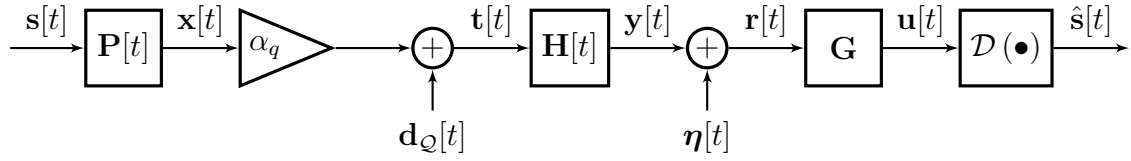


Fig. 4.2: Downlink system model with Bussgang decomposition.

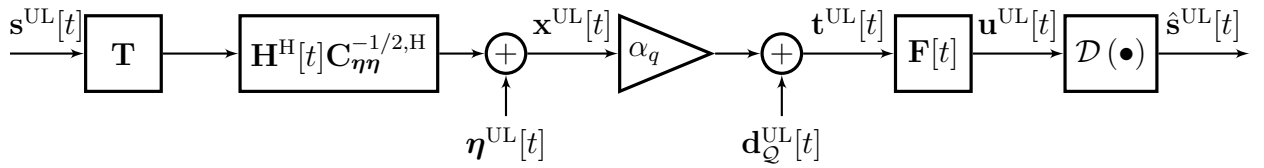


Fig. 4.3: Uplink system model with Bussgang decomposition.

Part I

Linear Transmit Signal Processing

5. Flat-Fading Channels

5.1 Input-Output Relationship

In this section, we consider linear precoding for flat-fading channels; that is $L = 1$. This implies that the precoder impulse response matrix $\mathbf{P}[t]$ reduces to the matrix \mathbf{P} . Accordingly, it holds that

$$\mathbf{x}[t] = \mathbf{P}\mathbf{s}[t]. \quad (5.1)$$

The time indexes are dropped out in the sequel. We recall the input-output relationship in (4.8)

$$\begin{aligned} \mathbf{u} &= \mathbf{G}(\mathbf{H}\mathbf{t} + \boldsymbol{\eta}) \\ &= \mathbf{G}(\mathbf{H}\mathcal{Q}_{\text{CE}}(\mathbf{P}\mathbf{s}) + \boldsymbol{\eta}). \end{aligned} \quad (5.2)$$

5.2 Optimization Problem

The precoding task consists of finding the optimal linear precoder \mathbf{P}_{opt} and the diagonal positive real-valued receive filter \mathbf{G}_{opt} that minimize the MSE between the desired and the received signals, \mathbf{s} and \mathbf{u} under the transmit power constraint given in (4.20). The MMSE optimization problem is expressed as

$$\{\mathbf{P}_{\text{opt}}, \mathbf{G}_{\text{opt}}\} = \arg \min_{\mathbf{P}, \mathbf{G}} \text{E} [\|\mathbf{u} - \mathbf{s}\|_2^2] \quad \text{s.t.} \quad \alpha_q \text{tr}(\mathbf{C}_{\mathbf{xx}}) \leq P_{\text{tx}}. \quad (5.3)$$

5.2.1 MSE

In general, the MSE is expressed as

$$\text{E} [\|\mathbf{u} - \mathbf{s}\|_2^2] = \text{tr}(\mathbf{G}\mathbf{H}\mathbf{C}_{\text{tt}}\mathbf{H}^{\text{H}}\mathbf{G} + \mathbf{G}\mathbf{C}_{\boldsymbol{\eta}\boldsymbol{\eta}}\mathbf{G}) - \text{tr}(\mathbf{G}\mathbf{H}\mathbf{C}_{\text{ts}}) - \text{tr}(\mathbf{C}_{\text{st}}\mathbf{H}^{\text{H}}\mathbf{G}) + \text{tr}(\mathbf{C}_{\text{ss}}). \quad (5.4)$$

Due to the central limit theorem [68], the entries in \mathbf{x} are approximately Gaussian distributed for massive MIMO systems, where a large number of users is served by a large number of transmit antennas. Hence, Price's theorem can be applied to obtain the expressions of the covariance matrices \mathbf{C}_{tt} and \mathbf{C}_{tx} according to (3.66) and (3.65). Consequently, the MSE expression calculates to

$$\begin{aligned} \text{E} [\|\mathbf{u} - \mathbf{s}\|_2^2] &= \frac{Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \text{tr}(\arcsin(\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\}) \boldsymbol{\Xi}\mathbf{H}^{\text{H}}\mathbf{G}^2\mathbf{H}\boldsymbol{\Xi}) \\ &\quad + \text{tr}(\mathbf{G}\mathbf{C}_{\boldsymbol{\eta}\boldsymbol{\eta}}\mathbf{G}) - \alpha_q \text{tr}(\mathbf{G}\mathbf{H}\mathbf{P}\mathbf{C}_{\text{ss}}) - \alpha_q \text{tr}(\mathbf{C}_{\text{ss}}\mathbf{P}^{\text{H}}\mathbf{H}^{\text{H}}\mathbf{G}) + \text{tr}(\mathbf{C}_{\text{ss}}). \end{aligned} \quad (5.5)$$

According to (5.1), the covariance matrix $\mathbf{C}_{\mathbf{xx}}$ is expressed in terms of the precoder \mathbf{P} as

$$\begin{aligned}\mathbf{C}_{\mathbf{xx}} &= \mathbb{E} [\mathbf{xx}^H] \\ &= \mathbf{P}\mathbf{C}_{\mathbf{ss}}\mathbf{P}^H.\end{aligned}\quad (5.6)$$

Therefore, we obtain the following expressions

$$\mathbf{R}_{\mathbf{xx}} = \text{diag} (\mathbf{P}\mathbf{C}_{\mathbf{ss}}\mathbf{P}^H)^{-1/2} \mathbf{P}\mathbf{C}_{\mathbf{ss}}\mathbf{P}^H \text{diag} (\mathbf{P}\mathbf{C}_{\mathbf{ss}}\mathbf{P}^H)^{-1/2}, \quad (5.7)$$

$$\mathbf{\Xi} = \text{diag} (\mathbf{P}\mathbf{C}_{\mathbf{ss}}\mathbf{P}^H)^{1/2}. \quad (5.8)$$

5.2.2 Transmit Power Constraint

After plugging (5.6) in the transmit power constraint in (4.20) and for maximal exploitation of the available power, we get

$$\alpha_q \text{tr} (\mathbf{P}\mathbf{C}_{\mathbf{ss}}\mathbf{P}^H) = P_{\text{tx}}. \quad (5.9)$$

As explained in Section 4.5, we differentiate between two cases

- equal power allocation, i.e. $\text{diag} (\mathbf{P}\mathbf{C}_{\mathbf{ss}}\mathbf{P}^H) = \frac{P_{\text{tx}}}{\alpha_q N} \mathbf{I}_N$,
- and
- unequal power allocation, i.e. $\text{tr} (\mathbf{P}\mathbf{C}_{\mathbf{ss}}\mathbf{P}^H) = \frac{P_{\text{tx}}}{\alpha_q}$.

5.3 Precoder Designs in the Primal Domain

5.3.1 Precoder Design Based on Gradient Projection Algorithm

Since the MSE expression in (5.5) is highly non-linear in $\mathbf{C}_{\mathbf{xx}} = \mathbf{P}\mathbf{C}_{\mathbf{ss}}\mathbf{P}^H$ and thus in \mathbf{P} , we cannot find a closed-form expression for \mathbf{P}_{opt} . Therefore, we use the Gradient-Projection algorithm as described in Algorithm 1 and Algorithm 2 for equal and unequal power allocation, respectively. To this end, we have to compute the derivatives of the MSE w.r.t. \mathbf{P} and \mathbf{G} . First, the derivative of the MSE w.r.t. \mathbf{P} is expressed as

$$\begin{aligned}\frac{\partial \mathbb{E} [\|\mathbf{u} - \mathbf{s}\|_2^2]}{\partial \mathbf{P}} &= \sum_{n=1}^N \sum_{m=1}^M \frac{\partial \mathbb{E} [\|\mathbf{u} - \mathbf{s}\|_2^2]}{\partial \mathbf{e}_n^T \mathbf{P} \mathbf{e}_m} \mathbf{e}_n \mathbf{e}_m^T \\ &= \frac{Q}{\pi} \sin^2(\psi) \\ &\quad \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \sum_{n=1}^N \sum_{m=1}^M \frac{\partial \text{tr} (\arcsin (\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\}) \mathbf{\Xi} \mathbf{H}^H \mathbf{G}^2 \mathbf{H} \mathbf{\Xi})}{\partial \mathbf{e}_n^T \mathbf{P} \mathbf{e}_m} \mathbf{e}_n \mathbf{e}_m^T \\ &\quad - \alpha_q \mathbf{H}^T \mathbf{G} \mathbf{C}_{\mathbf{ss}}^*.\end{aligned}\quad (5.10)$$

The challenging term in the above derivative expression is

$$\begin{aligned}\frac{\partial \text{tr} (\arcsin (\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\}) \mathbf{\Xi} \mathbf{H}^H \mathbf{G}^2 \mathbf{H} \mathbf{\Xi})}{\partial \mathbf{e}_n^T \mathbf{P} \mathbf{e}_m} &= \text{tr} \left(\frac{\partial \arcsin (\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\})}{\partial \mathbf{e}_n^T \mathbf{P} \mathbf{e}_m} \mathbf{\Xi} \mathbf{H}^H \mathbf{G}^2 \mathbf{H} \mathbf{\Xi} \right) \\ &\quad + \text{tr} \left(\arcsin (\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\}) \frac{\partial \mathbf{\Xi} \mathbf{H}^H \mathbf{G}^2 \mathbf{H} \mathbf{\Xi}}{\partial \mathbf{e}_n^T \mathbf{P} \mathbf{e}_m} \right).\end{aligned}\quad (5.11)$$

Algorithm 1 Gradient Projection Algorithm to obtain the QWF-Price precoder with equal power allocation.

- 1: Initialization
 $\mathbf{P}_{(0)}, \mathbf{G}_{(0)} = \mathbf{g}_{\text{opt}}(\mathbf{P}_{(0)}), \mu = 10, \epsilon = 10^{-4}$ and $n = 0$
 - 2: **repeat**
 - 3: $\mathbf{P}_{(n+1)} = \mathbf{P}_{(n)} - \mu \left(\frac{\partial \text{MSE}_{(n)}}{\partial \mathbf{P}} \right)^*$
 - 4: $\mathbf{P}_{(n+1)} = \sqrt{\frac{P_{\text{tx}}}{\alpha_q N}} \text{diag} \left(\mathbf{P}_{(n+1)} \mathbf{C}_{\text{ss}} \mathbf{P}_{(n+1)}^{\text{H}} \right)^{-1/2} \mathbf{P}_{(n+1)}$ {Equal power allocation constraint}
 - 5: $\mathbf{G}_{(n+1)} = \mathbf{g}_{\text{opt}}(\mathbf{P}_{(n+1)})$
 - 6: **if** $\text{MSE}_{(n+1)} > \text{MSE}_{(n)}$ **then**
 - 7: $\mu = \mu/2$
 - 8: **else**
 - 9: $n = n + 1$
 - 10: **end if**
 - 11: **until** $\frac{|\text{MSE}_{(n+1)} - \text{MSE}_{(n)}|}{|\text{MSE}_{(n)}|} \leq \epsilon$
-

Algorithm 2 Gradient Projection Algorithm to obtain the QWF-Price precoder with unequal power allocation.

- 1: Initialization
 $\mathbf{P}_{(0)}, \mathbf{G}_{(0)} = \mathbf{g}_{\text{opt}}(\mathbf{P}_{(0)}), \mu = 10, \epsilon = 10^{-4}$ and $n = 0$
 - 2: **repeat**
 - 3: $\mathbf{P}_{(n+1)} = \mathbf{P}_{(n)} - \mu \left(\frac{\partial \text{MSE}_{(n)}}{\partial \mathbf{P}} \right)^*$
 - 4: $\mathbf{P}_{(n+1)} = \sqrt{\frac{P_{\text{tx}}}{\alpha_q}} \text{tr} \left(\mathbf{P}_{(n+1)} \mathbf{C}_{\text{ss}} \mathbf{P}_{(n+1)}^{\text{H}} \right)^{-1/2} \mathbf{P}_{(n+1)}$ {Unequal power allocation constraint}
 - 5: $\mathbf{G}_{(n+1)} = \mathbf{g}_{\text{opt}}(\mathbf{P}_{(n+1)})$
 - 6: **if** $\text{MSE}_{(n+1)} > \text{MSE}_{(n)}$ **then**
 - 7: $\mu = \mu/2$
 - 8: **else**
 - 9: $n = n + 1$
 - 10: **end if**
 - 11: **until** $\frac{|\text{MSE}_{(n+1)} - \text{MSE}_{(n)}|}{|\text{MSE}_{(n)}|} \leq \epsilon$
-

According to A5, the derivative term related to $\arcsin(\bullet)$ calculates to

$$\begin{aligned}
\frac{\partial \arcsin(\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\})}{\partial \mathbf{e}_n^T \mathbf{P} \mathbf{e}_m} &= \text{nondiag} \left(\mathbf{1}_{N,N} - \text{nondiag} \left(\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\} \right)^{\circ 2} \right)^{\circ -1/2} \circ \\
&\quad \frac{\partial \Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\}}{\partial \mathbf{e}_n^T \mathbf{P} \mathbf{e}_m} \\
&= \frac{1}{2} \text{nondiag} \left(\mathbf{1}_{N,N} - \text{nondiag} \left(\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\} \right)^{\circ 2} \right)^{\circ -1/2} \circ \\
&\quad \left(\text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \mathbf{e}_n \mathbf{e}_m^T e^{-j2\Delta k\psi} \mathbf{C}_{\mathbf{ss}} \mathbf{P}^H \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \right. \\
&\quad + \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \mathbf{P}^* \mathbf{C}_{\mathbf{ss}}^* e^{j2\Delta k\psi} \mathbf{e}_m \mathbf{e}_n^T \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \\
&\quad - (\mathbf{e}_n^T \mathbf{C}_{\mathbf{xx}} \mathbf{e}_n)^{-3/2} \mathbf{e}_m^T \mathbf{C}_{\mathbf{ss}} \mathbf{P}^H \mathbf{e}_n \mathbf{e}_n \mathbf{e}_n^T \Re\{\mathbf{C}_{\mathbf{xx}} e^{-j2\Delta k\psi}\} \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \\
&\quad \left. - \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \Re\{\mathbf{C}_{\mathbf{xx}} e^{-j2\Delta k\psi}\} (\mathbf{e}_n^T \mathbf{C}_{\mathbf{xx}} \mathbf{e}_n)^{-3/2} \mathbf{e}_m^T \mathbf{C}_{\mathbf{ss}} \mathbf{P}^H \mathbf{e}_n \mathbf{e}_n \mathbf{e}_n^T \right). \tag{5.12}
\end{aligned}$$

Moreover, we have

$$\begin{aligned}
\frac{\partial \Xi \mathbf{H}^H \mathbf{G}^2 \mathbf{H} \Xi}{\partial \mathbf{e}_n^T \mathbf{P} \mathbf{e}_m} &= \alpha_q \frac{1}{2} (\mathbf{e}_n^T \mathbf{C}_{\mathbf{xx}} \mathbf{e}_n)^{-1/2} \mathbf{e}_m^T \mathbf{C}_{\mathbf{ss}} \mathbf{P}^H \mathbf{e}_n \mathbf{e}_n \mathbf{e}_n^T \mathbf{H}^H \mathbf{G}^2 \mathbf{H} \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2} \\
&\quad + \alpha_q \frac{1}{2} \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2} \mathbf{H}^H \mathbf{G}^2 \mathbf{H} (\mathbf{e}_n^T \mathbf{C}_{\mathbf{xx}} \mathbf{e}_n)^{-1/2} \mathbf{e}_m^T \mathbf{C}_{\mathbf{ss}} \mathbf{P}^H \mathbf{e}_n \mathbf{e}_n \mathbf{e}_n^T. \tag{5.13}
\end{aligned}$$

Plugging (5.12) and (5.13) in (5.11) and applying the property

$$\text{tr}((\mathbf{C} \circ \mathbf{v}_1 \mathbf{v}_2^T) \mathbf{D}) = \text{tr}(\mathbf{v}_2^T (\mathbf{D} \circ \mathbf{C}^T) \mathbf{v}_1), \tag{5.14}$$

we obtain

$$\begin{aligned}
&\frac{\partial \text{tr}(\arcsin(\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\}) \Xi \mathbf{H}^H \mathbf{G}^2 \mathbf{H} \Xi)}{\partial \mathbf{e}_n^T \mathbf{P} \mathbf{e}_m} = \\
&\quad \frac{1}{2} \left(\mathbf{e}_m^T e^{-j2\Delta k\psi} \mathbf{C}_{\mathbf{ss}} \mathbf{P}^H \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \Omega \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \mathbf{e}_n \right. \\
&\quad + \mathbf{e}_n^T e^{j2\Delta k\psi} \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \Omega \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \mathbf{P}^* \mathbf{C}_{\mathbf{ss}}^* \mathbf{e}_m \\
&\quad - \mathbf{e}_n^T \Re\{\mathbf{C}_{\mathbf{xx}} e^{-j2\Delta k\psi}\} \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \Omega (\mathbf{e}_n^T \mathbf{C}_{\mathbf{xx}} \mathbf{e}_n)^{-3/2} \mathbf{e}_m^T \mathbf{C}_{\mathbf{ss}} \mathbf{P}^H \mathbf{e}_n \mathbf{e}_n \\
&\quad \left. - \mathbf{e}_n^T \Omega \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \Re\{\mathbf{C}_{\mathbf{xx}} e^{-j2\Delta k\psi}\} (\mathbf{e}_n^T \mathbf{C}_{\mathbf{xx}} \mathbf{e}_n)^{-3/2} \mathbf{e}_m^T \mathbf{C}_{\mathbf{ss}} \mathbf{P}^H \mathbf{e}_n \mathbf{e}_n \right) \\
&\quad + \alpha_q \frac{1}{2} \left(\mathbf{e}_n^T \mathbf{H}^H \mathbf{G}^2 \mathbf{H} \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2} \arcsin(\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\}) \mathbf{e}_n (\mathbf{e}_n^T \mathbf{C}_{\mathbf{xx}} \mathbf{e}_n)^{-1/2} \mathbf{e}_n^T \mathbf{P}^* \mathbf{C}_{\mathbf{ss}}^* \mathbf{e}_m \right. \\
&\quad \left. + \mathbf{e}_m^T \mathbf{C}_{\mathbf{ss}} \mathbf{P}^H \mathbf{e}_n (\mathbf{e}_n^T \mathbf{C}_{\mathbf{xx}} \mathbf{e}_n)^{-1/2} \mathbf{e}_n^T \arcsin(\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\}) \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2} \mathbf{H}^H \mathbf{G}^2 \mathbf{H} \mathbf{e}_n \right), \tag{5.15}
\end{aligned}$$

where

$$\Omega = (\Xi \mathbf{H}^H \mathbf{G}^2 \mathbf{H} \Xi) \circ \text{nondiag} \left(\mathbf{1}_{N,N} - \text{nondiag} \left(\Re\{\mathbf{R}_{\mathbf{xx}}^* e^{-j2\Delta k\psi}\} \right)^{\circ 2} \right)^{\circ -1/2}. \tag{5.16}$$

Thus, the derivative of the MSE w.r.t. \mathbf{P} is given by

$$\begin{aligned}
 \frac{\partial E [\|\mathbf{u} - \mathbf{s}\|_2^2]}{\partial \mathbf{P}} &= \frac{1}{2} \frac{Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \left(\text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \left(e^{-j2\Delta k\psi} \boldsymbol{\Omega}^T + e^{j2\Delta k\psi} \boldsymbol{\Omega} \right. \right. \\
 &\quad - \text{diag}(\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\} \boldsymbol{\Omega}) - \text{diag}(\boldsymbol{\Omega} \Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\}) \\
 &\quad + \text{diag}(\boldsymbol{\Xi} \mathbf{H}^H \mathbf{G}^2 \mathbf{H} \boldsymbol{\Xi} \arcsin(\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\})) \\
 &\quad \left. \left. + \text{diag}(\arcsin(\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\}) \boldsymbol{\Xi} \mathbf{H}^H \mathbf{G}^2 \mathbf{H} \boldsymbol{\Xi}) \right) \right) \\
 &\quad \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \mathbf{P}^* \mathbf{C}_{\mathbf{ss}}^* \\
 &\quad - \alpha_q \mathbf{H}^T \mathbf{G} \mathbf{C}_{\mathbf{ss}}^*. \tag{5.17}
 \end{aligned}$$

Second, the derivative of the MSE expression in (5.5) w.r.t. \mathbf{G} calculates to

$$\frac{\partial E [\|\hat{\mathbf{u}} - \mathbf{s}\|_2^2]}{\partial \mathbf{G}} = 2 \text{diag}(\mathbf{H} \mathbf{C}_{\mathbf{tt}} \mathbf{H}^H + \mathbf{C}_{\eta\eta}) - \text{diag}(\mathbf{H} \mathbf{C}_{\mathbf{ts}}) - \text{diag}(\mathbf{H}^* \mathbf{C}_{\mathbf{ts}}^*). \tag{5.18}$$

Thus, the optimal filter \mathbf{G}_{opt} is obtained by setting (5.18) equal to a zero matrix and is given by

$$\mathbf{G}_{\text{opt}} = \mathbf{g}_{\text{opt}}(\mathbf{P}) = \left| \text{diag}(\mathbf{H} \mathbf{C}_{\mathbf{tt}} \mathbf{H}^H + \mathbf{C}_{\eta\eta})^{-1} \text{diag}(\Re\{\mathbf{H} \mathbf{C}_{\mathbf{ts}}\}) \right|, \tag{5.19}$$

where the operator $|\bullet|$ is applied element-wise to the matrix entries.

5.3.2 Precoder Design Based on LCA

In this section, we aim at getting a closed-form approximation for the MMSE precoder by using the LCA to obtain a linear expression of the MSE. We recall the linear precoder design in [25] and apply it to the case of QCE transmit signals. The precoder aims at minimizing the MSE under a transmit power constraint, where the receive processing is assumed to be a scaled identity matrix, i.e. $\mathbf{G} = g \mathbf{I}_M$. Using the LCA, the MSE expression simplifies to

$$\begin{aligned}
 E [\|\mathbf{u} - \mathbf{s}\|_2^2] &= \text{tr} \left(|g|^2 \mathbf{H} (\alpha_q^2 \mathbf{P} \mathbf{C}_{\mathbf{ss}} \mathbf{P}^H + \alpha_q \rho_q \text{diag}(\mathbf{P} \mathbf{C}_{\mathbf{ss}} \mathbf{P}^H)) \mathbf{H}^H \right. \\
 &\quad \left. - g \mathbf{H} \alpha_q \mathbf{P} \mathbf{C}_{\mathbf{ss}} - g^* \alpha_q \mathbf{C}_{\mathbf{ss}} \mathbf{P}^H \mathbf{H}^H + |g|^2 \mathbf{C}_{\eta\eta} + \mathbf{C}_{\mathbf{ss}} \right). \tag{5.20}
 \end{aligned}$$

The Lagrangian function is expressed by

$$\begin{aligned}
 \mathcal{L}(\mathbf{P}, g, \lambda) &= \text{tr} \left(|g|^2 \mathbf{H} (\alpha_q^2 \mathbf{P} \mathbf{C}_{\mathbf{ss}} \mathbf{P}^H + \alpha_q \beta_q \text{diag}(\mathbf{P} \mathbf{C}_{\mathbf{ss}} \mathbf{P}^H)) \mathbf{H}^H \right. \\
 &\quad \left. - g \alpha_q \mathbf{H} \mathbf{P} \mathbf{C}_{\mathbf{ss}} - g^* \alpha_q \mathbf{C}_{\mathbf{ss}} \mathbf{P}^H \mathbf{H}^H + |g|^2 \mathbf{C}_{\eta\eta} + \mathbf{C}_{\mathbf{ss}} \right) + \lambda (\alpha_q \text{tr}(\mathbf{P} \mathbf{C}_{\mathbf{ss}} \mathbf{P}^H) - P_{\text{tx}}). \tag{5.21}
 \end{aligned}$$

The Karush-Kuhn-Tucker (KKT) equations are then given by

$$\begin{aligned} \frac{\partial \mathcal{L}(\mathbf{P}, g, \lambda)}{\partial \mathbf{P}} &= |g|^2 (\alpha_q^2 \mathbf{H}^T \mathbf{H}^* \mathbf{P}^* \mathbf{C}_{\text{ss}}^T + \alpha_q \beta_q \text{diag}(\mathbf{H}^T \mathbf{H}^*) \mathbf{P}^* \mathbf{C}_{\text{ss}}^T) - g \alpha_q \mathbf{H}^T \mathbf{C}_{\text{ss}}^T + \lambda \alpha_q \mathbf{P}^* \mathbf{C}_{\text{ss}}^T \\ &= \mathbf{0}_{N,M}, \end{aligned} \quad (5.22)$$

$$\frac{\partial \mathcal{L}(\mathbf{P}, g, \lambda)}{\partial g} = g^* \text{tr}(\mathbf{H} (\alpha_q^2 \mathbf{P} \mathbf{C}_{\text{ss}} \mathbf{P}^H + \alpha_q \beta_q \text{diag}(\mathbf{P} \mathbf{C}_{\text{ss}} \mathbf{P}^H) \mathbf{H}^H + \mathbf{C}_{\eta\eta})) - \alpha_q \text{tr}(\mathbf{H} \mathbf{P} \mathbf{C}_{\text{ss}}) = 0, \quad (5.23)$$

and

$$\frac{\partial \mathcal{L}(\mathbf{P}, g, \lambda)}{\partial \lambda} = \alpha_q \text{tr}(\mathbf{P} \mathbf{C}_{\text{ss}} \mathbf{P}^H) - P_{\text{tx}} = 0. \quad (5.24)$$

Multiplying (5.22) by \mathbf{P}^T from the left side and taking the trace leads to

$$\begin{aligned} |g|^2 \text{tr}(\alpha_q^2 \mathbf{P}^T \mathbf{H}^T \mathbf{H}^* \mathbf{P}^* \mathbf{C}_{\text{ss}}^T + \alpha_q \beta_q \mathbf{P}^T \text{diag}(\mathbf{H}^T \mathbf{H}^*) \mathbf{P}^* \mathbf{C}_{\text{ss}}^T) - g \alpha_q \text{tr}(\mathbf{P}^T \mathbf{H}^T \mathbf{C}_{\text{ss}}^T) \\ + \lambda \alpha_q \text{tr}(\mathbf{P}^T \mathbf{P}^* \mathbf{C}_{\text{ss}}^T) = 0. \end{aligned} \quad (5.25)$$

From (5.23), we get

$$\alpha_q \text{tr}(\mathbf{H} \mathbf{P} \mathbf{C}_{\text{ss}}) = g^* \text{tr}(\mathbf{H} (\alpha_q^2 \mathbf{P} \mathbf{C}_{\text{ss}} \mathbf{P}^H + \alpha_q \beta_q \text{diag}(\mathbf{P} \mathbf{C}_{\text{ss}} \mathbf{P}^H) \mathbf{H}^H + \mathbf{C}_{\eta\eta})), \quad (5.26)$$

which when inserted in (5.25) gives the expression of λ

$$\lambda = \frac{|g|^2 \text{tr}(\mathbf{C}_{\eta\eta})}{\alpha_q \text{tr}(\mathbf{P} \mathbf{C}_{\text{ss}} \mathbf{P}^H)} = \frac{|g|^2 \text{tr}(\mathbf{C}_{\eta\eta})}{P_{\text{tx}}}, \quad (5.27)$$

where we used (5.24) and the properties $\text{tr}(\mathbf{A} \text{diag}(\mathbf{B})) = \text{tr}(\text{diag}(\mathbf{A}) \mathbf{B})$ and $\text{tr}(\mathbf{A}) = \text{tr}(\mathbf{A})^T$. Inserting (5.27) in (5.22) and solving it for \mathbf{P} , we obtain

$$\mathbf{P} = \frac{1}{g} \left(\alpha_q \mathbf{H}^H \mathbf{H} + \beta_q \text{diag}(\mathbf{H}^H \mathbf{H}) + \frac{\text{tr}(\mathbf{C}_{\eta\eta})}{P_{\text{tx}}} \mathbf{I}_N \right)^{-1} \mathbf{H}^H. \quad (5.28)$$

The optimal g is found by satisfying (5.24) with \mathbf{P} from (5.28); that is

$$g = \sqrt{\frac{\alpha_q}{P_{\text{tx}}}} \sqrt{\text{tr} \left(\left(\alpha_q \mathbf{H}^H \mathbf{H} + \beta_q \text{diag}(\mathbf{H}^H \mathbf{H}) + \frac{\text{tr}(\mathbf{C}_{\eta\eta})}{P_{\text{tx}}} \mathbf{I}_N \right)^{-2} \mathbf{H}^H \mathbf{C}_{\text{ss}} \mathbf{H} \right)}. \quad (5.29)$$

This precoder design was already derived in [25] for low resolution cartesian DACs at the transmitter. Whether we use polar or cartesian DACs, the expression of the WF precoder that takes into account the quantization distortions based on LCA turns to be the same. The coefficients α_q and β_q have to be defined accordingly.

5.4 Dual Optimization Problem

5.4.1 Does a Dual Problem Exist?

In the above sections, the MMSE precoder was designed based on Price's theorem for equal and unequal power allocation and on LCA for only unequal power allocation. With Price's

theorem, no closed-form expression of the precoder could be found but the optimal precoder is found by an iterative approach. However, with LCA an approximate expression of the optimal MMSE precoder is given in a closed-form. One might ask whether it is possible to get a closed-form expression for the exact optimal MMSE in the dual domain; that is when applying Price's theorem. To this end, we first have to find the dual uplink system model. Second, find the optimal dual filter, hopefully in closed-form expression. Finally, convert the expressions in the primal domain, i.e. downlink scenario, and obtain the expression of the MMSE precoder.

According to Fig. 4.2, the corresponding MSE expression for flat-fading channels, after dropping out the time index, is given by

$$\begin{aligned}
\text{MSE}^{\text{DL}} &= \text{E} \left[\|\mathbf{u} - \mathbf{s}\|_2^2 \right] \\
&= \alpha_q^2 \text{tr}(\mathbf{G}\mathbf{H}\mathbf{P}\mathbf{C}_{\text{ss}}\mathbf{P}^{\text{H}}\mathbf{H}^{\text{H}}\mathbf{G}) + \text{tr}(\mathbf{G}\mathbf{H}\mathbf{C}_{\text{d}_Q\text{d}_Q}\mathbf{H}^{\text{H}}\mathbf{G}) \\
&\quad - \alpha_q \text{tr}(\mathbf{G}\mathbf{H}\mathbf{P}\mathbf{C}_{\text{ss}}) - \alpha_q \text{tr}(\mathbf{C}_{\text{ss}}\mathbf{P}^{\text{H}}\mathbf{H}^{\text{H}}\mathbf{G}) + \text{tr}(\mathbf{G}\mathbf{C}_{\eta\eta}\mathbf{G}) + \text{tr}(\mathbf{C}_{\text{ss}}) \\
&\stackrel{(3.85)}{=} \alpha_q^2 \text{tr}(\mathbf{G}\mathbf{H}\mathbf{P}\mathbf{C}_{\text{ss}}\mathbf{P}^{\text{H}}\mathbf{H}^{\text{H}}\mathbf{G}) + \alpha_q\beta_q \text{tr}(\mathbf{G}\mathbf{H} \text{diag}(\mathbf{P}\mathbf{C}_{\text{ss}}\mathbf{P}^{\text{H}}) \mathbf{H}^{\text{H}}\mathbf{G}) \\
&\quad + \frac{Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \text{tr} \left(\mathbf{G}\mathbf{H}\mathbf{\Xi} \arcsin \left(\Re\{\text{nondiag}(\mathbf{R}_{\text{xx}}) e^{-j(2\Delta k\psi)}\} \right) \mathbf{\Xi}\mathbf{H}^{\text{H}}\mathbf{G} \right) \\
&\quad - \alpha_q^2 \text{tr}(\mathbf{G}\mathbf{H} \text{nondiag}(\mathbf{P}\mathbf{C}_{\text{ss}}\mathbf{P}^{\text{H}}) \mathbf{H}^{\text{H}}\mathbf{G}) \\
&\quad + \text{tr}(\mathbf{G}\mathbf{C}_{\eta\eta}\mathbf{G}) - \alpha_q \text{tr}(\mathbf{G}\mathbf{H}\mathbf{P}\mathbf{C}_{\text{ss}}) - \alpha_q \text{tr}(\mathbf{C}_{\text{ss}}\mathbf{P}^{\text{H}}\mathbf{H}^{\text{H}}\mathbf{G}) + \text{tr}(\mathbf{C}_{\text{ss}}) \\
&= \alpha_q \text{tr}(\mathbf{G}\mathbf{H} \text{diag}(\mathbf{P}\mathbf{C}_{\text{ss}}\mathbf{P}^{\text{H}}) \mathbf{H}^{\text{H}}\mathbf{G}) \\
&\quad + \frac{Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \text{tr} \left(\mathbf{G}\mathbf{H}\mathbf{\Xi} \arcsin \left(\Re\{\text{nondiag}(\mathbf{R}_{\text{xx}}) e^{-j(2\Delta k\psi)}\} \right) \mathbf{\Xi}\mathbf{H}^{\text{H}}\mathbf{G} \right) \\
&\quad + \text{tr}(\mathbf{G}\mathbf{C}_{\eta\eta}\mathbf{G}) - \alpha_q \text{tr}(\mathbf{G}\mathbf{H}\mathbf{P}\mathbf{C}_{\text{ss}}) - \alpha_q \text{tr}(\mathbf{C}_{\text{ss}}\mathbf{P}^{\text{H}}\mathbf{H}^{\text{H}}\mathbf{G}) + \text{tr}(\mathbf{C}_{\text{ss}}). \tag{5.30}
\end{aligned}$$

The potential dual uplink system model is illustrated in Fig. 4.3, where $\mathbf{d}_Q^{\text{UL}}[t] = \mathbf{d}_Q^{\text{UL}}$ can be freely designed and is not necessarily defined by a CEQ at the receiver. The statistical properties of \mathbf{d}_Q^{UL} have to be chosen such that an MSE duality between both systems can be achieved. The MSE expression for the uplink system for flat-fading channels is expressed as

$$\begin{aligned}
\text{MSE}^{\text{UL}} &= \text{E} \left[\|\mathbf{u}^{\text{UL}} - \mathbf{s}^{\text{UL}}\|_2^2 \right] \\
&= \alpha_q^2 \text{tr}(\mathbf{F}\mathbf{H}^{\text{H}}\mathbf{C}_{\eta\eta}^{-1/2,\text{H}}\mathbf{T}\mathbf{C}_{\text{ss}}^{\text{UL}}\mathbf{T}^{\text{H}}\mathbf{C}_{\eta\eta}^{-1/2}\mathbf{H}\mathbf{F}^{\text{H}}) + \text{tr}(\mathbf{F}\mathbf{C}_{\text{d}_Q\text{d}_Q}^{\text{UL}}\mathbf{F}^{\text{H}}) + \alpha_q^2 \text{tr}(\mathbf{F}\mathbf{C}_{\eta\eta}^{\text{UL}}\mathbf{F}^{\text{H}}) \\
&\quad - \alpha_q \text{tr}(\mathbf{F}\mathbf{H}^{\text{H}}\mathbf{C}_{\eta\eta}^{-1/2,\text{H}}\mathbf{T}\mathbf{C}_{\text{ss}}^{\text{UL}}) - \alpha_q \text{tr}(\mathbf{C}_{\text{ss}}^{\text{UL}}\mathbf{T}^{\text{H}}\mathbf{C}_{\eta\eta}^{-1/2}\mathbf{H}\mathbf{F}^{\text{H}}) + \text{tr}(\mathbf{C}_{\text{ss}}^{\text{UL}}), \tag{5.31}
\end{aligned}$$

where

$$\mathbf{u}^{\text{UL}} = \mathbf{u}^{\text{UL}}[t] \tag{5.32}$$

$$\mathbf{s}^{\text{UL}} = \mathbf{s}^{\text{UL}}[t], \tag{5.33}$$

and

$$\mathbf{F} = \mathbf{F}_0. \tag{5.34}$$

If we assume that $\mathbf{d}_{\mathcal{Q}}^{\text{UL}}$ is the quantization noise that results from quantizing $\mathbf{x}^{\text{UL}} = \mathbf{x}^{\text{UL}}[t]$ and is uncorrelated with it and hence we use the expression in (3.85) for $\mathbf{C}_{\mathbf{d}_{\mathcal{Q}}^{\text{UL}}}^{\text{UL}}$, we get

$$\begin{aligned}
\text{MSE}^{\text{UL}} &= \alpha_q^2 \text{tr} \left(\mathbf{F} \mathbf{H}^{\text{H}} \mathbf{C}_{\eta\eta}^{-1/2, \text{H}} \mathbf{T} \mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{T}^{\text{H}} \mathbf{C}_{\eta\eta}^{-1/2} \mathbf{H} \mathbf{F}^{\text{H}} \right) \\
&+ \alpha_q \beta_q \text{tr} \left(\mathbf{F} \text{diag} \left(\mathbf{H}^{\text{H}} \mathbf{C}_{\eta\eta}^{-1/2, \text{H}} \mathbf{T} \mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{T}^{\text{H}} \mathbf{C}_{\eta\eta}^{-1/2} \mathbf{H} + \mathbf{C}_{\eta\eta}^{\text{UL}} \right) \mathbf{F}^{\text{H}} \right) \\
&+ \frac{Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \text{tr} \left(\mathbf{F} \mathbf{\Xi}^{\text{UL}} \arcsin \left(\Re \{ \text{nondiag} \left(\mathbf{R}_{\text{xx}}^{\text{UL}} \right) e^{-j(2\Delta k\psi)} \} \right) \mathbf{\Xi}^{\text{UL}} \mathbf{F}^{\text{H}} \right) \\
&- \alpha_q^2 \text{tr} \left(\mathbf{F} \text{nondiag} \left(\mathbf{C}_{\text{xx}}^{\text{UL}} \right) \mathbf{F}^{\text{H}} \right) \\
&+ \alpha_q^2 \text{tr} \left(\mathbf{F} \mathbf{C}_{\eta\eta}^{\text{UL}} \mathbf{F}^{\text{H}} \right) \\
&- \alpha_q \text{tr} \left(\mathbf{F} \mathbf{H}^{\text{H}} \mathbf{C}_{\eta\eta}^{-1/2, \text{H}} \mathbf{T} \mathbf{C}_{\text{ss}}^{\text{UL}} \right) - \alpha_q \text{tr} \left(\mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{T}^{\text{H}} \mathbf{C}_{\eta\eta}^{-1/2} \mathbf{H} \mathbf{F}^{\text{H}} \right) + \text{tr} \left(\mathbf{C}_{\text{ss}}^{\text{UL}} \right) \\
&\stackrel{(h)}{=} \alpha_q \text{tr} \left(\mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{T}^{\text{H}} \mathbf{C}_{\eta\eta}^{-1/2, \text{H}} \mathbf{H} \text{diag} \left(\mathbf{F}^{\text{H}} \mathbf{F} \right) \mathbf{H}^{\text{H}} \mathbf{C}_{\eta\eta}^{-1/2} \mathbf{T} \right) \\
&+ \frac{Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \text{tr} \left(\mathbf{F} \mathbf{\Xi}^{\text{UL}} \arcsin \left(\Re \{ \text{nondiag} \left(\mathbf{R}_{\text{xx}}^{\text{UL}} \right) e^{-j(2\Delta k\psi)} \} \right) \mathbf{\Xi}^{\text{UL}} \mathbf{F}^{\text{H}} \right) \\
&+ \alpha_q \text{tr} \left(\mathbf{F} \mathbf{C}_{\eta\eta}^{\text{UL}} \mathbf{F}^{\text{H}} \right) \\
&- \alpha_q \text{tr} \left(\mathbf{F} \mathbf{H}^{\text{H}} \mathbf{C}_{\eta\eta}^{-1/2, \text{H}} \mathbf{T} \mathbf{C}_{\text{ss}}^{\text{UL}} \right) - \alpha_q \text{tr} \left(\mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{T}^{\text{H}} \mathbf{C}_{\eta\eta}^{-1/2} \mathbf{H} \mathbf{F}^{\text{H}} \right) + \text{tr} \left(\mathbf{C}_{\text{ss}}^{\text{UL}} \right), \tag{5.35}
\end{aligned}$$

where

$$\mathbf{C}_{\text{xx}}^{\text{UL}} = \mathbf{H}^{\text{H}} \mathbf{C}_{\eta\eta}^{-1/2, \text{H}} \mathbf{T} \mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{T}^{\text{H}} \mathbf{C}_{\eta\eta}^{-1/2} \mathbf{H}, \tag{5.36}$$

$$\mathbf{R}_{\text{xx}}^{\text{UL}} = \text{diag} \left(\mathbf{C}_{\text{xx}}^{\text{UL}} \right)^{-1/2} \mathbf{C}_{\text{xx}}^{\text{UL}} \text{diag} \left(\mathbf{C}_{\text{xx}}^{\text{UL}} \right)^{-1/2}, \tag{5.37}$$

$$\mathbf{\Xi}^{\text{UL}} = \sqrt{\alpha_q} \text{diag} \left(\mathbf{C}_{\text{xx}}^{\text{UL}} \right)^{1/2}, \tag{5.38}$$

and in (h) we used the equality

$$\alpha_q^2 \text{tr} \left(\mathbf{F} \text{nondiag} \left(\mathbf{C}_{\text{xx}}^{\text{UL}} \right) \mathbf{F}^{\text{H}} \right) = \alpha_q^2 \text{tr} \left(\mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{T}^{\text{H}} \mathbf{C}_{\eta\eta}^{-1/2, \text{H}} \mathbf{H} \text{nondiag} \left(\mathbf{F}^{\text{H}} \mathbf{F} \right) \mathbf{H}^{\text{H}} \mathbf{C}_{\eta\eta}^{-1/2} \mathbf{T} \right). \tag{5.39}$$

We aim to match both MSE expressions in (5.30) and (5.35). Applying the identities in (4.39), (4.40), (4.41) and (4.42) for the flat-fading case, i.e. $L = L_p = 1$, we get

$$\mathbf{F} = \frac{1}{\beta} \mathbf{P}^{\text{H}}, \tag{5.40}$$

$$\mathbf{T} = \beta \mathbf{C}_{\eta\eta}^{1/2, \text{H}} \mathbf{G}^{\text{H}}, \tag{5.41}$$

$$\mathbf{C}_{\text{ss}} = \mathbf{C}_{\text{ss}}^{\text{UL}} = \sigma_s^2 \mathbf{I}_M, \tag{5.42}$$

and

$$\mathbf{C}_{\eta\eta}^{\text{UL}} = \mathbf{I}_N. \tag{5.43}$$

Thus, (5.35) calculates to

$$\begin{aligned}
\text{MSE}^{\text{UL}} &= \alpha_q \text{tr} \left(\mathbf{G} \mathbf{H} \text{diag} \left(\mathbf{P} \mathbf{C}_{\text{ss}} \mathbf{P}^{\text{H}} \right) \mathbf{H}^{\text{H}} \mathbf{G} \right) \\
&+ \frac{1}{\beta^2} \frac{Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \text{tr} \left(\mathbf{P}^{\text{H}} \mathbf{\Xi}^{\text{UL}} \arcsin \left(\Re \{ \text{nondiag} \left(\mathbf{R}_{\text{xx}}^{\text{UL}} \right) e^{-j(2\Delta k\psi)} \} \right) \mathbf{\Xi}^{\text{UL}} \mathbf{P} \right) \\
&+ \frac{\alpha_q}{\sigma_s^2 \beta^2} \text{tr} \left(\mathbf{P} \mathbf{C}_{\text{ss}} \mathbf{P}^{\text{H}} \right) - \alpha_q \text{tr} \left(\mathbf{P}^{\text{H}} \mathbf{H}^{\text{H}} \mathbf{G} \mathbf{C}_{\text{ss}} \right) - \alpha_q \text{tr} \left(\mathbf{C}_{\text{ss}} \mathbf{G} \mathbf{H} \mathbf{P} \right) + \text{tr} \left(\mathbf{C}_{\text{ss}} \right). \tag{5.44}
\end{aligned}$$

As we can see, all terms of (5.30) and (5.44) match except of the second and third terms. To match the third terms, we choose β as

$$\begin{aligned}\beta &= \sqrt{\frac{\alpha_q \text{tr}(\mathbf{P}\mathbf{C}_{\text{ss}}\mathbf{P}^{\text{H}})}{\sigma_s^2 \text{tr}(\mathbf{G}\mathbf{C}_{\eta\eta}\mathbf{G})}} \\ &\stackrel{(5.9)}{=} \sqrt{\frac{P_{\text{tx}}}{\sigma_s^2 \text{tr}(\mathbf{G}\mathbf{C}_{\eta\eta}\mathbf{G})}}.\end{aligned}\quad (5.45)$$

However, we cannot match the second terms due to the non-linear expressions. Thus, there exists no MSE duality if $\mathbf{d}_{\mathcal{Q}}^{\text{UL}}$ is the quantization noise that is uncorrelated with \mathbf{x}^{UL} and results from quantizing \mathbf{x}^{UL} . The problem arises with the non-diagonal elements of $\mathbf{C}_{\mathbf{d}_{\mathcal{Q}}\mathbf{d}_{\mathcal{Q}}}^{\text{UL}}$ that should be chosen in an appropriate way to match both MSE expressions.

5.4.2 Exact Dual Problem

To find the dual MSE expression, we are changing the quantization function at the receiver to replace it with a different dual non-linear function; that is the distortion term $\mathbf{d}_{\mathcal{Q}}$ has the following covariance matrix

$$\begin{aligned}\mathbf{C}_{\mathbf{d}_{\mathcal{Q}}\mathbf{d}_{\mathcal{Q}}}^{\text{UL}} &= \alpha_q \beta_q \text{diag}(\mathbf{H}^{\text{H}}\mathbf{C}_{\eta\eta}^{-1/2, \text{H}}\mathbf{T}\mathbf{C}_{\text{ss}}^{\text{UL}}\mathbf{T}^{\text{H}}\mathbf{C}_{\eta\eta}^{-1/2}\mathbf{H} + \mathbf{C}_{\eta\eta}^{\text{UL}}) \\ &\quad - \alpha_q^2 \text{nondiag}(\mathbf{H}^{\text{H}}\mathbf{C}_{\eta\eta}^{-1/2, \text{H}}\mathbf{T}\mathbf{C}_{\text{ss}}^{\text{UL}}\mathbf{T}^{\text{H}}\mathbf{C}_{\eta\eta}^{-1/2}\mathbf{H}) \\ &\quad + \frac{Q}{\pi} \sin^2(\psi) \alpha_q \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \mathbf{F}^{\text{H}} (\mathbf{F}\mathbf{F}^{\text{H}})^{-1} \mathbf{T}^{\text{H}}\mathbf{C}_{\eta\eta}^{-1/2}\mathbf{H} \text{diag}(\mathbf{F}^{\text{H}}\mathbf{C}_{\text{ss}}^{\text{UL}}\mathbf{F})^{1/2} \\ &\quad \arcsin\left(\Re\{\text{nondiag}\left(\text{diag}(\mathbf{F}^{\text{H}}\mathbf{C}_{\text{ss}}^{\text{UL}}\mathbf{F})^{-1/2} \mathbf{F}^{\text{H}}\mathbf{C}_{\text{ss}}^{\text{UL}}\mathbf{F} \text{diag}(\mathbf{F}^{\text{H}}\mathbf{C}_{\text{ss}}^{\text{UL}}\mathbf{F})^{-1/2}\right) e^{-j(2\Delta k\psi)}\}\right) \\ &\quad \text{diag}(\mathbf{F}^{\text{H}}\mathbf{C}_{\text{ss}}^{\text{UL}}\mathbf{F})^{1/2} \mathbf{H}^{\text{H}}\mathbf{C}_{\eta\eta}^{-1/2, \text{H}}\mathbf{T} (\mathbf{F}\mathbf{F}^{\text{H}})^{-1} \mathbf{F}.\end{aligned}\quad (5.46)$$

The resulting MSE expression reads as

$$\begin{aligned}\text{MSE}^{\text{UL}} &= \alpha_q \text{tr}(\mathbf{F} \text{diag}(\mathbf{H}^{\text{H}}\mathbf{C}_{\eta\eta}^{-1/2, \text{H}}\mathbf{T}\mathbf{C}_{\text{ss}}^{\text{UL}}\mathbf{T}^{\text{H}}\mathbf{C}_{\eta\eta}^{-1/2}\mathbf{H}) \mathbf{F}^{\text{H}}) + \alpha_q \text{tr}(\mathbf{F}\mathbf{C}_{\eta\eta}^{\text{UL}}\mathbf{F}^{\text{H}}) \\ &\quad + \frac{Q}{\pi} \sin^2(\psi) \alpha_q \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \text{tr}\left(\mathbf{T}^{\text{H}}\mathbf{C}_{\eta\eta}^{-1/2}\mathbf{H} \text{diag}(\mathbf{F}^{\text{H}}\mathbf{C}_{\text{ss}}^{\text{UL}}\mathbf{F})^{1/2}\right. \\ &\quad \left.\arcsin\left(\Re\{\text{nondiag}\left(\text{diag}(\mathbf{F}^{\text{H}}\mathbf{C}_{\text{ss}}^{\text{UL}}\mathbf{F})^{-1/2} \mathbf{F}^{\text{H}}\mathbf{C}_{\text{ss}}^{\text{UL}}\mathbf{F} \text{diag}(\mathbf{F}^{\text{H}}\mathbf{C}_{\text{ss}}^{\text{UL}}\mathbf{F})^{-1/2}\right) e^{-j(2\Delta k\psi)}\}\right)\right) \\ &\quad \text{diag}(\mathbf{F}^{\text{H}}\mathbf{C}_{\text{ss}}^{\text{UL}}\mathbf{F})^{1/2} \mathbf{H}^{\text{H}}\mathbf{C}_{\eta\eta}^{-1/2, \text{H}}\mathbf{T}) \\ &\quad - \alpha_q \text{tr}(\mathbf{F}\mathbf{H}^{\text{H}}\mathbf{C}_{\eta\eta}^{-1/2, \text{H}}\mathbf{T}\mathbf{C}_{\text{ss}}^{\text{UL}}) - \alpha_q \text{tr}(\mathbf{C}_{\text{ss}}^{\text{UL}}\mathbf{T}^{\text{H}}\mathbf{C}_{\eta\eta}^{-1/2}\mathbf{H}\mathbf{F}^{\text{H}}) + \text{tr}(\mathbf{C}_{\text{ss}}^{\text{UL}}).\end{aligned}\quad (5.47)$$

By applying the equalities in (5.40), (5.41), (5.42), (5.43) and (5.45), we obtain

$$\text{MSE}^{\text{UL}} = \text{MSE}^{\text{DL}}.\quad (5.48)$$

The dual uplink system is not simply having the quantization at the receiver in the uplink scenario.

In the dual uplink system, we end up with two filters \mathbf{F} and \mathbf{T} to optimize. To simplify the optimization task, we assume that \mathbf{G} is a scaled identity, i.e. $\mathbf{G} = g\mathbf{I}_M$ and thus from (5.41) and (5.45) we get

$$\mathbf{T} = \sqrt{\frac{P_{\text{tx}}}{\sigma_s^2 \text{tr}(\mathbf{C}_{\eta\eta})}} \mathbf{C}_{\eta\eta}^{1/2, \text{H}}. \quad (5.49)$$

Consequently, the transmit power is the same as in the primal domain since

$$\text{tr}(\mathbf{T}\mathbf{C}_{\text{ss}}^{\text{UL}}\mathbf{T}^{\text{H}}) = P_{\text{tx}}. \quad (5.50)$$

To find the dual filter \mathbf{F} , we should compute the derivative of MSE^{UL} w.r.t. \mathbf{F} . In analogy to the derivations in Section 5.3.1, we get

$$\begin{aligned} \frac{\partial \text{MSE}^{\text{UL}}}{\partial \mathbf{F}} &= \frac{Q}{\pi} \sin^2(\psi) \alpha_q \sum_{\Delta k=0}^{Q/2-1} e^{-j(2\Delta k\psi)} \frac{1}{2} \mathbf{C}_{\text{ss}}^{\text{UL}*} \mathbf{F}^* \left(\text{diag}(\mathbf{F}^{\text{H}} \mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{F})^{-1/2} \left(e^{-j2\Delta k\psi} \boldsymbol{\Omega} + e^{j2\Delta k\psi} \boldsymbol{\Omega}^{\text{T}} \right. \right. \\ &\quad \left. \left. - \text{diag}(\Re\{\boldsymbol{\Sigma} e^{-j2\Delta k\psi}\} \boldsymbol{\Omega}^*) - \text{diag}(\boldsymbol{\Omega}^* \Re\{\boldsymbol{\Sigma} e^{-j2\Delta k\psi}\}) \right) \right. \\ &\quad \left. + \frac{P_{\text{tx}}}{\sigma_s^2 \text{tr}(\mathbf{C}_{\eta\eta})} \text{diag}(\boldsymbol{\Lambda} \mathbf{H}^{\text{T}} \mathbf{C}_{\eta\eta}^{-1/2, \text{T}} \mathbf{C}_{\eta\eta}^{-1/2, *} \mathbf{H}^* \boldsymbol{\Lambda} \arcsin(\Re\{\boldsymbol{\Sigma} e^{-j2\Delta k\psi}\})) \right) \\ &\quad \left. + \frac{P_{\text{tx}}}{\sigma_s^2 \text{tr}(\mathbf{C}_{\eta\eta})} \text{diag}(\arcsin(\Re\{\boldsymbol{\Sigma} e^{-j2\Delta k\psi}\}) \boldsymbol{\Lambda} \mathbf{H}^{\text{T}} \mathbf{C}_{\eta\eta}^{-1/2, \text{T}} \mathbf{C}_{\eta\eta}^{-1/2, *} \mathbf{H}^* \boldsymbol{\Lambda}) \right) \\ &\quad \text{diag}(\mathbf{F}^{\text{H}} \mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{F})^{-1/2} \left. \right) \\ &\quad + \sqrt{\frac{P_{\text{tx}}}{\sigma_s^2 \text{tr}(\mathbf{C}_{\eta\eta})}} \alpha_q \mathbf{F}^* \mathbf{C}_{\eta\eta}^{\text{UL}, \text{T}} - \alpha_q \mathbf{C}_{\text{ss}}^{\text{UL}, \text{T}} \mathbf{C}_{\eta\eta}^{-1/2, *} \mathbf{H}^*, \end{aligned} \quad (5.51)$$

where

$$\boldsymbol{\Sigma} = \text{diag}(\mathbf{F}^{\text{H}} \mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{F})^{-1/2} \mathbf{F}^{\text{H}} \mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{F} \text{diag}(\mathbf{F}^{\text{H}} \mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{F})^{-1/2}, \quad (5.52)$$

$$\boldsymbol{\Omega} = (\boldsymbol{\Lambda} \mathbf{H}^{\text{H}} \mathbf{C}_{\eta\eta}^{-1/2, \text{H}} \mathbf{T} \mathbf{T}^{\text{H}} \mathbf{C}_{\eta\eta}^{-1/2} \mathbf{H} \boldsymbol{\Lambda}) \circ \text{nondiag} \left(\mathbf{1}_{N, N} - \text{nondiag}(\Re\{\boldsymbol{\Sigma}^* e^{-j2\Delta k\psi}\})^{\circ 2} \right)^{\circ -1/2}, \quad (5.53)$$

and

$$\boldsymbol{\Lambda} = \text{diag}(\mathbf{F}^{\text{H}} \mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{F})^{1/2}. \quad (5.54)$$

Unfortunately, we cannot find a closed-form expression for \mathbf{F} and we have to find the optimal filter \mathbf{F} again with iterative algorithms, i.e. gradient descent algorithm as in the primal problem. However, the advantage here is that no power constraint has to be considered. The gradient descent algorithm is given in Algorithm 3. After obtaining the optimal filter \mathbf{F} , the duality scaling factor β in (4.44) can be expressed by

$$\beta = \sqrt{\frac{P_{\text{tx}}}{\alpha_q \text{tr}(\mathbf{F}^{\text{H}} \mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{F})}}. \quad (5.55)$$

Hence, the dual filters \mathbf{P} and \mathbf{G} can be extracted by applying the equalities in (5.40) and (5.41).

Algorithm 3 Gradient Descent Algorithm to obtain the dual WFQ-Price with unequal power allocation.

- 1: Initialization
 $\mathbf{F}_{(0)}$, $\mu = 10$, $\epsilon = 10^{-4}$ and $n = 0$
 - 2: **repeat**
 - 3: $\mathbf{F}_{(n+1)} = \mathbf{F}_{(n)} - \mu \left(\frac{\partial \text{MSE}_{(n)}^{\text{UL}}}{\partial \mathbf{F}} \right)^*$
 - 4: **if** $\text{MSE}_{(n+1)}^{\text{UL}} > \text{MSE}_{(n)}^{\text{UL}}$ **then**
 - 5: $\mu = \mu/2$
 - 6: **else**
 - 7: $n = n + 1$
 - 8: **end if**
 - 9: **until** $\frac{|\text{MSE}_{(n+1)}^{\text{UL}} - \text{MSE}_{(n)}^{\text{UL}}|}{|\text{MSE}_{(n)}^{\text{UL}}|} \leq \epsilon$
-

5.4.3 Approximate Dual Problem

We still aim to get a closed-form expression for the filter \mathbf{F} and hence for the precoder \mathbf{P} . To this end, we assume that the MSE expressions in (5.35) and in (5.30) are equal. This assumption holds true if the Bussgang decomposition is applied with the LCA; that is $\mathbf{C}_{\mathbf{d}_q \mathbf{d}_q} = \alpha_q \beta_q \text{diag} \left(\mathbf{H}^H \mathbf{C}_{\eta\eta}^{-1/2, \text{H}} \mathbf{T} \mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{T}^H \mathbf{C}_{\eta\eta}^{-1/2} \mathbf{H} + \mathbf{C}_{\eta\eta}^{\text{UL}} \right)$. The MSE^{UL} expression is in general expressed as

$$\text{MSE}^{\text{UL}} = \text{tr}(\mathbf{F} \mathbf{C}_{\text{tt}}^{\text{UL}} \mathbf{F}^H) - \text{tr}(\mathbf{F} \mathbf{C}_{\text{ts}}^{\text{UL}}) - \text{tr}(\mathbf{C}_{\text{st}}^{\text{UL}} \mathbf{F}^H) + \text{tr}(\mathbf{C}_{\text{ss}}^{\text{UL}}). \quad (5.56)$$

The optimal filter \mathbf{F} that minimizes (5.56) reads as

$$\mathbf{F} = \mathbf{C}_{\text{st}}^{\text{UL}} (\mathbf{C}_{\text{tt}}^{\text{UL}})^{-1}. \quad (5.57)$$

To compute $\mathbf{C}_{\text{tt}}^{\text{UL}}$ we apply Price's theorem in (3.66). The optimal filter \mathbf{T} is given in (5.49). Hence, the optimal filters in the primal domain are obtained by applying (5.40) and (5.41), where β is given in (4.44). This method was introduced in [69] for the one-bit quantization case, i.e. $Q = 4$.

5.5 Simulation Results

In this section, we compare the performance of the different linear precoding techniques that were introduced above with the ideal WF precoder. We denote the precoders by

- WF: the ideal WF, [70], where no CEQ is applied in the system.
- QWF: the Quantized WF that was introduced in [25] for cartesian DACs and was extended to the case of polar DACs in Section 5.3.2.
- QWF-Apprx_{Dual}: the approximate dual Quantized WF that applies unequal power allocation at the antennas. The resulting receive filter \mathbf{G} is a scaled identity. This precoder is derived in Section 5.4.3.

- QWF-Price_{eqPA}: the Quantized WF based on Price's theorem that applies equal power allocation at the antennas. This precoder is detailed in Algorithm 1. The start value is the WF precoder that is projected to fulfill the equal transmit power constraint.
- QWF-Price_{neqPA}: the Quantized WF based on Price's theorem that applies unequal power allocation at the antennas. This precoder is detailed in Algorithm 2. The start value is the WF precoder.
- QWF-Price_{Dual}: the dual Quantized WF based on Price's theorem that applies unequal power allocation at the antennas. The resulting receive filter \mathbf{G} is a scaled identity. This precoder is detailed in Algorithm 3. The start value is the QWF-Apprx_{Dual} precoder.

To this end, we assume a BS with $N = 64$ or $N = 128$ antennas serving $M = 8$ single-antenna users with 4-PSK or 16-QAM signals, respectively. The numerical results are obtained with Monte Carlo simulations of 100 independent flat-fading channel realizations from the i.i.d. channel model and the mmW sparse channel model described in Section 4.6. For the mmW sparse channel, we assume that $N_{\text{cl}} = 2$ and $N_{\text{ray}} = 10$. The AWGN is also i.i.d. with variance one at each antenna. The performance metric is the uncoded BER averaged over the single-antenna users. For the blind estimation of the coefficients g_m we use a block length of $T = 128$. The numerical results are plotted in Fig. 5.1, Fig. 5.2, Fig. 5.3 and Fig. 5.4.

First, it can be deduced that all proposed linear precoders perform almost the same for i.i.d. channels. There is a moderate gain compared to the QWF precoder. The loss compared to the ideal WF reduces with increased quantization resolution.

Second, the proposed precoders perform differently for mmW sparse channels. We can see that the QWF-Price_{neqPA} precoder performs the best followed by QWF-Price_{eqPA}, QWF-Price_{Dual}, QWF-Apprx_{Dual} and last QWF. This behavior is expected since the QWF-Price_{neqPA} precoder is computed in the primal domain without any approximation of Price's theorem. Additionally, having unequal power allocation at the antennas offers more degrees of freedom for the precoder design, which explains the performance loss of QWF-Price_{eqPA} compared to QWF-Price_{neqPA}. The performance loss of QWF-Price_{Dual} can be explained by the assumption of \mathbf{G} being a scaled identity, which again reduces the degrees of freedom in the precoder design. Although the QWF-Apprx_{Dual} precoder is more efficient in terms of convergence time and computational complexity, it performs worse than the linear precoders based on iterative algorithms. This is due to the approximate duality that we applied for this precoder design and the assumption on \mathbf{G} being a scaled identity. Last but not least, the QWF precoder performs the worst due to the LCA and the assumption on \mathbf{G} being a scaled identity.

Third, doubling the number of antennas with the use of the QWF-Price_{neqPA} precoder leads to almost the same BER performance as in the ideal unquantized case for $Q = 4$ and with 4-PSK signaling. With increased resolution, e.g. $Q = 8$, less than the double number of antennas is required in the quantized case to obtain the same performance as in the ideal case. For 16-QAM signaling the latter statement holds for low values of the transmit power.

In summary, the performance improvement with the iterative and more accurate linear precoding algorithms that take into account the quantization distortions is very moderate compared to the QWF that is based on approximate statistics. However, this comparison is restricted to flat-fading channels. So how does the performance comparison look like in the case of frequency-selective channels?

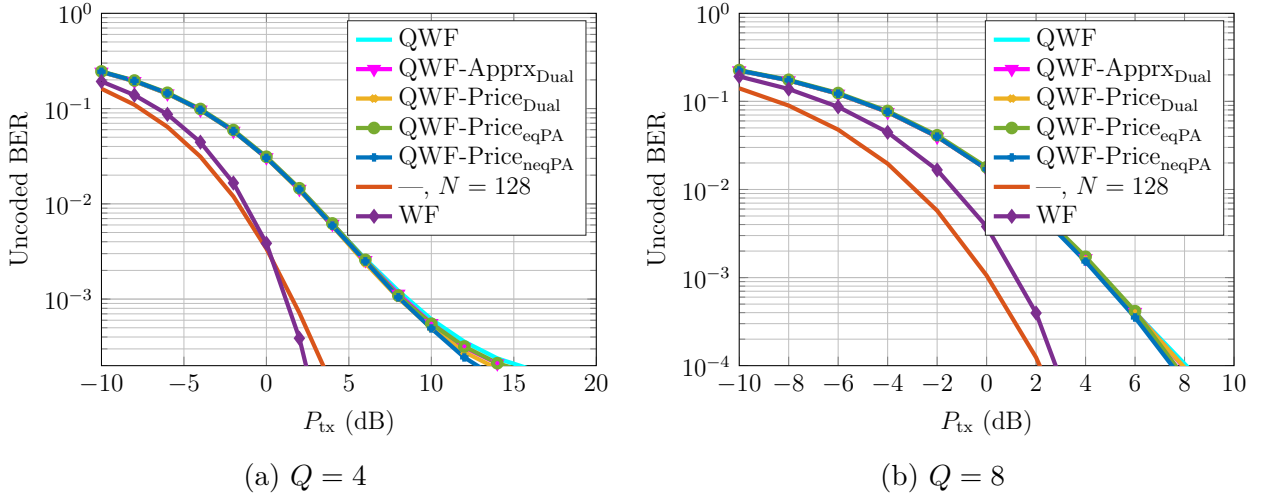


Fig. 5.1: Comparison between different linear precoders for 4-PSK signaling: $N = 64$, $M = 8$ with the i.i.d. channel.

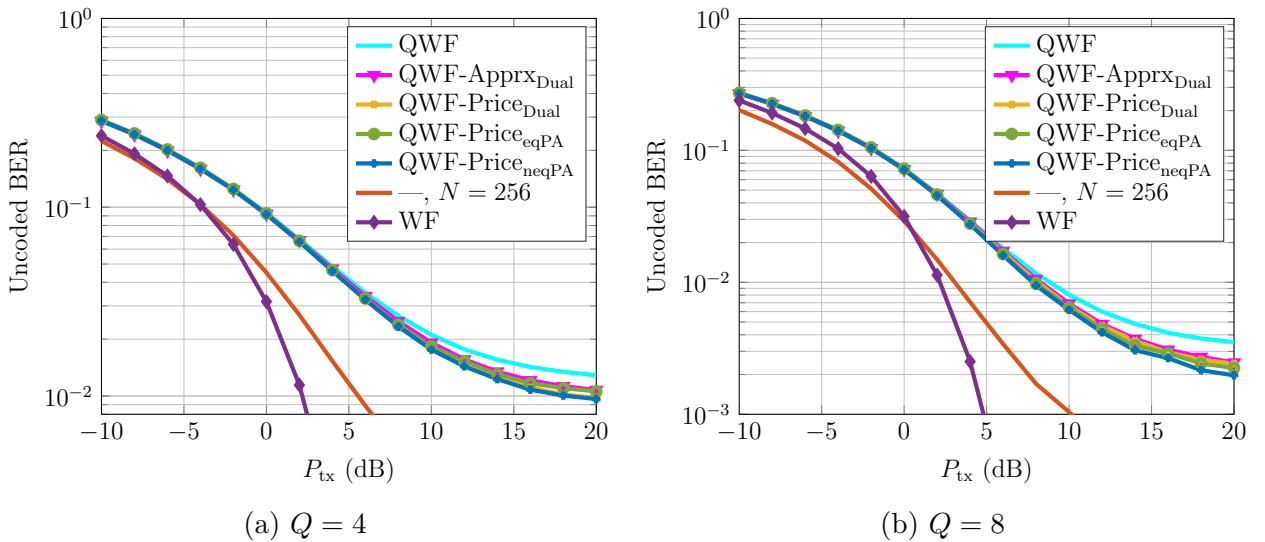


Fig. 5.2: Comparison between different linear precoders for 16-QAM signaling: $N = 128$, $M = 8$ with the i.i.d. channel.

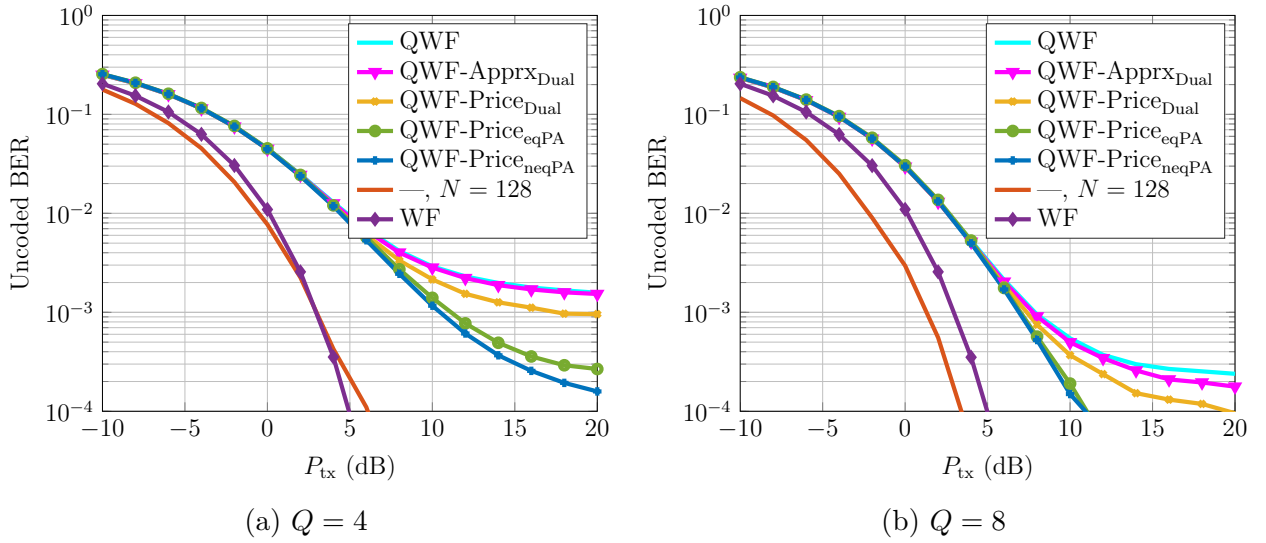


Fig. 5.3: Comparison between different linear precoders for 4-PSK signaling: $N = 64$, $M = 8$ with the mmW sparse channel.

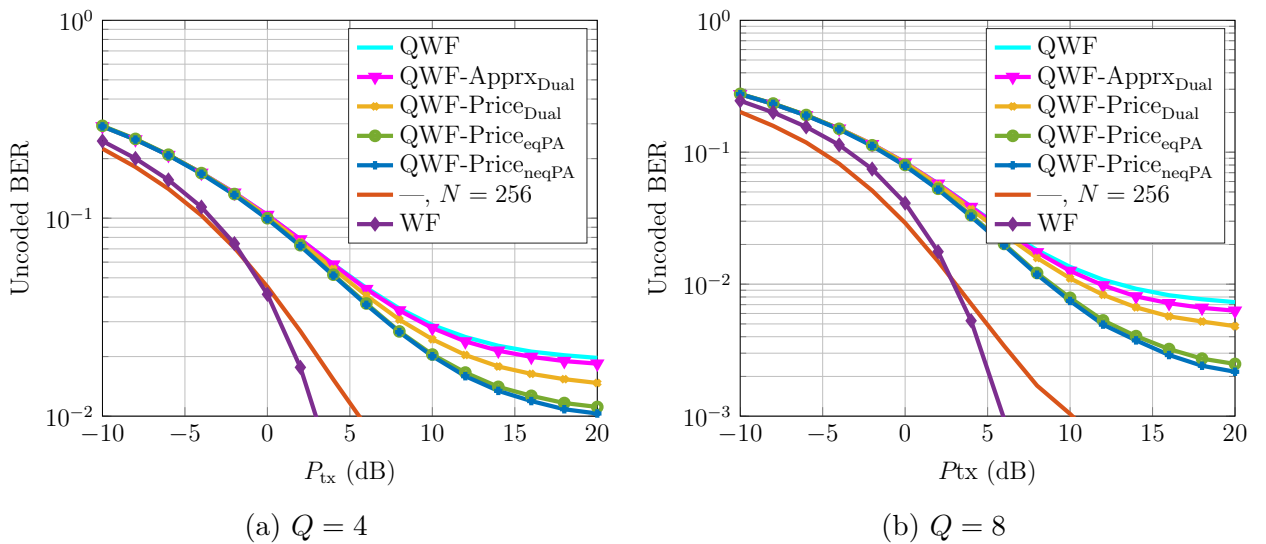


Fig. 5.4: Comparison between different linear precoders for 16-QAM signaling: $N = 128$, $M = 8$ with the mmW sparse channel.

6. Frequency-Selective Channels

6.1 Input-Output Relationship

In this section, we consider linear precoding for frequency-selective channels; that is $L > 1$. This implies that more than one precoder tap is required to mitigate the Inter-Symbol Interference (ISI); that is $L_p > 1$. Thus, the precoded signal vector is given by

$$\begin{aligned}
 \mathbf{x}[t] &= \sum_{\ell'=0}^{L_p-1} \mathbf{P}_{\ell'} \mathbf{s}[t - \ell'] \\
 &= \sum_{\ell'=0}^{L_p-1} \sum_{n=1}^N \sum_{m=1}^M \mathbf{e}_n \mathbf{e}_n^T \mathbf{P}_{\ell'} \mathbf{e}_m \mathbf{e}_m^T \mathbf{s}[t - \ell'] \\
 &= \sum_{\ell'=0}^{L_p-1} \sum_{n=1}^N \sum_{m=1}^M \mathbf{e}_n \mathbf{e}_m^T \mathbf{P}_{\ell'}^T \mathbf{e}_n \mathbf{e}_m^T \mathbf{s}[t - \ell'] \\
 &= \sum_{n=1}^N \mathbf{e}_n \mathbf{p}^T (\mathbf{I}_{ML_p} \otimes \mathbf{e}_n) \begin{bmatrix} \mathbf{s}[t] \\ \mathbf{s}[t-1] \\ \vdots \\ \mathbf{s}[t-L_p+1] \end{bmatrix}, \tag{6.1}
 \end{aligned}$$

where

$$\mathbf{p}^T = [\mathbf{e}_1^T \mathbf{P}_1^T \quad \cdots \quad \mathbf{e}_M^T \mathbf{P}_1^T \quad \cdots \quad \mathbf{e}_1^T \mathbf{P}_{L_p-1}^T \quad \cdots \quad \mathbf{e}_M^T \mathbf{P}_{L_p-1}^T]. \tag{6.2}$$

For the vector \mathbf{x} defined in (4.11) we get

$$\mathbf{x} = \sum_{i=1}^{NL} \mathbf{e}_i \mathbf{p}^T \mathbf{\Gamma}_i \begin{bmatrix} \mathbf{s}[t] \\ \vdots \\ \mathbf{s}[t-L-L_p+1] \end{bmatrix} \tag{6.3}$$

$$= (\mathbf{I}_{NL} \otimes \mathbf{p}^T) \mathbf{\Gamma} \mathbf{s}, \tag{6.4}$$

where

$$\mathbf{\Gamma}_i = \mathbf{S} \left(\left[\frac{i}{N} \right] - 1, L_p, L - 1 \right) \otimes \mathbf{I}_M \otimes \mathbf{e}_{\text{mod}(i-1, N)+1}, \tag{6.5}$$

$$\mathbf{S}(i, j, \ell) = [\mathbf{0}_{j,i}, \mathbf{I}_j, \mathbf{0}_{j,\ell-i}] \tag{6.6}$$

$$\mathbf{\Gamma} = \sum_{i=1}^{NL} \mathbf{e}_i \otimes \mathbf{\Gamma}_i, \tag{6.7}$$

and

$$\mathbf{s} = [\mathbf{s}[t]^T \quad \mathbf{s}[t-1]^T \quad \cdots \quad \mathbf{s}[t-L-L_p+1]^T]^T. \quad (6.8)$$

Thus, we recall the input-output relationship in (4.8) and use the identity in (6.4) to get

$$\begin{aligned} \mathbf{u} &= \mathbf{G}(\mathbf{H}\mathbf{t} + \boldsymbol{\eta}) \\ &= \mathbf{G}(\mathbf{H}\mathcal{Q}_{\text{CE}}((\mathbf{I}_{NL} \otimes \mathbf{p}^T)\boldsymbol{\Gamma}\mathbf{s}) + \boldsymbol{\eta}). \end{aligned} \quad (6.9)$$

6.2 Optimization Problem

The precoding task consists of finding the optimal linear precoding vector \mathbf{p}_{opt} and the diagonal positive real-valued receive filter \mathbf{G}_{opt} that minimize the MSE between the desired signal before a time delay τ and the received signal, i.e. $\mathbf{s}[t-\tau]$ and $\mathbf{u}[t]$, under the transmit power constraint given in (4.20). The MMSE optimization problem is expressed as

$$\begin{aligned} \{\mathbf{p}_{\text{opt}}, \mathbf{G}_{\text{opt}}, \tau_{\text{opt}}\} &= \arg \min_{\mathbf{p}, \mathbf{G}, \tau} \mathbb{E} [\|\mathbf{u}[t] - \mathbf{s}[t-\tau]\|_2^2] \quad \text{s.t.} \quad \alpha_q \text{tr}(\mathbf{C}_{\mathbf{xx}}) \leq P_{\text{tx}} \\ &= \arg \min_{\mathbf{p}, \mathbf{G}, \tau} \mathbb{E} [\|\mathbf{u} - (\mathbf{e}_{\tau+1}^T \otimes \mathbf{I}_M) \mathbf{s}\|_2^2] \quad \text{s.t.} \quad \alpha_q \text{tr}(\mathbf{C}_{\mathbf{xx}}) \leq P_{\text{tx}}. \end{aligned} \quad (6.10)$$

6.2.1 MSE

In general, the MSE expression is given by

$$\begin{aligned} \text{MSE} &= \text{tr}(\mathbf{G}\mathbf{H}\mathbf{C}_{\text{tt}}\mathbf{H}^H\mathbf{G}) + \text{tr}(\mathbf{G}\mathbf{C}_{\eta\eta}\mathbf{G}) - \text{tr}(\mathbf{G}\mathbf{H}\mathbf{C}_{\text{ts}}(\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M)) \\ &\quad - \text{tr}((\mathbf{e}_{\tau+1}^T \otimes \mathbf{I}_M)\mathbf{C}_{\text{st}}\mathbf{H}^H\mathbf{G}) + \text{tr}((\mathbf{e}_{\tau+1}^T \otimes \mathbf{I}_M)\mathbf{C}_{\text{ss}}(\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M)). \end{aligned} \quad (6.11)$$

In massive MIMO systems, the entries in \mathbf{x} are approximately Gaussian distributed due to the central limit theorem [68]. Hence, Price's theorem can be applied to compute the covariance matrices \mathbf{C}_{tt} and \mathbf{C}_{tx} as given in (3.66) and (3.65). Consequently, the MSE expression calculates to

$$\begin{aligned} \text{MSE} &= \frac{Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \text{tr}(\arcsin(\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\}) \boldsymbol{\Xi}\mathbf{H}^H\mathbf{G}^2\mathbf{H}\boldsymbol{\Xi}) \\ &\quad + \text{tr}(\mathbf{G}\mathbf{C}_{\eta\eta}\mathbf{G}) - \alpha_q \text{tr}(\mathbf{G}\mathbf{H}\mathbf{C}_{\text{xs}}(\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M)) - \alpha_q \text{tr}((\mathbf{e}_{\tau+1}^T \otimes \mathbf{I}_M)\mathbf{C}_{\text{sx}}\mathbf{H}^H\mathbf{G}) \\ &\quad + \text{tr}((\mathbf{e}_{\tau+1}^T \otimes \mathbf{I}_M)\mathbf{C}_{\text{ss}}(\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M)). \end{aligned} \quad (6.12)$$

According to (6.4), the covariance matrix $\mathbf{C}_{\mathbf{xx}}$ is expressed in terms of the precoder \mathbf{p} as

$$\mathbf{C}_{\mathbf{xx}} = \sum_{i=1}^{NL} \sum_{j=1}^{NL} \mathbf{p}^T \boldsymbol{\Gamma}_i \mathbf{C}_{\text{ss}} \boldsymbol{\Gamma}_j^T \mathbf{p}^* \mathbf{e}_i \mathbf{e}_j^T \quad (6.13)$$

$$= (\mathbf{I}_{NL} \otimes \mathbf{p}^T) \boldsymbol{\Gamma} \mathbf{C}_{\text{ss}} \boldsymbol{\Gamma}^T (\mathbf{I}_{NL} \otimes \mathbf{p}^*). \quad (6.14)$$

Therefore, we obtain the following expressions

$$\begin{aligned} \mathbf{R}_{\mathbf{xx}} &= \text{diag}((\mathbf{I}_{NL} \otimes \mathbf{p}^T) \boldsymbol{\Gamma} \mathbf{C}_{\text{ss}} \boldsymbol{\Gamma}^T (\mathbf{I}_{NL} \otimes \mathbf{p}^*))^{-1/2} \\ &\quad (\mathbf{I}_{NL} \otimes \mathbf{p}^T) \boldsymbol{\Gamma} \mathbf{C}_{\text{ss}} \boldsymbol{\Gamma}^T (\mathbf{I}_{NL} \otimes \mathbf{p}^*) \text{diag}((\mathbf{I}_{NL} \otimes \mathbf{p}^T) \boldsymbol{\Gamma} \mathbf{C}_{\text{ss}} \boldsymbol{\Gamma}^T (\mathbf{I}_{NL} \otimes \mathbf{p}^*))^{-1/2}, \end{aligned} \quad (6.15)$$

and

$$\boldsymbol{\Xi} = \text{diag}((\mathbf{I}_{NL} \otimes \mathbf{p}^T) \boldsymbol{\Gamma} \mathbf{C}_{\text{ss}} \boldsymbol{\Gamma}^T (\mathbf{I}_{NL} \otimes \mathbf{p}^*))^{1/2}. \quad (6.16)$$

6.2.2 Transmit Power Constraint

After plugging (6.14) in the transmit power constraint in (4.20) and for maximal exploitation of the available power, we get

$$\alpha_q \text{tr} \left((\mathbf{I}_{NL} \otimes \mathbf{p}^T) \mathbf{\Gamma} \mathbf{C}_{\text{ss}} \mathbf{\Gamma}^T (\mathbf{I}_{NL} \otimes \mathbf{p}^*) \right) = LP_{\text{tx}}, \quad (6.17)$$

which simplifies to

$$\alpha_q \mathbf{p}^T \sum_{i=1}^{NL} \mathbf{\Gamma}_i \mathbf{C}_{\text{ss}} \mathbf{\Gamma}_i^T \mathbf{p}^* = LP_{\text{tx}}. \quad (6.18)$$

Based on the assumption that $\mathbf{C}_{\text{ss}} = \sigma_s^2 \mathbf{I}_{M(L+L_p-1)}$ and using the identity

$$\mathbf{S}(i, j, \ell) \mathbf{S}(i, j, \ell)^T = \mathbf{I}_j \quad (6.19)$$

that leads to $\mathbf{\Gamma}_i \mathbf{C}_{\text{ss}} \mathbf{\Gamma}_i^T = \mathbf{\Gamma}_{kN+i} \mathbf{C}_{\text{ss}} \mathbf{\Gamma}_{kN+i}^T$, $k \in \mathbb{Z}^+$, we can further simplify the transmit power constraint and obtain

$$\alpha_q \mathbf{p}^T \sum_{i=1}^N \mathbf{\Gamma}_i \mathbf{C}_{\text{ss}} \mathbf{\Gamma}_i^T \mathbf{p}^* = P_{\text{tx}}. \quad (6.20)$$

6.2.2.1 Equal Power Allocation

For equal power allocation, it must holds that

$$\mathbf{p}^T \mathbf{\Gamma}_i \mathbf{C}_{\text{ss}} \mathbf{\Gamma}_i^T \mathbf{p}^* = \frac{P_{\text{tx}}}{\alpha_q N}, \quad i = 1, \dots, N. \quad (6.21)$$

According to the definition of $\mathbf{\Gamma}_i$ in (6.5) and the identity in (6.19), we get

$$\mathbf{\Gamma}_i \mathbf{C}_{\text{ss}} \mathbf{\Gamma}_i^T = \sigma_s^2 (\mathbf{I}_{L_p M} \otimes (\mathbf{e}_i \mathbf{e}_i^T)). \quad (6.22)$$

Thus, the transmit power constraint is expressed as

$$\mathbf{p}^T (\mathbf{I}_{L_p M} \otimes (\mathbf{e}_i \mathbf{e}_i^T)) \mathbf{p}^* = \frac{P_{\text{tx}}}{\sigma_s^2 \alpha_q N}, \quad i = 1, \dots, N. \quad (6.23)$$

6.2.2.2 Unequal Power Allocation

For unequal power allocation, the transmit power constraint is given in (6.20). From (6.22), we can conclude that

$$\sum_{i=1}^N \mathbf{\Gamma}_i \mathbf{C}_{\text{ss}} \mathbf{\Gamma}_i^T = \sigma_s^2 \mathbf{I}_{MNL_p}. \quad (6.24)$$

Hence, the transmit power constraint simplifies to

$$\|\mathbf{p}\|_2^2 \leq \frac{P_{\text{tx}}}{\sigma_s^2 \alpha_q}. \quad (6.25)$$

6.3 Precoder Designs in the Primal Domain

6.3.1 Precoder Design Based on Gradient Projection Algorithm

Since the MSE expression in (6.12) is highly non-linear in $\mathbf{C}_{\mathbf{xx}} = (\mathbf{I}_{NL} \otimes \mathbf{p}^T) \mathbf{\Gamma} \mathbf{C}_{\mathbf{ss}} \mathbf{\Gamma}^T (\mathbf{I}_{NL} \otimes \mathbf{p}^*)$ and thus in \mathbf{p} , we cannot find a closed-form expression for \mathbf{p}_{opt} . Therefore, we use the Gradient Projection algorithm. To this end, we have to compute the derivatives of MSE w.r.t. \mathbf{p} and \mathbf{G} .

$$\begin{aligned} \frac{\partial \text{MSE}}{\partial \mathbf{p}} &= \frac{Q}{\pi} \sin^2(\psi) \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \frac{\partial \text{tr}(\arcsin(\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\}) \mathbf{\Xi} \mathbf{H}^H \mathbf{G}^2 \mathbf{H} \mathbf{\Xi})}{\partial \mathbf{p}} \\ &\quad - \frac{Q^2}{4\pi} \sin^2(\psi) \frac{\partial \text{tr}((\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M) \mathbf{G} \mathbf{H} \mathbf{C}_{\mathbf{xs}})}{\partial \mathbf{p}}, \end{aligned} \quad (6.26)$$

where

$$\begin{aligned} &\frac{\partial \text{tr}(\arcsin(\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\}) \mathbf{\Xi} \mathbf{H}^H \mathbf{G}^2 \mathbf{H} \mathbf{\Xi})}{\partial \mathbf{p}} \\ &= \frac{1}{2} \sum_{i=1}^{NL} \sum_{j=1}^{NL} \mathbf{e}_i^T \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} (\mathbf{\Omega}^T e^{-j2\Delta k\psi} + \mathbf{\Omega} e^{j2\Delta k\psi}) \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \mathbf{e}_j \mathbf{\Gamma}_i \mathbf{C}_{\mathbf{ss}} \mathbf{\Gamma}_j^T \mathbf{p}^* \\ &+ \frac{1}{2} \sum_{i=1}^{NL} \mathbf{e}_i^T \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \mathbf{\Xi} \mathbf{H}^H \mathbf{G}^2 \mathbf{H} \mathbf{\Xi} \arcsin(\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\}) \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \mathbf{e}_i \mathbf{\Gamma}_i \mathbf{C}_{\mathbf{ss}} \mathbf{\Gamma}_i^T \mathbf{p}^* \\ &+ \frac{1}{2} \sum_{i=1}^{NL} \mathbf{e}_i^T \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \arcsin(\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\}) \mathbf{\Xi} \mathbf{H}^H \mathbf{G}^2 \mathbf{H} \mathbf{\Xi} \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \mathbf{e}_i \mathbf{\Gamma}_i \mathbf{C}_{\mathbf{ss}} \mathbf{\Gamma}_i^T \mathbf{p}^* \\ &- \frac{1}{2} \sum_{i=1}^{NL} \mathbf{e}_i^T \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} (\Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\} \mathbf{\Omega} + \mathbf{\Omega} \Re\{\mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi}\}) \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \mathbf{e}_i \mathbf{\Gamma}_i \mathbf{C}_{\mathbf{ss}} \mathbf{\Gamma}_i^T \mathbf{p}^*, \end{aligned} \quad (6.27)$$

$\mathbf{\Omega}$ is defined in (5.16) and

$$\frac{\partial \text{tr}((\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M) \mathbf{G} \mathbf{H} \mathbf{C}_{\mathbf{xs}})}{\partial \mathbf{p}} = \sum_{i=1}^{NL} \mathbf{\Gamma}_i \mathbf{C}_{\mathbf{ss}} (\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M) \mathbf{G} \mathbf{H} \mathbf{e}_i. \quad (6.28)$$

The derivative of the MSE expression in (6.12) w.r.t. \mathbf{G} is expressed as

$$\frac{\partial \text{MSE}}{\partial \mathbf{G}} = 2 \text{diag}(\mathbf{H} \mathbf{C}_{\mathbf{tt}} \mathbf{H}^H + \mathbf{C}_{\eta\eta}) - \text{diag}(\mathbf{H} \mathbf{C}_{\mathbf{ts}} (\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M)) - \text{diag}(\mathbf{H}^* \mathbf{C}_{\mathbf{ts}}^* (\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M)). \quad (6.29)$$

Thus, the optimal filter \mathbf{G}_{opt} is obtained by setting (6.29) equal to zero and is given by

$$\mathbf{G}_{\text{opt}} = \mathbf{g}_{\text{opt}}(\mathbf{p}, \tau) = \left| \text{diag}(\mathbf{H} \mathbf{C}_{\mathbf{tt}} \mathbf{H}^H + \mathbf{C}_{\eta\eta})^{-1} \text{diag}(\Re\{\mathbf{H} \mathbf{C}_{\mathbf{ts}} (\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M)\}) \right|, \quad (6.30)$$

where the $|\bullet|$ operator is applied element-wise to the matrix entries. The Gradient Projection algorithms for equal and unequal power allocation are given in Algorithm 4 and Algorithm 5, respectively. In each iteration, the MSE expression is computed for different values of τ . The value of τ that leads to the minimal MSE is the optimal τ .

Algorithm 4 Gradient Projection Algorithm to obtain the FIR-QWF-Price precoder with equal power allocation.

- 1: Initialization
 $\mathbf{p}_{(0)}, \mathbf{G}_{(0)} = \mathbf{g}_{\text{opt}}(\mathbf{p}_{(0)}), \mu = 10$ and $n = 0$
 - 2: **repeat**
 - 3: $\mathbf{p}_{(n+1)} = \mathbf{p}_{(n)} - \mu \left(\frac{\partial \text{MSE}_{(n)}}{\partial \mathbf{p}} \right)^*$
 - 4: $\mathbf{p}_{(n+1)} = \sqrt{\frac{P_{\text{tx}}}{\alpha_q N}} \text{diag} \left((\mathbf{I}_{NML_p} \otimes \mathbf{p}^T) \mathbf{\Gamma}' \mathbf{C}_{\text{ss}} \mathbf{\Gamma}'^T (\mathbf{I}_{NML_p} \otimes \mathbf{p}^*) \right)^{-1/2} \mathbf{p}_{(n+1)}$,
 where $\mathbf{\Gamma}' = \sum_{i=1}^{NML_p} \mathbf{e}_i \otimes \mathbf{\Gamma}_i$ {Equal power allocation constraint}
 - 5: $\tau = \arg \min_{\tau} \text{MSE}(\mathbf{p}_{(n+1)}, \tau)$
 - 6: $\mathbf{G}_{(n+1)} = \mathbf{g}_{\text{opt}}(\mathbf{p}_{(n+1)}, \tau)$
 - 7: **if** $\text{MSE}_{(n+1)} > \text{MSE}_{(n)}$ **then**
 - 8: $\mu = \mu/2$
 - 9: **else**
 - 10: $n = n + 1$
 - 11: **end if**
 - 12: **until** $\frac{|\text{MSE}_{(n+1)} - \text{MSE}_{(n)}|}{\text{MSE}_{(n)}} \leq \epsilon$
-

Algorithm 5 Gradient Projection Algorithm to obtain the FIR-QWF-Price precoder with unequal power allocation.

- 1: Initialization
 $\mathbf{p}_{(0)}, \mathbf{G}_{(0)} = \mathbf{g}_{\text{opt}}(\mathbf{p}_{(0)}), \mu = 10$ and $n = 0$
 - 2: **repeat**
 - 3: $\mathbf{p}_{(n+1)} = \mathbf{p}_{(n)} - \mu \left(\frac{\partial \text{MSE}_{(n)}}{\partial \mathbf{p}} \right)^*$
 - 4: $\mathbf{p}_{(n+1)} = \sqrt{\frac{P_{\text{tx}}}{\sigma_s^2 \alpha_q}} \|\mathbf{p}_{(n+1)}\|_2^{-1} \mathbf{p}_{(n+1)}$ {Unequal power allocation constraint}
 - 5: $\tau = \arg \min_{\tau} \text{MSE}(\mathbf{p}_{(n+1)}, \tau)$
 - 6: $\mathbf{G}_{(n+1)} = \mathbf{g}_{\text{opt}}(\mathbf{p}_{(n+1)}, \tau)$
 - 7: **if** $\text{MSE}_{(n+1)} > \text{MSE}_{(n)}$ **then**
 - 8: $\mu = \mu/2$
 - 9: **else**
 - 10: $n = n + 1$
 - 11: **end if**
 - 12: **until** $\frac{|\text{MSE}_{(n+1)} - \text{MSE}_{(n)}|}{\text{MSE}_{(n)}} \leq \epsilon$
-

6.3.2 Precoder Design Based on LCA

In this section, we recall the linear precoder design in [25] and apply it to the case of QCE transmit signals for frequency-selective channels. The precoder aims at minimizing the MSE under the constraint of unequal power allocation, where the receive processing is assumed to be a scaled identity matrix, i.e. $\mathbf{G} = g\mathbf{I}_M$. Note that the case of equal power allocation will not be considered, since it must be solved with iterative algorithms and hence no benefits in terms of computational complexity are obtained by using the LCA. However, applying unequal power allocation and assuming a scaled identity processing at the receiver will lead to a closed-form expression of the precoder, as detailed in this section.

Using the LCA and the appropriate expressions in (3.74) and (3.75), we get

$$\begin{aligned} \text{MSE} &= \text{tr} \left(|g|^2 \mathbf{H} (\alpha_q^2 \mathbf{C}_{\mathbf{xx}} + \alpha_q \beta_q \text{diag}(\mathbf{C}_{\mathbf{xx}})) \mathbf{H}^{\text{H}} \right. \\ &\quad - g \mathbf{H} \alpha_q \mathbf{C}_{\mathbf{xs}} (\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M) - g^* \alpha_q (\mathbf{e}_{\tau+1}^{\text{T}} \otimes \mathbf{I}_M) \mathbf{C}_{\mathbf{sx}} \mathbf{H}^{\text{H}} + |g|^2 \mathbf{C}_{\eta\eta} \left. \right) \\ &\quad + \text{tr} \left((\mathbf{e}_{\tau+1}^{\text{T}} \otimes \mathbf{I}_M) \mathbf{C}_{\mathbf{ss}} (\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M) \right). \end{aligned} \quad (6.31)$$

The Lagrangian function is expressed by

$$\begin{aligned} \mathcal{L}(\mathbf{P}, g, \lambda) &= \text{tr} \left(|g|^2 \mathbf{H} (\alpha_q^2 \mathbf{C}_{\mathbf{xx}} + \alpha_q \beta_q \text{diag}(\mathbf{C}_{\mathbf{xx}})) \mathbf{H}^{\text{H}} \right. \\ &\quad - g \mathbf{H} \alpha_q \mathbf{C}_{\mathbf{xs}} (\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M) - g^* \alpha_q (\mathbf{e}_{\tau+1}^{\text{T}} \otimes \mathbf{I}_M) \mathbf{C}_{\mathbf{sx}} \mathbf{H}^{\text{H}} + |g|^2 \mathbf{C}_{\eta\eta} \left. \right) \\ &\quad + \text{tr} \left((\mathbf{e}_{\tau+1}^{\text{T}} \otimes \mathbf{I}_M) \mathbf{C}_{\mathbf{ss}} (\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M) \right) + \lambda \left(\|\mathbf{p}\|_2^2 - \frac{P_{\text{tx}}}{\sigma_s^2 \alpha_q} \right) \\ &\stackrel{(5.6), (6.3)}{=} \text{tr} \left(|g|^2 \mathbf{H} \left(\alpha_q^2 \sum_{i=1}^{NL} \sum_{j=1}^{NL} \mathbf{p}^{\text{T}} \mathbf{\Gamma}_i \mathbf{C}_{\mathbf{ss}} \mathbf{\Gamma}_j^{\text{T}} \mathbf{p}^* \mathbf{e}_i \mathbf{e}_j^{\text{T}} + \alpha_q \beta_q \sum_{i=1}^{NL} \mathbf{p}^{\text{T}} \mathbf{\Gamma}_i \mathbf{C}_{\mathbf{ss}} \mathbf{\Gamma}_i^{\text{T}} \mathbf{p}^* \mathbf{e}_i \mathbf{e}_i^{\text{T}} \right) \mathbf{H}^{\text{H}} \right. \\ &\quad - g \mathbf{H} \alpha_q \sum_{i=1}^{NL} \mathbf{e}_i \mathbf{p}^{\text{T}} \mathbf{\Gamma}_i \mathbf{C}_{\mathbf{ss}} (\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M) - g^* \alpha_q (\mathbf{e}_{\tau+1}^{\text{T}} \otimes \mathbf{I}_M) \mathbf{C}_{\mathbf{ss}} \sum_{i=1}^{NL} \mathbf{\Gamma}_i^{\text{T}} \mathbf{p}^* \mathbf{e}_i^{\text{T}} \mathbf{H}^{\text{H}} \left. \right) \\ &\quad + \text{tr} (|g|^2 \mathbf{C}_{\eta\eta}) + \text{tr} \left((\mathbf{e}_{\tau+1}^{\text{T}} \otimes \mathbf{I}_M) \mathbf{C}_{\mathbf{ss}} (\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M) \right) + \lambda \left(\|\mathbf{p}\|_2^2 - \frac{P_{\text{tx}}}{\sigma_s^2 \alpha_q} \right). \end{aligned} \quad (6.32)$$

The KKT equations are then given by

$$\begin{aligned} \frac{\partial \mathcal{L}(\mathbf{p}, g, \lambda)}{\partial \mathbf{p}} &= |g|^2 \alpha_q^2 \sum_{i=1}^{NL} \sum_{j=1}^{NL} \mathbf{e}_i^{\text{T}} \mathbf{H}^{\text{T}} \mathbf{H}^* \mathbf{e}_j \mathbf{\Gamma}_i \mathbf{C}_{\mathbf{ss}} \mathbf{\Gamma}_j^{\text{T}} \mathbf{p}^* + |g|^2 \alpha_q \beta_q \sum_{i=1}^{NL} \mathbf{e}_i^{\text{T}} \mathbf{H}^{\text{T}} \mathbf{H}^* \mathbf{e}_i \mathbf{\Gamma}_i \mathbf{C}_{\mathbf{ss}} \mathbf{\Gamma}_i^{\text{T}} \mathbf{p}^* \\ &\quad - g \alpha_q \sum_{i=1}^{NL} \mathbf{\Gamma}_i \mathbf{C}_{\mathbf{ss}} (\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M) \mathbf{H} \mathbf{e}_i + \lambda \mathbf{p}^* = \mathbf{0}_{NML_p \times 1}, \end{aligned} \quad (6.33)$$

$$\frac{\partial \mathcal{L}(\mathbf{p}, g, \lambda)}{\partial g} = g^* \text{tr} \left(\mathbf{H} (\alpha_q^2 \mathbf{C}_{\mathbf{xx}} + \alpha_q \beta_q \text{diag}(\mathbf{C}_{\mathbf{xx}})) \mathbf{H}^{\text{H}} + \mathbf{C}_{\eta\eta} \right) - \alpha_q \text{tr} \left(\mathbf{H} \mathbf{C}_{\mathbf{xs}} (\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M) \right) = 0, \quad (6.34)$$

and

$$\frac{\partial \mathcal{L}(\mathbf{p}, g, \lambda)}{\partial \lambda} = \left(\|\mathbf{p}\|_2^2 - \frac{P_{\text{tx}}}{\sigma_s^2 \alpha_q} \right) = 0. \quad (6.35)$$

Multiplying (6.33) by \mathbf{p}^T from the left side and taking the trace leads to

$$\begin{aligned}
 & |g|^2 \alpha_q^2 \operatorname{tr} \left(\sum_{i=1}^{NL} \sum_{j=1}^{NL} \mathbf{e}_i^T \mathbf{H}^T \mathbf{H}^* \mathbf{e}_j \mathbf{p}^T \Gamma_i \mathbf{C}_{\text{ss}} \Gamma_j^T \mathbf{p}^* \right) + |g|^2 \alpha_q \beta_q \operatorname{tr} \left(\sum_{i=1}^{NL} \mathbf{e}_i^T \mathbf{H}^T \mathbf{H}^* \mathbf{e}_i \mathbf{p}^T \Gamma_i \mathbf{C}_{\text{ss}} \Gamma_i^T \mathbf{p}^* \right) \\
 & - g \alpha_q \operatorname{tr} \left(\sum_{i=1}^{NL} \mathbf{p}^T \Gamma_i \mathbf{C}_{\text{ss}} (\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M) \mathbf{H} \mathbf{e}_i \right) - \lambda \|\mathbf{p}\|_2^2 = 0.
 \end{aligned} \tag{6.36}$$

From (6.34), we get

$$\begin{aligned}
 \alpha_q \operatorname{tr} (\mathbf{H} \mathbf{C}_{\text{xs}} (\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M)) &= \alpha_q \operatorname{tr} \left(\sum_{i=1}^{NL} \mathbf{H} \mathbf{e}_i \mathbf{p}^T \Gamma_i \mathbf{C}_{\text{ss}} (\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M) \right) \\
 &= g^* \operatorname{tr} (\mathbf{H} (\alpha_q^2 \mathbf{C}_{\text{xx}} + \alpha_q \beta_q \operatorname{diag}(\mathbf{C}_{\text{xx}}) \mathbf{H}^H + \mathbf{C}_{\eta\eta})),
 \end{aligned} \tag{6.37}$$

which when inserted in (6.36) gives the expression of λ

$$\lambda = \frac{|g|^2 \operatorname{tr}(\mathbf{C}_{\eta\eta})}{\|\mathbf{p}\|_2^2} = \frac{\sigma_s^2 \alpha_q |g|^2 \operatorname{tr}(\mathbf{C}_{\eta\eta})}{P_{\text{tx}}}, \tag{6.38}$$

where we used (6.35). Inserting (6.38) in (6.33) and solving it for \mathbf{p} , we obtain

$$\begin{aligned}
 \mathbf{p} &= \frac{1}{g} \left(\alpha_q \sum_{i=1}^{NL} \sum_{j=1}^{NL} \mathbf{e}_i^T \mathbf{H}^H \mathbf{H} \mathbf{e}_j \Gamma_i \mathbf{C}_{\text{ss}}^* \Gamma_j^T + \beta_q \sum_{i=1}^{NL} \mathbf{e}_i^T \mathbf{H}^H \mathbf{H} \mathbf{e}_i \Gamma_i \mathbf{C}_{\text{ss}}^* \Gamma_i^T + \frac{\sigma_s^2 \operatorname{tr}(\mathbf{C}_{\eta\eta})}{P_{\text{tx}}} \mathbf{I}_{NML_p} \right)^{-1} \\
 & \sum_{i=1}^{NL} \Gamma_i \mathbf{C}_{\text{ss}}^* (\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M) \mathbf{H}^* \mathbf{e}_i.
 \end{aligned} \tag{6.39}$$

The optimal g is found by satisfying (6.35) with \mathbf{p} from (6.39); that is

$$\begin{aligned}
 g &= \sqrt{\frac{\sigma_s^2 \alpha_q}{P_{\text{tx}}}} \left(\sum_{i=1}^{NL} \sum_{j=1}^{NL} \mathbf{e}_j^T \mathbf{H}^T (\mathbf{e}_{\tau+1}^T \otimes \mathbf{I}_M) \mathbf{C}_{\text{ss}}^* \Gamma_j^T \right. \\
 & \left(\alpha_q \sum_{i=1}^{NL} \sum_{j=1}^{NL} \mathbf{e}_i^T \mathbf{H}^H \mathbf{H} \mathbf{e}_j \Gamma_i \mathbf{C}_{\text{ss}}^* \Gamma_j^T + \beta_q \sum_{i=1}^{NL} \mathbf{e}_i^T \mathbf{H}^H \mathbf{H} \mathbf{e}_i \Gamma_i \mathbf{C}_{\text{ss}}^* \Gamma_i^T + \frac{\sigma_s^2 \operatorname{tr}(\mathbf{C}_{\eta\eta})}{P_{\text{tx}}} \mathbf{I}_{NML_p} \right)^{-2} \\
 & \left. \Gamma_i \mathbf{C}_{\text{ss}}^* (\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M) \mathbf{H}^* \mathbf{e}_i \right)^{1/2}.
 \end{aligned} \tag{6.40}$$

Each value of τ determines a different precoder vector \mathbf{p} and a different scalar g , which in turn determine the value of MSE. Therefore, the MSE expression is computed for different values of τ . The value of τ and the corresponding filters \mathbf{p} and g that lead to the minimal MSE are the optimal solutions.

6.4 Dual Optimization Problem

6.4.1 Exact Dual Problem

To find the exact dual optimization problem, we first state the input-output relationship in the uplink system model illustrated in Fig. 4.3; that is

$$\mathbf{u}^{\text{UL}}[t] = \sum_{\ell'=0}^{L_p-1} \mathbf{F}_{\ell'} \left(\alpha_q \left(\sum_{\ell=0}^{L-1} \mathbf{H}_{\ell}^{\text{H}} \mathbf{C}_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-1/2, \text{H}} \mathbf{T} \mathbf{s}^{\text{UL}}[t - \ell' - \ell] + \boldsymbol{\eta}^{\text{UL}}[t - \ell'] \right) + \mathbf{d}_{\mathcal{Q}}^{\text{UL}}[t - \ell'] \right). \quad (6.41)$$

Hence, the uplink MSE expression is given by

$$\begin{aligned} \text{MSE}^{\text{UL}} &= \text{E} \left[\left\| \hat{\mathbf{s}}^{\text{UL}}[t] - \mathbf{s}^{\text{UL}}[t - \tau] \right\|_2^2 \right] \\ &= \alpha_q^2 \text{tr} \left(\sum_{\ell'_1=0}^{L_p-1} \sum_{\ell'_2=0}^{L_p-1} \sum_{\ell_1=0}^{L-1} \sum_{\ell_2=0}^{L-1} \mathbf{F}_{\ell'_1}^{\text{H}} \mathbf{H}_{\ell_1}^{\text{H}} \mathbf{C}_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-1/2, \text{H}} \mathbf{T} \text{E} \left[\mathbf{s}^{\text{UL}}[t - \ell_1 - \ell'_1] \mathbf{s}^{\text{UL}, \text{H}}[t - \ell_2 - \ell'_2] \right] \mathbf{T} \mathbf{C}_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-1/2} \mathbf{H}_{\ell_2} \mathbf{F}_{\ell'_2}^{\text{H}} \right) \\ &+ \alpha_q^2 \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \mathbf{F}_{\ell'} \mathbf{C}_{\boldsymbol{\eta}\boldsymbol{\eta}}^{\text{UL}} \mathbf{F}_{\ell'}^{\text{H}} \right) + \text{tr} \left(\sum_{\ell'_1=0}^{L_p-1} \sum_{\ell'_2=0}^{L_p-1} \mathbf{F}_{\ell'_1} \text{E} \left[\mathbf{d}_{\mathcal{Q}}^{\text{UL}}[t - \ell'_1] \mathbf{d}_{\mathcal{Q}}^{\text{UL}, \text{H}}[t - \ell'_2] \right] \mathbf{F}_{\ell'_2}^{\text{H}} \right) \\ &- \alpha_q \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \sum_{\ell=0}^{L-1} \mathbf{F}_{\ell'} \mathbf{H}_{\ell}^{\text{H}} \mathbf{C}_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-1/2, \text{H}} \mathbf{T} \text{E} \left[\mathbf{s}^{\text{UL}}[t - \ell - \ell'] \mathbf{s}^{\text{UL}, \text{H}}[t - \tau] \right] \right) \\ &- \alpha_q \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \sum_{\ell=0}^{L-1} \text{E} \left[\mathbf{s}^{\text{UL}}[t - \tau] \mathbf{s}^{\text{UL}, \text{H}}[t - \ell - \ell'] \right] \mathbf{T} \mathbf{C}_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-1/2} \mathbf{H}_{\ell} \mathbf{F}_{\ell'}^{\text{H}} \right) \\ &+ \text{tr} \left(\mathbf{C}_{\text{ss}}^{\text{UL}} \right). \end{aligned} \quad (6.42)$$

Since it holds that $\text{E} \left[\mathbf{s}[t_1] \mathbf{s}^{\text{H}}[t_2] \right] = \mathbf{0}_{M, M}$, if $t_1 \neq t_2$, the uplink MSE expression recalculates to

$$\begin{aligned} \text{MSE}^{\text{UL}} &= \alpha_q^2 \text{tr} \left(\sum_{\ell'_1=0}^{L_p-1} \sum_{\ell'_2=0}^{L_p-1} \sum_{\Delta\ell=\ell'_1}^{L-1+\ell'_1} \mathbf{F}_{\ell'_1}^{\text{H}} \mathbf{H}_{\Delta\ell-\ell'_1}^{\text{H}} \mathbf{C}_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-1/2, \text{H}} \mathbf{T} \mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{T} \mathbf{C}_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-1/2} \mathbf{H}_{\Delta\ell-\ell'_2} \mathbf{F}_{\ell'_2}^{\text{H}} \right) \\ &+ \alpha_q^2 \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \mathbf{F}_{\ell'} \mathbf{C}_{\boldsymbol{\eta}\boldsymbol{\eta}}^{\text{UL}} \mathbf{F}_{\ell'}^{\text{H}} \right) + \text{tr} \left(\sum_{\ell'_1=0}^{L_p-1} \sum_{\ell'_2=0}^{L_p-1} \mathbf{F}_{\ell'_1} \text{E} \left[\mathbf{d}_{\mathcal{Q}}^{\text{UL}}[t - \ell'_1] \mathbf{d}_{\mathcal{Q}}^{\text{UL}, \text{H}}[t - \ell'_2] \right] \mathbf{F}_{\ell'_2}^{\text{H}} \right) \\ &- \alpha_q \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \mathbf{F}_{\ell'} \mathbf{H}_{\tau-\ell'}^{\text{H}} \mathbf{C}_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-1/2, \text{H}} \mathbf{T} \mathbf{C}_{\text{ss}}^{\text{UL}} \right) - \alpha_q \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{T} \mathbf{C}_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-1/2} \mathbf{H}_{\tau-\ell'} \mathbf{F}_{\ell'}^{\text{H}} \right) \\ &+ \text{tr} \left(\mathbf{C}_{\text{ss}}^{\text{UL}} \right). \end{aligned} \quad (6.43)$$

For comparison, we need to express the downlink MSE. Therefore, we state the downlink input-output relationship as

$$\mathbf{u}[t] = \sum_{\ell=0}^{L-1} \mathbf{G} \mathbf{H}_{\ell} \left(\alpha_q \sum_{\ell'=0}^{L_p-1} \mathbf{P}_{\ell'} \mathbf{s}[t - \ell - \ell'] + \mathbf{d}_{\mathcal{Q}}[t - \ell] \right) + \mathbf{G} \boldsymbol{\eta}[t]. \quad (6.44)$$

Thus, the downlink MSE expression reads as

$$\begin{aligned}
\text{MSE}^{\text{DL}} &= \text{E} \left[\|\hat{\mathbf{s}}[t] - \mathbf{s}[t - \tau]\|_2^2 \right] \\
&= \alpha_q^2 \text{tr} \left(\sum_{\ell'_1=0}^{L_p-1} \sum_{\ell'_2=0}^{L_p-1} \sum_{\Delta\ell=\ell'_1}^{L-1+\ell'_1} \mathbf{G} \mathbf{H}_{\Delta\ell-\ell'_1} \mathbf{P}_{\ell'_1} \mathbf{C}_{\text{ss}} \mathbf{P}_{\ell'_2}^{\text{H}} \mathbf{H}_{\Delta\ell-\ell'_2}^{\text{H}} \mathbf{G} \right) \\
&\quad + \text{tr} \left(\sum_{\ell_1=0}^{L-1} \sum_{\ell_2=0}^{L-1} \mathbf{G} \mathbf{H}_{\ell_1} \text{E} \left[\mathbf{d}_{\mathcal{Q}}[t - \ell_1] \mathbf{d}_{\mathcal{Q}}^{\text{H}}[t - \ell_2] \right] \mathbf{H}_{\ell_2}^{\text{H}} \mathbf{G} \right) + \text{tr}(\mathbf{G} \mathbf{C}_{\eta\eta} \mathbf{G}) \\
&\quad - \alpha_q \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \mathbf{G} \mathbf{H}_{\tau-\ell'} \mathbf{P}_{\ell'} \mathbf{C}_{\text{ss}} \right) - \alpha_q \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \mathbf{C}_{\text{ss}} \mathbf{P}_{\ell'}^{\text{H}} \mathbf{H}_{\tau-\ell'}^{\text{H}} \mathbf{G} \right) + \text{tr}(\mathbf{C}_{\text{ss}}). \quad (6.45)
\end{aligned}$$

Section 5.4.1 proves that for flat-fading channels there is no MSE duality between the downlink system with a CEQ at the transmitter and the uplink system with a CEQ at the receiver. According to Section 5.4.3, there exists only an approximate duality, when the LCA instead of Price's theorem is applied to compute the MSE expressions in the downlink and uplink scenarios. In Section 5.4.2, we modified the covariance matrix $\mathbf{C}_{\mathbf{d}_{\mathcal{Q}}^{\text{UL}} \mathbf{d}_{\mathcal{Q}}^{\text{UL}}}$ to achieve the exact MSE duality. Analogously, we draw the same conclusions for frequency-selective channels. To this end, we have to find the matching parts and the non-matching parts between the downlink and uplink MSE expressions.

First, we split the terms $\text{E} \left[\mathbf{d}_{\mathcal{Q}}^{\text{UL}}[t - \ell'_1] \mathbf{d}_{\mathcal{Q}}^{\text{UL,H}}[t - \ell'_2] \right]$ and $\text{E} \left[\mathbf{d}_{\mathcal{Q}}[t - \ell'_1] \mathbf{d}_{\mathcal{Q}}^{\text{H}}[t - \ell'_2] \right]$ in the uplink and the downlink systems into three terms

- $\text{diag} \left(\text{E} \left[\mathbf{d}_{\mathcal{Q}}^{\text{UL}}[t - \ell'_1] \mathbf{d}_{\mathcal{Q}}^{\text{UL,H}}[t - \ell'_2] \right] \right)$ and $\text{diag} \left(\text{E} \left[\mathbf{d}_{\mathcal{Q}}[t - \ell'_1] \mathbf{d}_{\mathcal{Q}}^{\text{H}}[t - \ell'_2] \right] \right)$, $\ell'_1 = \ell'_2$,
- $\text{nondiag} \left(\text{E} \left[\mathbf{d}_{\mathcal{Q}}^{\text{UL}}[t - \ell'_1] \mathbf{d}_{\mathcal{Q}}^{\text{UL,H}}[t - \ell'_2] \right] \right)$ and $\text{nondiag} \left(\text{E} \left[\mathbf{d}_{\mathcal{Q}}[t - \ell'_1] \mathbf{d}_{\mathcal{Q}}^{\text{H}}[t - \ell'_2] \right] \right)$, $\ell'_1 = \ell'_2$,

and

- $\text{E} \left[\mathbf{d}_{\mathcal{Q}}^{\text{UL}}[t - \ell'_1] \mathbf{d}_{\mathcal{Q}}^{\text{UL,H}}[t - \ell'_2] \right]$ and $\text{E} \left[\mathbf{d}_{\mathcal{Q}}[t - \ell'_1] \mathbf{d}_{\mathcal{Q}}^{\text{H}}[t - \ell'_2] \right]$, $\ell'_1 \neq \ell'_2$.

The above covariance matrices can be computed exactly with Price's theorem according to (3.87). Note that the LCA makes the second and third terms vanish.

Second, we compute the covariance matrices stated above according to Price's theorem and plug the resulting expressions in (6.43) and (6.45). Thus, the uplink and downlink MSE

expressions calculate to

$$\begin{aligned}
\text{MSE}^{\text{UL}} &= \alpha_q \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \sum_{\ell=0}^{L-1} \mathbf{F}_{\ell'} \text{diag} \left(\mathbf{H}_{\ell}^{\text{H}} \mathbf{C}_{\eta\eta}^{-1/2, \text{H}} \mathbf{T} \mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{T} \mathbf{C}_{\eta\eta}^{-1/2} \mathbf{H}_{\ell} \right) \mathbf{F}_{\ell'}^{\text{H}} \right) \\
&+ \alpha_q^2 \text{tr} \left(\sum_{\substack{\ell'_1=0 \\ \ell'_2 \neq \ell'_1}}^{L_p-1} \sum_{\substack{\ell'_2=0 \\ \Delta\ell=\ell'_1}}^{L_p-1} \sum_{\Delta\ell=\ell'_1}^{L-1+\ell'_1} \mathbf{F}_{\ell'_1} \mathbf{H}_{\Delta\ell-\ell'_1}^{\text{H}} \mathbf{C}_{\eta\eta}^{-1/2, \text{H}} \mathbf{T} \mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{T} \mathbf{C}_{\eta\eta}^{-1/2} \mathbf{H}_{\Delta\ell-\ell'_2} \mathbf{F}_{\ell'_2}^{\text{H}} \right) \\
&+ \frac{Q}{\pi} \sin^2(\psi) \alpha_q \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \mathbf{F}_{\ell'} \sum_{\Delta k=0}^{Q/2-1} e^{j2\Delta k\psi} \text{diag}(\mathbf{C}_{\text{xx}}^{\text{UL}})^{1/2} \right. \\
&\quad \left. \arcsin \left(\Re \left\{ \text{nondiag} \left(\mathbf{R}_{\text{xx}}^{\text{UL}} \right) e^{-j2\Delta k\psi} \right\} \right) \text{diag}(\mathbf{C}_{\text{xx}}^{\text{UL}})^{1/2} \mathbf{F}_{\ell'}^{\text{H}} \right) \\
&+ \text{tr} \left(\sum_{\substack{\ell'_1=0 \\ \ell'_2 \neq \ell'_1}}^{L_p-1} \sum_{\ell'_2=0}^{L_p-1} \mathbf{F}_{\ell'_1} \text{E} \left[\mathbf{d}_{\mathcal{Q}}^{\text{UL}}[t-\ell'_1] \mathbf{d}_{\mathcal{Q}}^{\text{UL,H}}[t-\ell'_2] \right] \mathbf{F}_{\ell'_2}^{\text{H}} \right) \\
&- \alpha_q \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \mathbf{F}_{\ell'} \mathbf{H}_{\tau-\ell'}^{\text{H}} \mathbf{C}_{\eta\eta}^{-1/2, \text{H}} \mathbf{T} \mathbf{C}_{\text{ss}}^{\text{UL}} \right) - \alpha_q \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{T} \mathbf{C}_{\eta\eta}^{-1/2} \mathbf{H}_{\tau-\ell'} \mathbf{F}_{\ell'}^{\text{H}} \right) \\
&+ \alpha_q \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \mathbf{F}_{\ell'} \mathbf{C}_{\eta\eta}^{\text{UL}} \mathbf{F}_{\ell'}^{\text{H}} \right) + \text{tr}(\mathbf{C}_{\text{ss}}^{\text{UL}}), \tag{6.46}
\end{aligned}$$

and

$$\begin{aligned}
\text{MSE}^{\text{DL}} &= \alpha_q \text{tr} \left(\sum_{\ell=0}^{L-1} \sum_{\ell'=0}^{L_p-1} \mathbf{G} \mathbf{H}_{\ell} \text{diag} \left(\mathbf{P}_{\ell'} \mathbf{C}_{\text{ss}} \mathbf{P}_{\ell'}^{\text{H}} \right) \mathbf{H}_{\ell}^{\text{H}} \mathbf{G} \right) \\
&+ \alpha_q^2 \text{tr} \left(\sum_{\substack{\ell'_1=0 \\ \ell'_2 \neq \ell'_1}}^{L_p-1} \sum_{\substack{\ell'_2=0 \\ \Delta\ell=\ell'_1}}^{L_p-1} \sum_{\Delta\ell=\ell'_1}^{L-1+\ell'_1} \mathbf{G} \mathbf{H}_{\Delta\ell-\ell'_1} \mathbf{P}_{\ell'_1} \mathbf{C}_{\text{ss}} \mathbf{P}_{\ell'_2}^{\text{H}} \mathbf{H}_{\Delta\ell-\ell'_2}^{\text{H}} \mathbf{G} \right) \\
&+ \frac{Q}{\pi} \sin^2(\psi) \alpha_q \text{tr} \left(\sum_{\ell=0}^{L-1} \mathbf{G} \mathbf{H}_{\ell} \sum_{\Delta k=0}^{Q/2-1} e^{j2\Delta k\psi} \text{diag}(\mathbf{C}_{\text{xx}})^{1/2} \right. \\
&\quad \left. \arcsin \left(\Re \left\{ \text{nondiag} \left(\mathbf{R}_{\text{xx}} \right) e^{-j2\Delta k\psi} \right\} \right) \text{diag}(\mathbf{C}_{\text{xx}})^{1/2} \mathbf{H}_{\ell}^{\text{H}} \mathbf{G} \right) \\
&+ \text{tr} \left(\sum_{\ell_1=0}^{L-1} \sum_{\ell_2=0, \ell_2 \neq \ell_1}^{L-1} \mathbf{G} \mathbf{H}_{\ell_1} \text{E} \left[\mathbf{d}_{\mathcal{Q}}[t-\ell_1] \mathbf{d}_{\mathcal{Q}}^{\text{H}}[t-\ell_2] \right] \mathbf{H}_{\ell_2}^{\text{H}} \mathbf{G} \right) \\
&- \alpha_q \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \mathbf{G} \mathbf{H}_{\tau-\ell'} \mathbf{P}_{\ell'} \mathbf{C}_{\text{ss}} \right) - \alpha_q \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \mathbf{C}_{\text{ss}} \mathbf{P}_{\ell'}^{\text{H}} \mathbf{H}_{\tau-\ell'}^{\text{H}} \mathbf{G} \right) \\
&+ \text{tr}(\mathbf{G} \mathbf{C}_{\eta\eta} \mathbf{G}) + \text{tr}(\mathbf{C}_{\text{ss}}). \tag{6.47}
\end{aligned}$$

By applying the identities in (4.39), (4.40), (4.41) and (4.42), the uplink MSE expression recalculates to

$$\begin{aligned}
\text{MSE}^{\text{UL}} &= \alpha_q \text{tr} \left(\sum_{\ell=0}^{L_p-1} \sum_{\ell'=0}^{L_p-1} \mathbf{G} \mathbf{H}_\ell \text{diag}(\mathbf{P}_{\ell'} \mathbf{C}_{\text{ss}} \mathbf{P}_{\ell'}^{\text{H}}) \mathbf{H}_\ell^{\text{H}} \mathbf{G} \right) \\
&+ \alpha_q^2 \text{tr} \left(\sum_{\ell'_1=0}^{L_p-1} \sum_{\substack{\ell'_2=0 \\ \ell'_2 \neq \ell'_1}}^{L_p-1} \sum_{\Delta \ell = \ell'_1}^{L-1+\ell'_1} \mathbf{G} \mathbf{H}_{\Delta \ell - \ell'_1} \mathbf{P}_{\ell'_1} \mathbf{C}_{\text{ss}} \mathbf{P}_{\ell'_2}^{\text{H}} \mathbf{H}_{\Delta \ell - \ell'_2}^{\text{H}} \mathbf{G} \right) \\
&+ \frac{Q}{\pi} \sin^2(\psi) \frac{\alpha_q}{\beta^2} \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \mathbf{P}_{\ell'}^{\text{H}} \sum_{\Delta k=0}^{Q/2-1} e^{j2\Delta k \psi} \text{diag}(\mathbf{C}_{\mathbf{xx}}^{\text{UL}})^{1/2} \right. \\
&\quad \left. \arcsin(\Re \{ \text{nondiag}(\mathbf{R}_{\mathbf{xx}}^{\text{UL}}) e^{-j2\Delta k \psi} \}) \text{diag}(\mathbf{C}_{\mathbf{xx}}^{\text{UL}})^{1/2} \mathbf{P}_{\ell'} \right) \\
&+ \frac{1}{\beta^2} \text{tr} \left(\sum_{\ell'_1=0}^{L_p-1} \sum_{\substack{\ell'_2=0 \\ \ell'_2 \neq \ell'_1}}^{L_p-1} \mathbf{P}_{\ell'_1}^{\text{H}} \mathbf{E} \left[\mathbf{d}_Q^{\text{UL}}[t - \ell'_1] \mathbf{d}_Q^{\text{UL,H}}[t - \ell'_2] \right] \mathbf{P}_{\ell'_2} \right) \\
&- \alpha_q \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \mathbf{G} \mathbf{H}_{\tau - \ell'} \mathbf{P}_{\ell'} \mathbf{C}_{\text{ss}} \right) - \alpha_q \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \mathbf{C}_{\text{ss}} \mathbf{P}_{\ell'}^{\text{H}} \mathbf{H}_{\tau - \ell'}^{\text{H}} \mathbf{G} \right) \\
&+ \frac{\alpha_q}{\beta^2 \sigma_s^2} \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \mathbf{P}_{\ell'} \mathbf{C}_{\text{ss}} \mathbf{P}_{\ell'}^{\text{H}} \right) + \text{tr}(\mathbf{C}_{\text{ss}}). \tag{6.48}
\end{aligned}$$

By comparing (6.48) to (6.47), we can see that all terms match except of the third, fourth and seventh terms. The seventh terms can be matched by choosing

$$\begin{aligned}
\beta &= \sqrt{\frac{\alpha_q \text{tr} \left(\sum_{\ell'=0}^{L_p-1} \mathbf{P}_{\ell'} \mathbf{C}_{\text{ss}} \mathbf{P}_{\ell'}^{\text{H}} \right)}{\sigma_s^2 \text{tr}(\mathbf{G} \mathbf{C}_{\eta\eta} \mathbf{G})}} \\
&\stackrel{(4.43)}{=} \sqrt{\frac{P_{\text{tx}}}{\text{tr}(\mathbf{G} \mathbf{C}_{\eta\eta} \mathbf{G})}}. \tag{6.49}
\end{aligned}$$

The obtained expression for β is the same as for flat-fading channels. In analogy to Section 5.4.2, we assume that \mathbf{G} is a scaled identity, i.e. $\mathbf{G} = g \mathbf{I}_M$ and thus from (4.40) and (6.49) we get

$$\mathbf{T} = \sqrt{\frac{P_{\text{tx}}}{\sigma_s^2 \text{tr}(\mathbf{C}_{\eta\eta})}} \mathbf{C}_{\eta\eta}^{1/2, \text{H}}. \tag{6.50}$$

Hence, the transmit power constraint in the uplink system holds true; that is

$$\text{tr}(\mathbf{T} \mathbf{C}_{\text{ss}}^{\text{UL}} \mathbf{T}^{\text{H}}) = P_{\text{tx}}. \tag{6.51}$$

For exact MSE duality, the third and fourth terms in the uplink and downlink MSE expressions must additionally match.

6.4.2 Approximate Dual Problem

By using the LCA instead of Price's theorem, the third and fourth terms in the uplink and downlink MSE expressions in (6.48) and (6.47), respectively, vanish. Thus, the MSE expressions match. In other words, we can achieve approximate duality by applying the LCA for the computation of the MSE.

After considering the approximate dual uplink system, the optimal MMSE equalizer $\mathbf{F}[t]$ has to be designed. To this end, we state the compact input-output relationship as

$$\mathbf{u}^{\text{UL}}[t] = \mathbf{u} = \mathbf{F}\mathbf{t}^{\text{UL}}, \quad (6.52)$$

where

$$\mathbf{F} = [\mathbf{F}_0 \quad \mathbf{F}_1 \quad \cdots \quad \mathbf{F}_{L_p-1}], \quad (6.53)$$

$$\mathbf{t}^{\text{UL}} = \mathcal{Q}_{\text{CE}}(\mathbf{x}^{\text{UL}}), \quad (6.54)$$

$$\begin{aligned} \mathbf{x}^{\text{UL}} &= [\mathbf{x}[t]^{\text{UL,T}} \quad \mathbf{x}[t-1]^{\text{UL,T}} \quad \cdots \quad \mathbf{x}[t-L_p+1]^{\text{UL,T}}]^{\text{T}} \\ &= \mathbf{H}^{\text{UL}}\mathbf{s}^{\text{UL}}, \end{aligned} \quad (6.55)$$

$$\mathbf{H}^{\text{UL}} = \sum_{\ell=0}^{L-1} \mathbf{S}(\ell, L, L_p-1) \otimes (\mathbf{H}_{\ell}^{\text{H}}\mathbf{C}_{\eta\eta}^{-1/2,\text{H}}\mathbf{T}), \quad (6.56)$$

and

$$\mathbf{s}^{\text{UL}} = [\mathbf{s}[t]^{\text{UL,T}} \quad \mathbf{s}[t-1]^{\text{UL,T}} \quad \cdots \quad \mathbf{s}[t-L_p-L+2]^{\text{UL,T}}]^{\text{T}}. \quad (6.57)$$

The uplink MSE is expressed in terms of the defined variables as

$$\begin{aligned} \text{MSE}^{\text{UL}} &= \text{E} \left[\|\mathbf{u}^{\text{UL}} - (\mathbf{e}_{\tau+1}^{\text{T}} \otimes \mathbf{I}_M) \mathbf{s}^{\text{UL}}\|_2^2 \right] \\ &= \text{tr}(\mathbf{F}\mathbf{C}_{\text{tt}}\mathbf{F}^{\text{H}}) - \text{tr}(\mathbf{F}\mathbf{C}_{\text{ts}}(\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M)) - \text{tr}((\mathbf{e}_{\tau+1}^{\text{T}} \otimes \mathbf{I}_M)\mathbf{C}_{\text{st}}\mathbf{F}^{\text{H}}) \\ &\quad + \text{tr}((\mathbf{e}_{\tau+1}^{\text{T}} \otimes \mathbf{I}_M)\mathbf{C}_{\text{ss}}(\mathbf{e}_{\tau+1} \otimes \mathbf{I}_M)). \end{aligned} \quad (6.58)$$

Hence, the optimal MMSE dual filter is then given by

$$\mathbf{F} = (\mathbf{e}_{\tau+1}^{\text{T}} \otimes \mathbf{I}_M) \mathbf{C}_{\text{st}}^{\text{UL}} \mathbf{C}_{\text{tt}}^{\text{UL},-1}. \quad (6.59)$$

First, we apply Price's theorem to compute $\mathbf{C}_{\text{tt}}^{\text{UL}}$ and $\mathbf{C}_{\text{st}}^{\text{UL}}$ and hence obtain the filter \mathbf{F} . The matrix \mathbf{T} is given in (6.50). Second, the optimal filters in the primal domain can be computed using the identities (4.39) and (4.40). The corresponding scaling factor β in (4.44) can be expressed as

$$\beta = \sqrt{\frac{P_{\text{tx}}}{\alpha_q \text{tr}(\mathbf{F}^{\text{H}}\mathbf{C}_{\text{ss}}\mathbf{F})}}. \quad (6.60)$$

6.5 Simulation Results

In this section, we compare the performance of the different linear precoding techniques that were introduced above with the ideal WF precoder while assuming full Channel State Information (CSI). We denote the precoders by

- FIR-WF: the ideal FIR-WF, [70], where no CEQ is applied in the system.
- FIR-QWF: the FIR-Quantized WF that was introduced in Section 6.3.2.
- FIR-QWF-Apprx_{Dual}: the approximate dual Quantized WF that applies unequal power allocation at the antennas. The resulting receive filter \mathbf{G} is a scaled identity. This precoder is derived in Section 6.4.2.
- FIR-QWF-Price_{eqPA}: the FIR-Quantized WF based on Price's theorem that applies equal power allocation at the antennas. This precoder is detailed in Algorithm 4. The start value is the FIR-WF precoder that is projected to fulfill the equal transmit power constraint.
- FIR-QWF-Price_{neqPA}: the FIR-Quantized WF based on Price's theorem that applies unequal power allocation at the antennas. This precoder is detailed in Algorithm 5. The start value is the FIR-WF precoder.

For all precoders, it holds that $L_p = L$. To this end, we assume a BS with $N = 64$ or $N = 128$ antennas serving $M = 8$ single-antenna users with 4-PSK or 16-QAM signals, respectively. The numerical results are obtained with Monte Carlo simulations of 100 independent frequency-selective channel realizations of $L = 3$ from the i.i.d. channel model and the mmW sparse channel model described in Section 4.6. For the mmW sparse channel, we assume that $N_{\text{cl}} = 2$ and $N_{\text{ray}} = 10$. The AWGN is also i.i.d. with variance one at each antenna. The performance metric is the uncoded BER averaged over the single-antenna users. For the blind estimation of the coefficients g_m we use a block length of $T = 128$. The numerical results are plotted in Fig. 6.1, Fig. 6.2, Fig. 6.3 and Fig. 6.4.

First, it can be deduced that all proposed linear precoders perform almost the same for i.i.d. channels. The loss compared to the ideal WF reduces with increased quantization resolution.

Second, the proposed precoders perform slightly different for mmW sparse channels. We can see that the QWF-Price_{neqPA} and the QWF-Price_{eqPA} precoders perform the best and have almost the same behavior. In the presence of the ISI, there is almost no much difference between equal and unequal power allocation at the BS antennas.

Third, as observed in the case of flat-fading channels, a QCE system requires at most the double number of antennas to obtain the same BER performance as in the ideal unquantized case.

In summary, the performance improvement between the proposed linear precoding algorithms that take into account the quantization distortions is too small and the performance gap for the same system settings to the ideal linear system is still large. Therefore, we decide to study non-linear precoding techniques with the hope that they offer a much better performance.

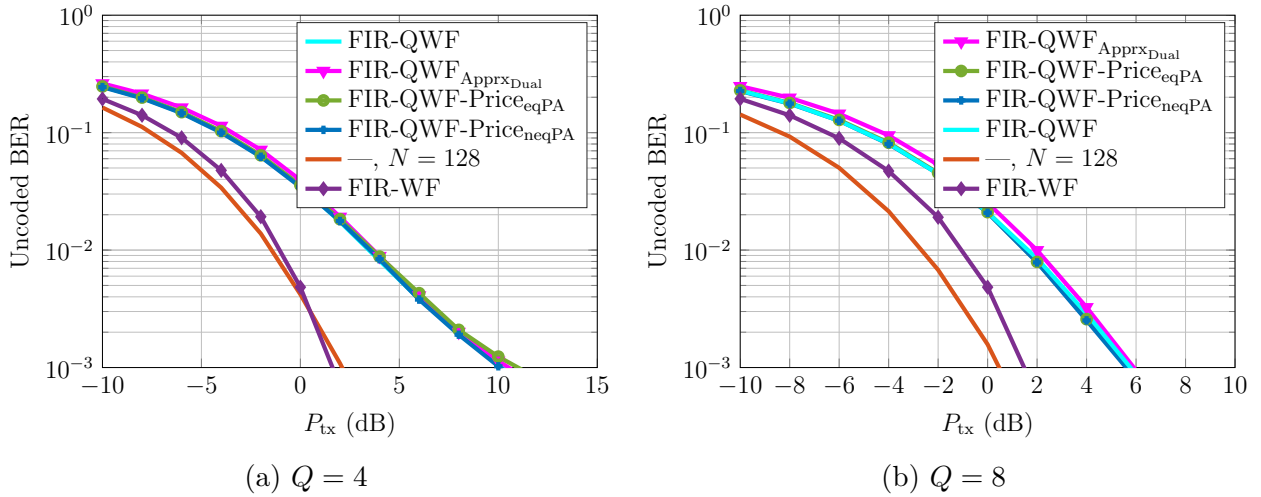


Fig. 6.1: Comparison between different linear precoders for 4-PSK signaling: $N = 64$, $M = 8$ with the i.i.d. channel with $L = 3$.

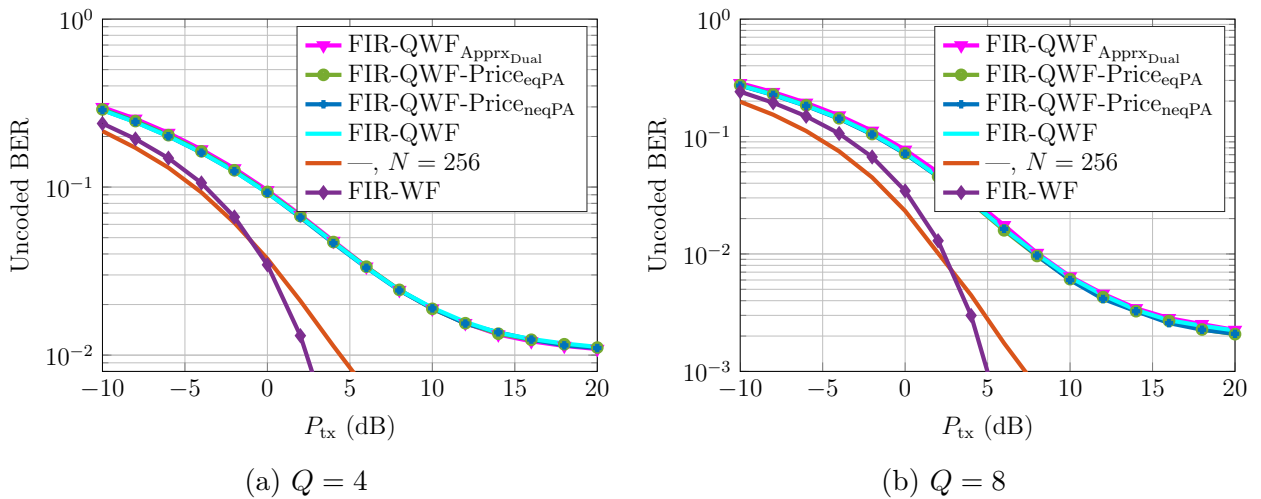


Fig. 6.2: Comparison between different linear precoders for 16-QAM signaling: $N = 128$, $M = 8$ with the i.i.d. channel with $L = 3$.

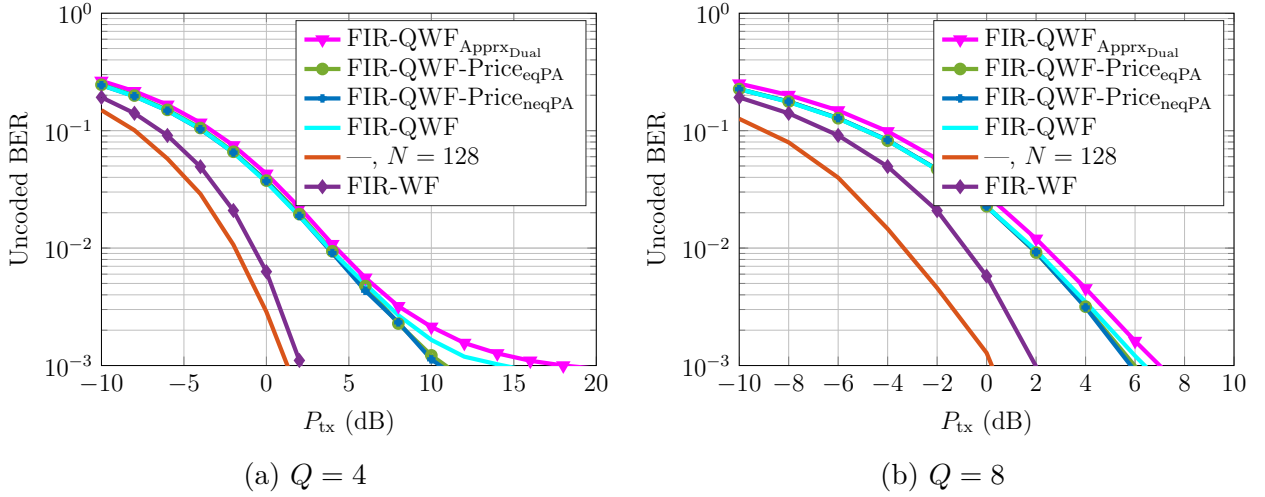


Fig. 6.3: Comparison between different linear precoders for 4-PSK signaling: $N = 64$, $M = 8$ with the mmW sparse channel with $L = 3$.

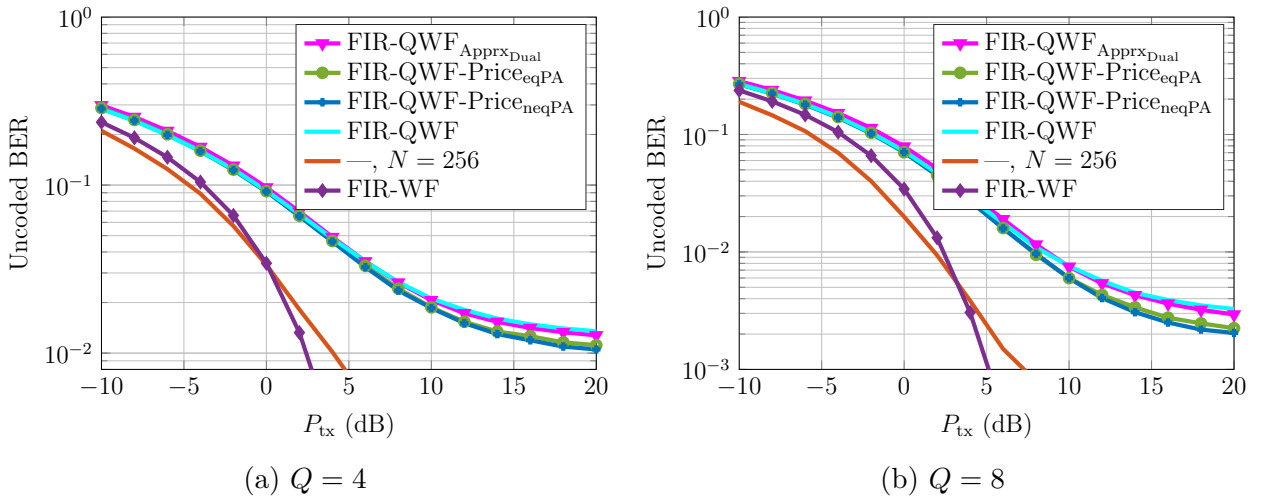


Fig. 6.4: Comparison between different linear precoders for 16-QAM signaling: $N = 128$, $M = 8$ with the mmW sparse channel with $L = 3$.

Part II

Non-linear Transmit Signal Processing

7. Flat-Fading Channels

As explained in Section 4.4.2, the signal vector \mathbf{x} is designed rather than the precoder matrix \mathbf{P} . For every given input signal \mathbf{s} and for each channel realization \mathbf{H} , the precoding task is to find

$$\mathbf{x} = \mathcal{P}(\mathbf{s}, \mathbf{H}). \quad (7.1)$$

The task consists in designing the transmit vector \mathbf{x} such that $\hat{\mathbf{s}} = \mathbf{s}$ holds true with high probability to reduce the detection error probability. The symbol-wise precoder aims to mitigate all sources of distortion

- the quantization distortions
- the channel distortions and the MUI, and
- the AWGN.

Our goal is to develop a problem formulation that jointly minimizes all three distortion sources.

7.1 Transmit Power Constraint

The transmit power is allocated equally among the antennas. From (4.23), it follows that

$$\mathbf{\Xi} = \sqrt{\frac{P_{\text{tx}}}{N}} \mathbf{I}_N. \quad (7.2)$$

Thus, the set \mathbb{T}_n , $n = 1, \dots, N$, reads as

$$\mathbb{T}_n = \mathbb{T} = \left\{ \sqrt{\frac{P_{\text{tx}}}{N}} \exp\left(j(2i-1)\frac{\pi}{Q}\right) : i = 1, \dots, Q \right\}. \quad (7.3)$$

7.2 Optimization Problem

To formulate the optimization problem, we have to consider all distortion sources and find a way to mitigate them one by one.

7.2.1 Mitigation of the Quantization Distortions

First, it is obvious that the quantization distortions can be omitted if we design the precoded vector \mathbf{x} such that

$$\mathbf{x} \in \mathbb{T}^N, \quad (7.4)$$

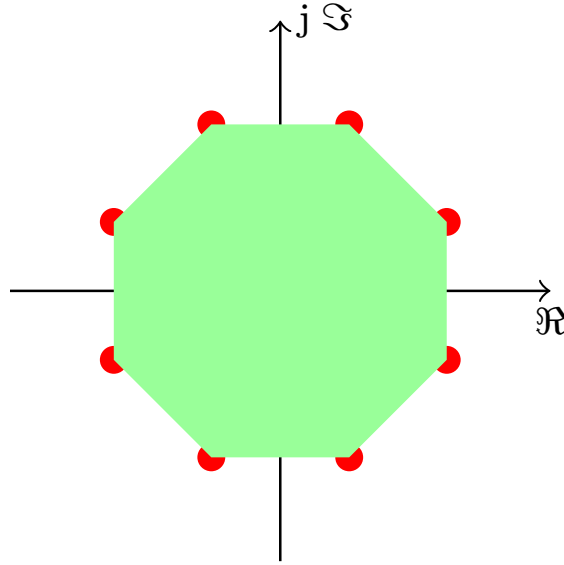


Fig. 7.1: Illustration of the relaxed polygon constraint for $Q=8$, © 2018 IEEE.

i.e. $\mathbb{X} = \mathbb{T}$. This would ensure that the quantizer $\mathcal{Q}_{\text{CE}}(\bullet)$ produces no distortion, and we would have an undistorted transmit signal $\mathbf{t} = \mathbf{x}$. However, the QCE constraint in (7.4) would lead to a discrete optimization problem due to the discrete nature of the set \mathbb{T} . To avoid this problem, we relax the discrete set \mathbb{T} to the convex set \mathbb{X} that represents the polygon built by the Q scaled PSK points of the set \mathbb{T} . Thus, the QCE constraint is relaxed to a convex constraint that we call the relaxed polygon constraint. Fig. 7.1 illustrates the relaxed polygon constraint for the case of $Q = 8$. Instead of designing $\mathbf{x} \in \mathbb{T}^N$ to completely eliminate the quantization distortions, we design $\mathbf{x} \in \mathbb{X}^N$ to minimize them.

The set \mathbb{X} can be mathematically described by a set of linear inequalities. For q -bit polar DACs, i.e., where the transmitted data are constrained to be Q scaled PSK symbols, the polygon can be constructed by the intersection of $Q/4$ squares that have an angular shift of $2\pi/Q$. To this end, we define the rotation matrix \mathbf{R}_i of angle $\beta_i = \frac{2\pi}{Q}(i-1)$ as

$$\mathbf{R}_i = \begin{bmatrix} \cos \beta_i & \sin \beta_i \\ -\sin \beta_i & \cos \beta_i \end{bmatrix} \otimes \mathbf{I}_N, \quad i = 1, \dots, Q/4. \quad (7.5)$$

The system of inequalities that considers the feasible set, i.e. the relaxed polygon constraint, and hence relaxes the constraint in (7.4) is given by

$$\left[\mathbf{R}_1^T - \mathbf{R}_1^T \cdots \mathbf{R}_{\frac{Q}{4}}^T - \mathbf{R}_{\frac{Q}{4}}^T \right]^T \bar{\mathbf{x}} \leq \sqrt{\frac{P_{\text{tx}}}{N}} \cos\left(\frac{\pi}{Q}\right) \mathbf{1}_{NQ}, \quad (7.6)$$

where $\bar{\mathbf{x}} = [\Re\{\mathbf{x}\}^T \quad \Im\{\mathbf{x}\}^T]^T$. Since $\mathbf{R}_1 = \mathbf{I}_{2N}$, the first $4N$ inequalities in (7.6) define the bounds for $\bar{\mathbf{x}}$. Hence, the relaxed polygon constraint, i.e. $\mathbf{x} \in \mathbb{X}^N$, is equivalent to

$$\begin{aligned} -\sqrt{\frac{P_{\text{tx}}}{N}} \cos\left(\frac{\pi}{Q}\right) \mathbf{1}_{2N} &\leq \bar{\mathbf{x}} \leq \sqrt{\frac{P_{\text{tx}}}{N}} \cos\left(\frac{\pi}{Q}\right) \mathbf{1}_{2N}, \\ \text{and } \underbrace{\left[\mathbf{R}_2^T - \mathbf{R}_2^T \cdots \mathbf{R}_{\frac{Q}{4}}^T - \mathbf{R}_{\frac{Q}{4}}^T \right]^T}_{=\mathbf{E}} \bar{\mathbf{x}} &\leq \sqrt{\frac{P_{\text{tx}}}{N}} \cos\left(\frac{\pi}{Q}\right) \mathbf{1}_{N(Q-4)}. \end{aligned} \quad (7.7)$$

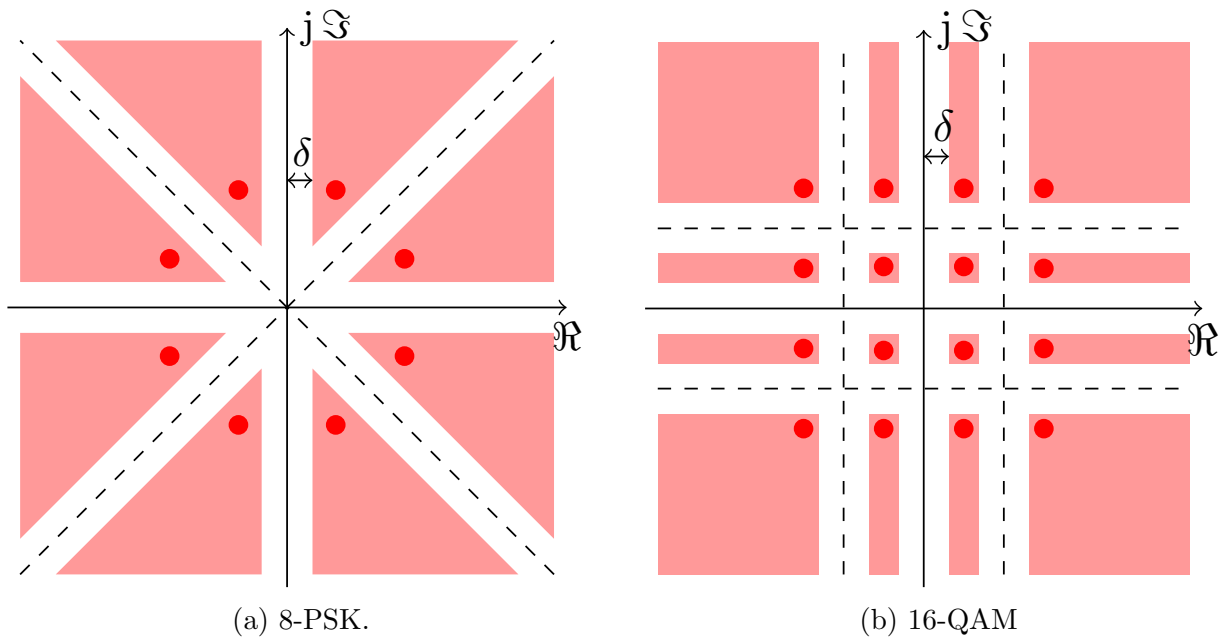


Fig. 7.2: Decision regions and SRs (in red) for different constellations, © 2018 IEEE.

This reformulation leads to significant computational savings since the final optimization problem will be written as a linear program with bounded variables. As discussed in Section 7.5, it is beneficial in terms of computational complexity to have fewer inequalities.

7.2.2 Mitigation of the Channel Distortions and the MUI

Second, to minimize the channel distortions and the MUI, we look deeper into the properties of the constellations. As illustrated in Fig. 7.2a and Fig. 7.2b, each constellation is defined by thresholds that separate the distinct decision regions of the constellation points. In total, we have as many contiguous decision regions as constellation points. For each outer constellation point, the decision region has at least one infinite boundary. Therefore, we make use of the idea of constructive interference optimization [42, 43]. When the downlink channel and all users' data are known at the transmitter, the instantaneous constructive MUI can be exploited to move the received signals further from the decision thresholds [43]. In contrast to this, conventional precoding methods, e.g. MMSE, aim at minimizing the total MUI such that the received signals lie as close as possible to the nominal constellation points. Constructive interference optimization exploits the larger symbol decision regions and thus leads to a more relaxed optimization.

Each constellation symbol lies within a Symbol Region (SR) that is a downscaled version of the decision region. In contrast to the decision region, the SR has a safety margin denoted by δ that separates it from the decision thresholds. When each entry of the noiseless received signal vector \mathbf{y} belongs to the correct SR and thus to the correct decision region, the channel distortions and the MUI are mitigated.

7.2.3 Mitigation of the AWGN

Finally, the safety margin δ has to be large enough such that, when \mathbf{y} is perturbed by the AWGN, the received signals do not jump to unintended neighboring decision regions.

7.2.4 General Problem Formulation

In summary, the problem formulation has to take into account the relaxed QCE constraint in (7.7), the SR constraint for each received signal and maximizing the safety margin δ . Thus, the optimization problem for the symbol-wise precoder, which we call the MSM precoder, can be written in general as follows

$$\max_{\mathbf{x}} \delta \quad (7.8)$$

$$\text{s.t. } y'_m \in \text{SR}_m, \forall m \quad (7.9)$$

$$\text{and } \mathbf{x} \in \mathbb{X}^N, \quad (7.10)$$

where

$$\mathbf{y}' = \mathbf{H}\mathbf{x} \quad (7.11)$$

represents the relaxation of \mathbf{y} due to the relaxation of (7.4) to (7.7). Exact expressions for the safety margin δ and the SRs as a function of \mathbf{x} are provided in Section 7.3 and Section 7.4 for PSK and QAM signaling, respectively.

This problem formulation depends on the input symbol vector \mathbf{s} , which determines the intended SR for each received signal and on the channel \mathbf{H} . The optimization is run for one specific value of the transmit power, i.e. $P_{\text{tx}} = N$. Since the optimization variables depend linearly on the square-root of the transmit power $\sqrt{P_{\text{tx}}}$, it is sufficient to solve the optimization problem for one transmit power value. Once the optimal relaxed solution \mathbf{x} is found, the transmit vector \mathbf{t} is obtained as stated in (4.1).

7.3 Problem Formulation for PSK Signaling

7.3.1 SR for PSK Signals

In this section, we assume that the input signals s_m , $m = 1, \dots, M$, belong to the S -PSK constellation. The set \mathbb{S} in this case is defined in (4.13).

Each SR in the PSK constellation, as shown in Fig. 7.2a, is a circular sector of infinite radius and angle 2θ . To find a mathematical expression for each SR, we rotate the original coordinate system by the phase of the symbol of interest s_m , as illustrated in Fig. 7.3, and introduce a new variable z_m that represents the signal y'_m in the new coordinate system such that

$$\begin{aligned} z_m &= y'_m \frac{s_m^*}{|s_m|}, \quad m = 1, \dots, M \\ &= y'_m s_m^*, \quad m = 1, \dots, M, \end{aligned} \quad (7.12)$$

since PSK signals have unit magnitude. The m -th SR, $m = 1, \dots, M$, can be hence described by

$$\Re\{z_m\} \geq \delta', \quad (7.13)$$

$$|\Im\{z_m\}| \leq (\Re\{z_m\} - \delta') \tan \theta, \quad (7.14)$$

where $\delta' = \frac{\delta}{\sin\theta}$. Note that the inequality in (7.13) is already fulfilled if the inequality in (7.14) is satisfied. In the vector notation, the vector \mathbf{z} can be expressed as

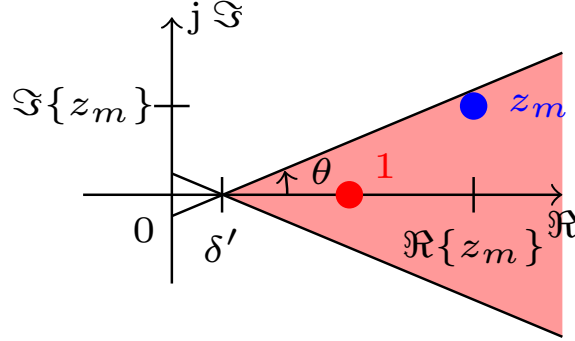


Fig. 7.3: Illustration of the PSK SR in the modified coordinate system.

$$\begin{aligned}\mathbf{z} &= \text{diag}(\mathbf{s}^*)\mathbf{H}\mathbf{x} \\ &= \tilde{\mathbf{H}}\mathbf{x},\end{aligned}\tag{7.15}$$

where

$$\tilde{\mathbf{H}} = \text{diag}(\mathbf{s}^*)\mathbf{H}.\tag{7.16}$$

Hence, all M SRs can be compactly described as

$$|\Im\{\mathbf{z}\}| \leq (\Re\{\mathbf{z}\} - \delta'\mathbf{1}_M) \tan\theta,\tag{7.17}$$

as also given in [43]. When using the following real-valued representation

$$\Re\{\mathbf{z}\} = \underbrace{\begin{bmatrix} \Re\{\tilde{\mathbf{H}}\} & -\Im\{\tilde{\mathbf{H}}\} \end{bmatrix}}_{=\mathbf{A}} \begin{bmatrix} \Re\{\mathbf{x}\} \\ \Im\{\mathbf{x}\} \end{bmatrix} = \mathbf{A}\bar{\mathbf{x}},\tag{7.18}$$

$$\Im\{\mathbf{z}\} = \underbrace{\begin{bmatrix} \Im\{\tilde{\mathbf{H}}\} & \Re\{\tilde{\mathbf{H}}\} \end{bmatrix}}_{=\mathbf{B}} \begin{bmatrix} \Re\{\mathbf{x}\} \\ \Im\{\mathbf{x}\} \end{bmatrix} = \mathbf{B}\bar{\mathbf{x}},\tag{7.19}$$

the constraint in (7.17) can be rewritten as

$$\begin{bmatrix} \mathbf{B} - \tan\theta\mathbf{A} & \frac{1}{\cos\theta}\mathbf{1}_M \\ -\mathbf{B} - \tan\theta\mathbf{A} & \frac{1}{\cos\theta}\mathbf{1}_M \end{bmatrix} \begin{bmatrix} \bar{\mathbf{x}} \\ \delta \end{bmatrix} \leq \mathbf{0}_{2M}.\tag{7.20}$$

7.3.2 Optimization Problem with the Relaxed Polygon Constraint

Finally, the optimization problem for the symbol-wise precoder with PSK signaling is obtained by combining (7.8), (7.20) and (7.7) and is expressed for the case of $P_{\text{tx}} = N$ as

$$\begin{aligned}\max_{\mathbf{v}} \quad & \begin{bmatrix} \mathbf{0}_{2N}^T & \mathbf{1} \end{bmatrix} \mathbf{v} \text{ s.t. } \begin{bmatrix} \mathbf{B} - \tan\theta\mathbf{A} & \frac{1}{\cos\theta}\mathbf{1}_M \\ -\mathbf{B} - \tan\theta\mathbf{A} & \frac{1}{\cos\theta}\mathbf{1}_M \\ \mathbf{E} & \mathbf{0}_{N(Q-4)} \end{bmatrix} \mathbf{v} \leq \begin{bmatrix} \mathbf{0}_{2M} \\ \cos\left(\frac{\pi}{Q}\right)\mathbf{1}_{N(Q-4)} \end{bmatrix}, \\ \text{and} \quad & \begin{bmatrix} -\cos\left(\frac{\pi}{Q}\right)\mathbf{1}_{2N} \\ 0 \end{bmatrix} \leq \mathbf{v} \leq \begin{bmatrix} \cos\left(\frac{\pi}{Q}\right)\mathbf{1}_{2N} \\ \infty \end{bmatrix},\end{aligned}\tag{7.21}$$

where $\mathbf{v}^T = [\bar{\mathbf{x}}^T \ \delta]$. The resulting optimization problem is a linear programming problem for which there exist very efficient solving methods [71]. In order to solve the problem for different P_{tx} values, it is sufficient to scale the solution of (7.21) by $\frac{P_{\text{tx}}}{N}$ due to the linearity of the problem.

When the optimization terminates, the optimal signal $\mathbf{x} \in \mathbb{X}^N$ is found. The signal \mathbf{t} that goes through the channel is obtained as described in (4.1). In other words, each entry in \mathbf{x} gets mapped to the corresponding CE point depending on the circular sector that it lies in.

7.3.3 Safety Margin for PSK Signals

The safety margin δ in (7.8) can be expressed for the PSK case as

$$\begin{aligned} \delta &= \min_m \delta_m \\ &= \min_m (\sin(\theta)\Re\{\mathbf{z}\} - \cos(\theta)|\Im\{\mathbf{z}\}|), \end{aligned} \quad (7.22)$$

where the operator $|\bullet|$ is applied element-wise to the entries of $\Im\{\mathbf{z}\}$. Note that an equivalent objective function was introduced in [22] in the context of continuous-phase CE precoding for PSK signaling. In [22], the strict CE constraint is relaxed to the convex unit circle, whereas in our work the QCE constraint is relaxed to the linear polygon constraint. Consequently, due to the linear objective function and the linear constraints, our optimization problem can be formulated as a linear programming problem unlike [22].

7.3.4 Interpretation of the Safety Margin δ for PSK Signals

The safety margin δ is a parameter that affects the receiver Signal-to-Noise Ratio (SNR) and the SER. These relationships will be given for the relaxed problem; that is the quantization is omitted and we consider the relaxed received signal \mathbf{y}' instead of \mathbf{y} .

7.3.4.1 Safety Margin vs. Receiver SNR

The receiver SNR at the m -th user is given by

$$\text{SNR}_m = \frac{\text{E}[|y'_m|^2]}{\sigma_n^2} = \text{E}[|y'_m|^2], \quad (7.23)$$

since we assume unit-variance AWGN. The expected value can be computed by averaging over N_s transmit signals. Hence, we get

$$\text{SNR}_m \geq \frac{1}{N_s} \sum_{i=1}^{N_s} \left(\frac{\delta^{(i)}}{\sin \theta} \right)^2. \quad (7.24)$$

Thus, we can conclude that maximizing the safety margin δ leads in turn to maximizing the lower bound of the receive SNR at each user.

7.3.4.2 Safety Margin vs. SER

The statistics of the AWGN in the modified coordinate system do not change, since it just consists of a rotation. We denote the modified noise vector by $\tilde{\boldsymbol{\eta}}$, where $\tilde{\boldsymbol{\eta}} \sim \mathcal{CN}_{\mathbb{C}}(\mathbf{0}_M, \mathbf{I}_M)$.

The SER represents the probability that the entries of $\mathbf{z} + \tilde{\boldsymbol{\eta}}$ lie outside the corresponding decision regions. The SER of the m -th user for the t -th transmission can be then expressed as

$$\begin{aligned}
\text{SER}_m^{(t)} &= \Pr(\Im\{z_m^{(t)}\} \geq 0) \left(\Pr\left(\Re\{\tilde{\eta}_m + z_m^{(t)}\} \geq 0\right) \wedge \left(\Im\{\tilde{\eta}_m\} \geq \frac{\delta_m^{(t)}}{\cos(\theta)} + \Re\{\tilde{\eta}_m\} \tan(\theta)\right) \right) \\
&\quad + \Pr\left(\Re\{\tilde{\eta}_m + z_m^{(t)}\} \geq 0\right) \wedge \left(\Im\{\tilde{\eta}_m\} \leq -2|\Im\{z_m^{(t)}\}| - \frac{\delta_m^{(t)}}{\cos(\theta)} - \Re\{\tilde{\eta}_m\} \tan(\theta)\right) \\
&\quad + \Pr(\Im\{z_m^{(t)}\} < 0) \left(\Pr\left(\Re\{\tilde{\eta}_m + z_m^{(t)}\} \geq 0\right) \wedge \left(\Im\{\tilde{\eta}_m\} \leq -\frac{\delta_m^{(t)}}{\cos(\theta)} - \Re\{\tilde{\eta}_m\} \tan(\theta)\right) \right) \\
&\quad + \Pr\left(\Re\{\tilde{\eta}_m + z_m^{(t)}\} \geq 0\right) \wedge \left(\Im\{\tilde{\eta}_m\} \geq 2|\Im\{z_m^{(t)}\}| + \frac{\delta_m^{(t)}}{\cos(\theta)} + \Re\{\tilde{\eta}_m\} \tan(\theta)\right) \\
&\quad + \Pr(\Re\{\tilde{\eta}_m + z_m^{(t)}\} < 0) \\
&= 1 - \int_{-\Re\{z_m^{(t)}\}}^{\infty} \frac{1}{2} \left(\text{erf}\left(\frac{\delta_m^{(t)}}{\cos(\theta)} + \Re\{\tilde{\eta}_m\} \tan(\theta)\right) \right. \\
&\quad \left. + \text{erf}\left(2|\Im\{z_m^{(t)}\}| + \frac{\delta_m^{(t)}}{\cos(\theta)} + \Re\{\tilde{\eta}_m\} \tan(\theta)\right) \right) \frac{1}{\sqrt{\pi}} e^{-\Re\{\tilde{\eta}_m\}^2} d\Re\{\tilde{\eta}_m\}. \tag{7.25}
\end{aligned}$$

Since the erf function is monotonically increasing and it holds that $\delta^{(t)} \leq \delta_m^{(t)}, \forall m$ and $\Re\{z_m^{(t)}\} \geq \frac{\delta_m^{(t)}}{\sin(\theta)}$, we obtain

$$\text{SER}_m^{(t)} \leq 1 - \int_{-\frac{\delta^{(t)}}{\sin \theta}}^{\infty} \text{erf}\left(\frac{\delta_m^{(t)}}{\cos(\theta)} + \Re\{\tilde{\eta}_m\} \tan(\theta)\right) \frac{1}{\sqrt{\pi}} e^{-\Re\{\tilde{\eta}_m\}^2} d\Re\{\tilde{\eta}_m\}. \tag{7.26}$$

Thus, the SER at the m -th user averaged over a transmission block of length N_s is upper bounded by

$$\text{SER}_m \leq 1 - \frac{1}{N_s} \sum_{t=1}^{N_s} \int_{-\frac{\delta^{(t)}}{\sin \theta}}^{\infty} \text{erf}\left(\frac{\delta^{(t)}}{\cos(\theta)} + \Re\{\tilde{\eta}_m\} \tan(\theta)\right) \frac{1}{\sqrt{\pi}} e^{-\Re\{\tilde{\eta}_m\}^2} d\Re\{\tilde{\eta}_m\}, \tag{7.27}$$

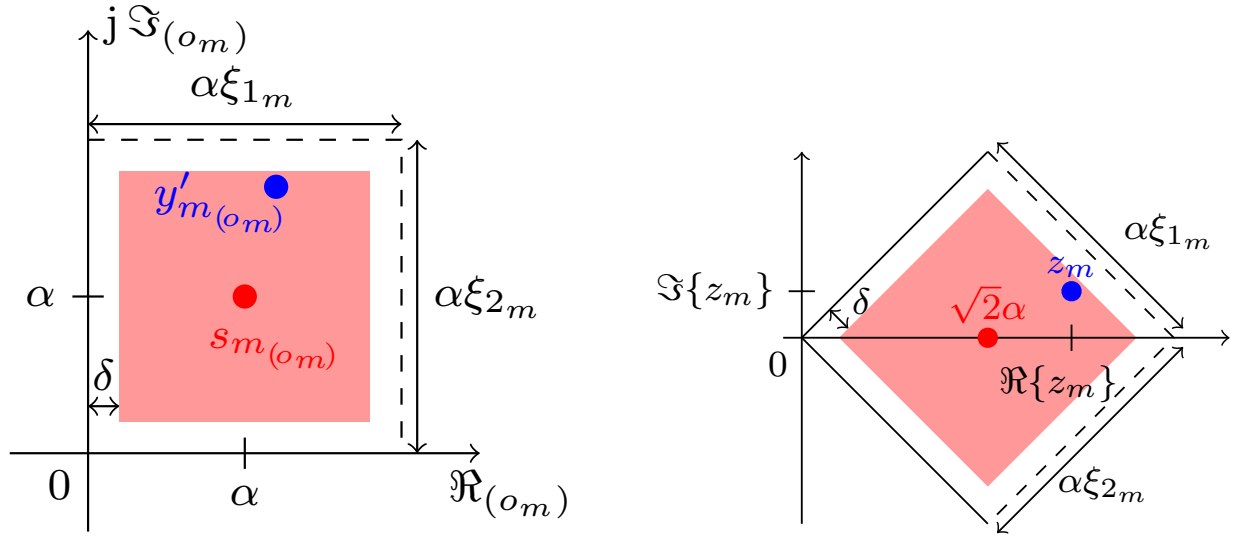
which means that maximizing δ minimizes the upper bound on the SER.

7.4 Problem Formulation for QAM Signaling

7.4.1 The Need for an Additional Degree of Freedom α

In this section, we assume that the input signals s_m , $m = 1, \dots, M$, belong to the S -QAM constellation that is defined in (4.14). As explained in Section 7.2, the safety margin δ has to be maximized such that the entries of the noiseless received signal \mathbf{y} belong to the intended SRs. The SRs in turn are determined by the constellation set \mathbb{S} and the safety margin δ . Hence, the maximal value is limited by the coordinates of the inner constellation points; that is

$$\delta \leq 1. \tag{7.28}$$



(a) Shifted coordinate system for $\Re\{s_m\} > 0$ and $\Im\{s_m\} > 0$. (b) Shifted and rotated coordinate system.

Fig. 7.4: Illustration of the QAM receiver SR in the shifted and in the shifted and rotated coordinate system : $\xi_{1/2_m} \in \{2, \infty\}$.

Independently of the available transmit power, the entries of \mathbf{y} cannot have a distance to the decision thresholds larger than 1. Hence, the available transmit power cannot be exploited to the fullest. This results already in a limitation of the problem formulation.

Thanks to the receive processing \mathbf{G} , we can introduce an additional degree of freedom α such that the entries of the received signal \mathbf{y} do not have to belong to the SRs of the set \mathcal{S} but rather to a scaled version of them; that is, the QAM constellation at each receiver gets scaled by α . The receive processing \mathbf{G} has to rescale the received signals to their nominal positions before the decision block. Thus, the constraint in (7.28) is replaced by

$$\delta \leq \alpha, \quad (7.29)$$

where α has to be jointly optimized with δ . Note that maximizing δ results in turn to maximizing α , which leads to a maximal exploitation of the available transmit power. Thus, the entries of the signal vector \mathbf{x} will get closer to the polygon corners, which decreases the variations between \mathbf{t} and \mathbf{x} .

The factor α denotes the expansion or shrinkage factor of the constellation at the receiver side depending on the available transmit power P_{tx} . As explained in Section 7.2, the optimization problem is formulated for the specific case, i.e. $P_{\text{tx}} = N$.

7.4.2 Scaled Symbol Region for QAM Signals

To describe the SRs for QAM signaling after considering α , we define a new coordinate system, that is a shifted and rotated version of the original coordinate system. First, the original receiver constellation system is shifted by o_m

$$o_m = \alpha (s_m - \text{sgn}(s_m)), \quad m = 1, \dots, M. \quad (7.30)$$

We get the following expressions for the received and the desired signal in the intermediate coordinate system depicted in Fig. 7.4b

$$y_{m(o_m)}' = y_m' - o_m \quad (7.31)$$

$$\begin{aligned} s_{m(o_m)} &= \alpha s_m - o_m \\ &\stackrel{(7.30)}{=} \alpha \operatorname{sgn}(s_m). \end{aligned} \quad (7.32)$$

Second, the intermediate coordinate system is rotated by the phase of the symbol of interest $s_{m(o_m)}$. So the signal z_m , which represents the received signal y_m' in the modified coordinate system, is given by

$$z_m = \frac{y_{m(o_m)}' s_{m(o_m)}^*}{|s_{m(o_m)}|}, \quad m = 1, \dots, M. \quad (7.33)$$

The m -th SR, $m = 1, \dots, M$, as shown in Fig. 7.4b, can be hence described by

$$\Re\{z_m\} \geq \sqrt{2}\delta, \quad (7.34)$$

$$\Re\{z_m\} \leq \sqrt{(\alpha\xi_{1_m} - \delta)^2 + (\alpha\xi_{2_m} - \delta)^2}, \quad (7.35)$$

$$|\Im\{z_m\}| \leq \left(\Re\{z_m\} - \sqrt{2}\delta \right), \quad (7.36)$$

$$\Im\{z_m\} \leq -\Re\{z_m\} + \sqrt{2}(\alpha\xi_{2_m} - \delta), \quad (7.37)$$

$$\Im\{z_m\} \geq \Re\{z_m\} - \sqrt{2}(\alpha\xi_{1_m} - \delta), \quad (7.38)$$

Note that ξ_{1_m} and $\xi_{2_m} \in \{2, \infty\}$ depending on which constellation point the symbol of interest s_m corresponds to. If s_m is one of the outer constellation points, then at least ξ_{1_m} or ξ_{2_m} must be equal to ∞ . Moreover, (7.34) and (7.35) are inherently fulfilled by (7.36), (7.37) and (7.38). In the vector notation, the vector \mathbf{z} can be expressed as

$$\mathbf{z} = \frac{1}{\sqrt{2}} \operatorname{diag}(\operatorname{sgn}(\mathbf{s}^*)) (\mathbf{H}\mathbf{x} - \alpha(\mathbf{s} - \operatorname{sgn}(\mathbf{s}))) \quad (7.39)$$

$$= \hat{\mathbf{H}}\mathbf{x} - \alpha\mathbf{c}, \quad (7.40)$$

where

$$\hat{\mathbf{H}} = \frac{1}{\sqrt{2}} \operatorname{diag}(\operatorname{sgn}(\mathbf{s}^*)) \mathbf{H}, \quad (7.41)$$

$$\mathbf{c} = \frac{1}{\sqrt{2}} \operatorname{diag}(\operatorname{sgn}(\mathbf{s}^*)) (\mathbf{s} - \operatorname{sgn}(\mathbf{s})). \quad (7.42)$$

Thus, all M SRs can be compactly described by

$$|\Im\{\mathbf{z}\}| \leq \left(\Re\{\mathbf{z}\} - \sqrt{2}\delta\mathbf{1}_M \right), \quad (7.43)$$

$$\Im\{\mathbf{z}\} \leq -\Re\{\mathbf{z}\} + \sqrt{2}(\alpha\xi_2 - \delta\mathbf{1}_M), \quad (7.44)$$

$$\Im\{\mathbf{z}\} \geq \Re\{\mathbf{z}\} - \sqrt{2}(\alpha\xi_1 - \delta\mathbf{1}_M). \quad (7.45)$$

When using the following real-valued representation

$$\mathbf{V} = [\Re\{\hat{\mathbf{H}}\} \quad -\Im\{\hat{\mathbf{H}}\}] \quad (7.46)$$

$$\mathbf{W} = [\Im\{\hat{\mathbf{H}}\} \quad \Re\{\hat{\mathbf{H}}\}], \quad (7.47)$$

the constraint in (7.45) can be rewritten as

$$\begin{bmatrix} \mathbf{W} - \mathbf{V} & \mathbf{1}_M & \Re\{\mathbf{c}\} - \Im\{\mathbf{c}\} \\ -\mathbf{W} - \mathbf{V} & \mathbf{1}_M & \Re\{\mathbf{c}\} + \Im\{\mathbf{c}\} \\ \mathbf{W} + \mathbf{V} & \mathbf{1}_M & -\Re\{\mathbf{c}\} - \Im\{\mathbf{c}\} - \sqrt{2}\boldsymbol{\xi}_2 \\ -\mathbf{W} + \mathbf{V} & \mathbf{1}_M & -\Re\{\mathbf{c}\} + \Im\{\mathbf{c}\} - \sqrt{2}\boldsymbol{\xi}_1 \end{bmatrix} \begin{bmatrix} \bar{\mathbf{x}} \\ \sqrt{2}\delta \\ \alpha \end{bmatrix} \leq \mathbf{0}_{4M}. \quad (7.48)$$

By adding the first line to the fourth and the second line to the third, the latter constraint recalculates to

$$\begin{bmatrix} \mathbf{W} - \mathbf{V} & \mathbf{1}_M & \Re\{\mathbf{c}\} - \Im\{\mathbf{c}\} \\ -\mathbf{W} - \mathbf{V} & \mathbf{1}_M & \Re\{\mathbf{c}\} + \Im\{\mathbf{c}\} \\ \mathbf{0}_{M,2N} & \sqrt{2} \mathbf{1}_M & -\boldsymbol{\xi}_2 \\ \mathbf{0}_{M,2N} & \sqrt{2} \mathbf{1}_M & -\boldsymbol{\xi}_1 \end{bmatrix} \begin{bmatrix} \bar{\mathbf{x}} \\ \sqrt{2}\delta \\ \alpha \end{bmatrix} \leq \mathbf{0}_{4M}. \quad (7.49)$$

7.4.3 Optimization Problem with the Relaxed Polygon Constraint

We are interested in maximizing the safety margin as presented in (7.8). In contrast to the PSK case, there is a constraint on δ in the QAM case, stated in (7.29), which is inherently fulfilled by (7.49). Combining (7.8) with the SRs constraint in (7.49) and the relaxed polygon constraint in (7.7), we get a linear programming problem for the design of the symbol-wise precoder for QAM signaling. The optimization problem for the case of $P_{\text{tx}} = N$ is given by

$$\begin{aligned} \max_{\mathbf{v}} [\mathbf{0}_{2N}^T \quad 1 \quad 0] \mathbf{v} \text{ s.t. } & \begin{bmatrix} \mathbf{W} - \mathbf{V} & \mathbf{1}_M & \Re\{\mathbf{c}\} - \Im\{\mathbf{c}\} \\ -\mathbf{W} - \mathbf{V} & \mathbf{1}_M & \Re\{\mathbf{c}\} + \Im\{\mathbf{c}\} \\ \mathbf{0}_{M,2N} & \sqrt{2} \mathbf{1}_M & -\boldsymbol{\xi}_2 \\ \mathbf{0}_{M,2N} & \sqrt{2} \mathbf{1}_M & -\boldsymbol{\xi}_1 \\ \mathbf{E} & \mathbf{0}_{N(Q-4)} & \mathbf{0}_{N(Q-4)} \end{bmatrix} \mathbf{v} \leq \begin{bmatrix} \mathbf{0}_{4M} \\ \cos\left(\frac{\pi}{Q}\right) \mathbf{1}_{N(Q-4)} \end{bmatrix} \\ \text{and } & \begin{bmatrix} -\cos\left(\frac{\pi}{Q}\right) \mathbf{1}_{2N} \\ 0 \\ 0 \end{bmatrix} \leq \mathbf{v} \leq \begin{bmatrix} \cos\left(\frac{\pi}{Q}\right) \mathbf{1}_{2N} \\ \infty \\ \infty \end{bmatrix}, \end{aligned} \quad (7.50)$$

where $\mathbf{v}^T = [\bar{\mathbf{x}}^T \quad \sqrt{2}\delta \quad \alpha]$. In order to solve the optimization problem for different P_{tx} values, it is sufficient to scale the optimal solution of (7.21) by $\frac{P_{\text{tx}}}{N}$ due to the linearity of the optimization problem.

Again the optimized vector $\mathbf{x} \in \mathbb{X}^N$ goes through the quantizer, as stated in (4.1), to obtain the transmit vector \mathbf{t} .

7.4.4 Safety Margin δ for QAM Signals

For the QAM case, the safety margin δ in (7.8) can be expressed as

$$\delta = \min_m \min \left(\frac{1}{\sqrt{2}} (\Re\{\mathbf{z}\} - |\Im\{\mathbf{z}\}|), \alpha \boldsymbol{\xi}_1 - \frac{1}{\sqrt{2}} (\Re\{\mathbf{z}\} - \Im\{\mathbf{z}\}), \alpha \boldsymbol{\xi}_2 - \frac{1}{\sqrt{2}} (\Re\{\mathbf{z}\} + \Im\{\mathbf{z}\}) \right), \quad (7.51)$$

where the operator $|\bullet|$ is applied element-wise to the entries of $\Im\{\mathbf{z}\}$.

7.4.5 Interpretation of the Safety Margin δ for QAM Signals

Again we consider the relaxed problem; that is the quantization is omitted and we consider the received signal \mathbf{y}' instead of \mathbf{y} .

7.4.5.1 Safety Margin vs. Receiver SNR

The receiver SNR at the m -th user can be approximated by

$$\text{SNR}_m \approx \frac{1}{N_s} \sum_{t=1}^{N_s} (\alpha^{(t)})^2 \sigma_s^2. \quad (7.52)$$

Since $\delta \leq \alpha$, we get

$$\text{SNR}_m \geq \frac{1}{N_s} \sum_{t=1}^{N_s} (\delta^{(t)})^2 \sigma_s^2. \quad (7.53)$$

Thus, maximizing δ results in maximizing the lower bound of the receiver SNR.

7.4.5.2 Safety Margin vs. SER

The statistics of the AWGN in the modified coordinate system do not change, since it just consists of a rotation. The shift is applied only to the noiseless received signal \mathbf{y} . We denote the modified noise vector by $\tilde{\boldsymbol{\eta}}$, where $\tilde{\boldsymbol{\eta}} \sim \mathcal{CN}_{\mathcal{C}}(\mathbf{0}_M, \mathbf{I}_M)$. The SER represents the probability that the entries of $\mathbf{z} + \tilde{\boldsymbol{\eta}}$ lie outside the thresholds that are represented by the black lines in Fig. 7.4b. We consider the worst case, where the decision region is bounded from all four sides. Additionally, we assume that the distribution of the entries of \mathbf{z} is uniform. Therefore, the SER of the m -th user for the t -th transmission is upper bounded by

$$\begin{aligned} \text{SER}_m^{(t)} &\leq \Pr \left(\left(0 \leq \Re\{z_m^{(t)}\} + \Re\{\tilde{n}_m\} \leq 2\sqrt{2}\alpha^{(t)} \right) \wedge \left(\Im\{\tilde{n}_m\} \geq \sqrt{2}\delta_m^{(t)} + \Re\{\tilde{n}_m\} \right) \right) \\ &\quad + \Pr \left(\left(0 \leq \Re\{z_m^{(t)}\} + \Re\{\tilde{n}_m\} \leq 2\sqrt{2}\alpha^{(t)} \right) \wedge \left(\Im\{\tilde{n}_m\} \leq -\sqrt{2}\delta_m^{(t)} - \Re\{\tilde{n}_m\} \right) \right) \\ &\quad + \Pr \left(\left(0 \leq \Re\{z_m^{(t)}\} + \Re\{\tilde{n}_m\} \leq 2\sqrt{2}\alpha^{(t)} \right) \wedge \left(\Im\{\tilde{n}_m\} \geq 2\sqrt{2}\alpha^{(t)} - \sqrt{2}\delta_m^{(t)} - \Re\{\tilde{n}_m\} \right) \right) \\ &\quad + \Pr \left(\left(0 \leq \Re\{z_m^{(t)}\} + \Re\{\tilde{n}_m\} \leq 2\sqrt{2}\alpha^{(t)} \right) \wedge \left(\Im\{\tilde{n}_m\} \leq -2\sqrt{2}\alpha + \sqrt{2}\delta_m^{(t)} + \Re\{\tilde{n}_m\} \right) \right) \\ &\quad + 1 - \Pr(0 \leq \Re\{z_m^{(t)}\} + \Re\{\tilde{n}_m\} \leq 2\sqrt{2}\alpha^{(t)}) \end{aligned} \quad (7.54)$$

$$\begin{aligned} &= 1 + \int_{-\Re\{z_m^{(t)}\}}^{2\sqrt{2}\alpha^{(t)} - \Re\{z_m^{(t)}\}} \left(1 - \left(\text{erf}\left(\sqrt{2}\delta_m^{(t)} + \Re\{\tilde{n}_m\}\right) + \text{erf}\left(2\sqrt{2}\alpha^{(t)} - \sqrt{2}\delta_m^{(t)} - \Re\{\tilde{n}_m\}\right) \right) \right) \\ &\quad \frac{1}{\sqrt{\pi}} e^{-\Re\{\tilde{\eta}_m\}^2} d\Re\{\tilde{\eta}_m\} \end{aligned} \quad (7.55)$$

Since the erf function is monotonically increasing and it holds that $\delta^{(t)} \leq \delta_m^{(t)} \leq \alpha^{(t)}$, $\forall m$ and $\sqrt{2}\delta_m^{(t)} \leq \Re\{z_m^{(t)}\} \leq 2\sqrt{2}\alpha^{(t)} - \sqrt{2}\delta_m^{(t)}$, we obtain

$$\text{SER}_m^{(t)} \leq 1 + \int_{-2\sqrt{2}\alpha^{(t)} + \sqrt{2}\delta_m^{(t)}}^{2\sqrt{2}\alpha^{(t)} - \sqrt{2}\delta_m^{(t)}} \left(1 - 2\text{erf}\left(\sqrt{2}\delta_m^{(t)} + \Re\{\tilde{\eta}_m\}\right) \right) \frac{1}{\sqrt{\pi}} e^{-\Re\{\tilde{\eta}_m\}^2} d\Re\{\tilde{\eta}_m\} \quad (7.56)$$

$$\leq 1 + \int_{-2\sqrt{2}\alpha^{(t)} + \sqrt{2}\delta^{(t)}}^{2\sqrt{2}\alpha^{(t)} - \sqrt{2}\delta^{(t)}} \left(1 - 2\text{erf}\left(\sqrt{2}\delta^{(t)} + \Re\{\tilde{\eta}_m\}\right) \right) \frac{1}{\sqrt{\pi}} e^{-\Re\{\tilde{\eta}_m\}^2} d\Re\{\tilde{\eta}_m\}. \quad (7.57)$$

Thus, the SER at the m -th user averaged over a transmission block of length N_s is upper bounded by

$$\text{SER}_m \leq 1 + \frac{1}{N_s} \sum_{t=1}^{N_s} \int_{-2\sqrt{2}\alpha^{(t)} + \sqrt{2}\delta^{(t)}}^{2\sqrt{2}\alpha^{(t)} - \sqrt{2}\delta^{(t)}} \left(1 - 2\text{erf}\left(\sqrt{2}\delta^{(t)} + \Re\{\tilde{\eta}_m\}\right)\right) \frac{1}{\sqrt{\pi}} e^{-\Re\{\tilde{\eta}_m\}^2} d\Re\{\tilde{\eta}_m\}. \quad (7.58)$$

Since erf is a monotonically increasing function, larger values of δ make the upper bound of the SER decrease. Note that larger values of δ lead inherently to larger values of α .

7.4.6 Symbol-wise Processing vs. Block-wise Processing

One might ask why we opt for symbol-wise processing and not block-wise processing. The factor α cannot be communicated to the receiver and hence has to be estimated. The estimation is based on averaging over a block of T received signals. Thus, one expects that the design of α at the transmitter has to be computed for the same block length, i.e. $B = T$. However, fixing α for a certain block length B means that B vectors \mathbf{x} have to be designed jointly with a single factor α instead of having a distinct factor α for every vector \mathbf{x} . Additionally, the joint optimization of B vectors results in a higher-dimensional linear programming problem, where the number of inequalities is increased by a factor of B . Hence, block-wise processing not only increases the computational complexity of the problem as can be deduced from Section 7.5 but also reduces the degrees of freedom of the optimization problem at the transmitter. This leads to the entries of the vector \mathbf{x} moving farther from the polygon corners, thus increasing the quantization distortions. This effect is illustrated in Fig. 7.5, where the entries of $\mathbf{e}_m^T \mathbf{H}\mathbf{x}$, $\mathbf{e}_m^T \mathbf{H}\mathbf{t}$ and $\frac{1}{\alpha} \mathbf{e}_m^T \mathbf{H}\mathbf{t}$ of an arbitrary user m are obtained by transmitting 1024 16-QAM signal vectors through an i.i.d. channel of coefficients $h_{mn} \sim \mathcal{CN}_{\mathbb{C}}(0, 1)$, $n = 1, \dots, N$, $m = 1, \dots, M$, where $N = 64$, $M = 8$ and $Q = 4$. The optimization is computed for both symbol-wise processing, i.e. $B = 1$, and block-wise processing with $B = 4$. As can be deduced from the plots, block-wise processing leads to a larger safety margin with the relaxed vector \mathbf{x} . However, after applying the quantization this gain is lost and the symbol-wise processing is more robust against the quantization operation. This can be further explained by the results in Table 7.1, which shows $\mathbb{E} \left[\frac{\|\mathbf{t} - \mathbf{x}\|_1}{N} \right]$, the percentage of entries of \mathbf{x} that are distorted due to the quantization and the MSE between \mathbf{t} and \mathbf{x} . We see that increasing B significantly increases the quantization distortion. Therefore, the

	$B = 1$	$B = 4$
$\mathbb{E} \left[\frac{\ \mathbf{t} - \mathbf{x}\ _1}{N} \right]$	0.2176	0.4432
$\mathbb{E} [\ \mathbf{t} - \mathbf{x}\ _2^2]$	2.5458	12.6429

Table 7.1: Quantization distortion vs. B , © 2018 IEEE.

symbol-wise processing is chosen in this contribution, i.e. an optimal value of α is designed for each vector \mathbf{x} .

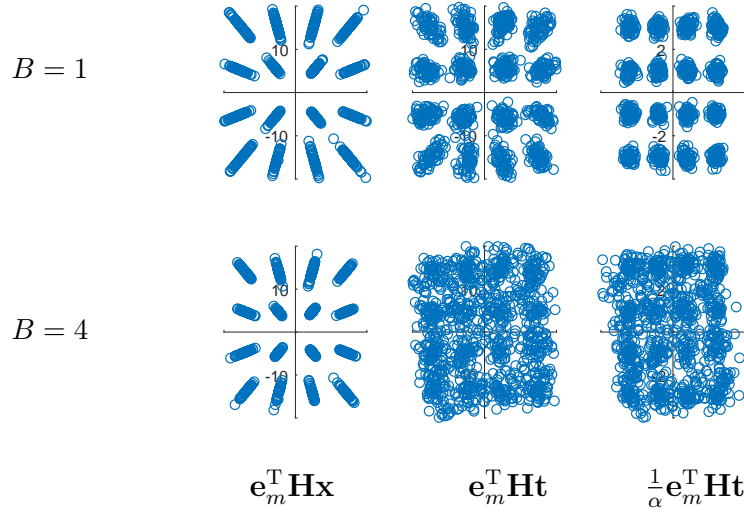


Fig. 7.5: The noiseless received symbols at one arbitrary user m for an arbitrary i.i.d. channel realization with $N = 64$, $M = 8$ and $Q = 4$, © 2018 IEEE.

7.4.7 One Joint α vs. M Distinct α 's for M Users

Symbol-wise transmit processing followed by block-wise receive processing is reliable only if the obtained values of α , i.e. $\alpha^{(i)}$, $i = 1, \dots, T$, do not vary much from one vector $\mathbf{x}^{(i)}$ to another. Otherwise, estimating the mean value of α at the receiver would not be sufficient for correct detection. This explains why we choose one joint α for all users. If a different value α_m per user is chosen, this will result in more degrees of freedom and the values α_m , $m = 1, \dots, M$, would fluctuate much more from one vector \mathbf{x} to another, which worsens the estimation result at the receiver. For a large number of users, the jointly designed α will not vary much, since the norm of the input vector \mathbf{s} will not fluctuate much from one realization to another. This behavior is explained as follows.

For a given channel realization, T symbols are transmitted and the MSM method is applied T times to get in the optimal case

$$\begin{aligned} \mathbf{H}\mathbf{x}^{(1)} &= \alpha^{(1)}\mathbf{s}^{(1)} \\ &\vdots \\ \mathbf{H}\mathbf{x}^{(T)} &= \alpha^{(T)}\mathbf{s}^{(T)}. \end{aligned} \quad (7.59)$$

The factor $\alpha^{(i)}$ can be then expressed as

$$\alpha^{(i)} = \frac{\|\mathbf{H}\mathbf{x}^{(i)}\|_2}{\|\mathbf{s}^{(i)}\|_2}, \quad i = 1, \dots, T \quad (7.60)$$

and hence is upper bounded by

$$\alpha^{(i)} \leq \frac{\|\mathbf{H}\|_2 \|\mathbf{x}^{(i)}\|_2}{\|\mathbf{s}^{(i)}\|_2} \leq \frac{\|\mathbf{H}\|_2 \|\mathbf{t}^{(i)}\|_2}{\|\mathbf{s}^{(i)}\|_2}. \quad (7.61)$$

Due to the constant envelope property, we get $\|\mathbf{t}^{(i)}\|_2 = \sqrt{N}$, $\forall i = 1, \dots, T$. Thus, the fluctuation of the upper bound of $\alpha^{(i)}$ is determined by $\|\mathbf{s}^{(i)}\|_2$. Since the entries of the

vector \mathbf{s} are i.i.d, the fluctuation of $\|\mathbf{s}^{(i)}\|_2$ from one realization to another and for $S > 4$ diminishes by increasing its dimension, i.e. the number of users M . We recall that the MSM method aims to maximize the safety margin δ that is upper bounded by α as stated in (7.29). Consequently, the MSM method will maximize α as well. As explained before, the upper bound of α fluctuates less from one signal realization to another for a large number of users and thus the optimal value of $\alpha^{(i)}$ has only small fluctuations in such cases.

To make the received signals belong to the nominal constellation, the receive processing \mathbf{G} , as explained in Section 4.7, is applied. For each time instant, the coefficients g_m , $m = 1, \dots, M$ should be equal to the inverse of α . However, the computation of the coefficients g_m is based on blind estimation over a block length T . Therefore, the resulting coefficients g_m are ideally equal to the inverse of the mean value of α , i.e. $g_m \approx 1/\mathbb{E}[\alpha]$. Hence, the designed SR at the transmitter in the noise-free case, which is illustrated in Fig. 7.4a, is scaled by $1/\mathbb{E}[\alpha]$ at the receiver before the decision operation. For the i -th transmission, the center of the obtained nominal SR_m corresponds to $s_m^{(i)}\alpha^{(i)}/\mathbb{E}[\alpha]$. If $\alpha^{(i)}$ corresponds to $\mathbb{E}[\alpha]$, the nominal SR_m after the receive filter \mathbf{G} is situated in the center of the corresponding nominal decision region, that is $s_m^{(i)}$. However, if $\alpha^{(i)}$ is smaller than $\mathbb{E}[\alpha]$, the center of the nominal SR_m shifts towards the left lower corner of the decision region of $s_m^{(i)}$ and might pass over it. Analogously, if $\alpha^{(i)}$ is larger than $\mathbb{E}[\alpha]$, the center of the nominal SR_m shifts towards the right upper corner of the decision region and might pass over it, too. To illustrate the passing over phenomenon, we introduce the nominal safety margin δ_{nom} , which is given by

$$\begin{aligned} \delta_{\text{nom}} &= \min \left(\frac{2\alpha_{\min} + \delta_{\alpha_{\min}}}{\mathbb{E}[\alpha]} - 2, 2 - \frac{2\alpha_{\max} - \delta_{\alpha_{\max}}}{\mathbb{E}[\alpha]} \right) \\ &= \min \left(\frac{\delta_{\alpha_{\min}}}{\mathbb{E}[\alpha]} - 2 \frac{\mathbb{E}[\alpha] - \alpha_{\min}}{\mathbb{E}[\alpha]}, \frac{\delta_{\alpha_{\max}}}{\mathbb{E}[\alpha]} - 2 \frac{\alpha_{\max} - \mathbb{E}[\alpha]}{\mathbb{E}[\alpha]} \right) \\ &\geq \frac{1}{\mathbb{E}[\alpha]} (\min(\delta_{\alpha_{\min}}, \delta_{\alpha_{\max}}) - 2 \max(\mathbb{E}[\alpha] - \alpha_{\min}, \alpha_{\max} - \mathbb{E}[\alpha])) \\ &= \frac{1}{\mathbb{E}[\alpha]} \delta_{\alpha_{\min}} - 2\Delta\alpha \\ \delta_{\text{nom}} &\geq \frac{\min_i \delta^{(i)}}{\sum_{i=1}^T \alpha^{(i)}/T} - 2\Delta\alpha, \end{aligned} \quad (7.62)$$

where α_{\min} , α_{\max} , $\delta_{\alpha_{\min}}$ and $\delta_{\alpha_{\max}}$ denote the minimal and maximal value of $\alpha^{(i)}$ and their corresponding safety margins, respectively. $\Delta\alpha$ represents the maximal relative fluctuation of α . For the optimization with one single α , it is defined as

$$\Delta\alpha = \max \left(1 - \frac{\min_i \alpha^{(i)}}{\sum_{i=1}^T \alpha^{(i)}/T}, \frac{\max_i \alpha^{(i)}}{\sum_{i=1}^T \alpha^{(i)}/T} - 1 \right), \quad (7.63)$$

whereas it is expressed as

$$\Delta\alpha = \max_m \max \left(1 - \frac{\min_i \alpha_m^{(i)}}{\sum_{i=1}^T \alpha_m^{(i)}/T}, \frac{\max_i \alpha_m^{(i)}}{\sum_{i=1}^T \alpha_m^{(i)}/T} - 1 \right) \quad (7.64)$$

for the optimization with M distinct α 's for the different users. In the case that δ_{nom} gets negative values due to the large fluctuation of α , i.e. $\Delta\alpha$, it implies that the nominal SRs overlap in the noise-free case. To ensure a large positive number for δ_{nom} , $\Delta\alpha$ must be kept as small as possible. In the best case, $\Delta\alpha$ must be equal to zero.

From (7.62), we can conclude that smaller fluctuations of α lead to larger values of δ_{nom} . This observation is justified by numerical results, where the MSM optimization is run for $T = 128$ 16-QAM symbols, for 100 i.i.d. channel realizations and for $Q = 4$. The obtained values for $\Delta\alpha$ and δ_{nom} are averaged over the channels and shown in Table 7.2 and Table 7.3. We can deduce from Table 7.2, where a common α is designed for all users, that the values

N	64		200	
M	$\Delta\alpha$	δ_{nom}	$\Delta\alpha$	δ_{nom}
2	1.3	-0.3	1.2	-0.2
8	0.5	0.2	0.4	0.3
14	0.3	0.2	0.3	0.4

Table 7.2: $\Delta\alpha$ and δ_{nom} vs. N and M : single α , © 2018 IEEE.

of α fluctuate less and hence δ_{nom} increases by increasing the number of users. However, no monotonic behavior of the fluctuations is noticed in Table 7.3, which shows relatively larger values of $\Delta\alpha$ and hence smaller values of δ_{nom} compared to Table 7.2.

7.5 Computational Complexity of MSM

7.5.1 On the Computational Complexity of General Linear Programming Problems

In this section, we study the computational complexity of the simplex method for a general linear programming problem with bounded variables in inequality form:

$$\begin{aligned} \max_{\mathbf{x}} \mathbf{c}^T \mathbf{x} \text{ s.t. } \mathbf{A}\mathbf{x} \leq \mathbf{b} \\ \text{and } \mathbf{l} \leq \mathbf{x} \leq \mathbf{u}, \end{aligned} \quad (7.65)$$

where \mathbf{c} , \mathbf{x} , \mathbf{l} and $\mathbf{u} \in \mathbb{R}^n$, $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$.

N	64		200	
M	$\Delta\alpha$	δ_{nom}	$\Delta\alpha$	δ_{nom}
2	1.2	-0.2	1.1	-0.1
8	1.1	-0.8	0.5	0.3
14	1.6	-2	0.7	-0.4

Table 7.3: $\Delta\alpha$ and δ_{nom} vs. N and M : M α 's, © 2018 IEEE.

First, we have to make sure that the entries of \mathbf{b} are non-negative. To this end, we change the signs of the inequalities that correspond to negative entries in \mathbf{b} . So we get

$$\begin{aligned} \min_{\mathbf{x}} \mathbf{c}^T \mathbf{x} \text{ s.t. } \tilde{\mathbf{A}} \mathbf{x} &\begin{cases} \geq \tilde{\mathbf{b}} \\ \leq \tilde{\mathbf{b}} \end{cases} \\ \text{and } \mathbf{l} &\leq \mathbf{x} \leq \mathbf{u}, \end{aligned} \quad (7.66)$$

where $\tilde{\mathbf{b}} \in \mathbb{R}_+^m$ and some inequalities hold with the sign \leq and others with the sign \geq .

Second, the linear programming problem is transformed to the canonical form by introducing m slack and surplus variables denoted by \mathbf{x}_s . Additionally, a artificial variables denoted by \mathbf{x}_a , with $0 \leq a \leq m$, are added to set up an initial feasible solution [72]. The equivalent enlarged problem reads as

$$\begin{aligned} \min_{\bar{\mathbf{x}}} \bar{\mathbf{c}}^T \bar{\mathbf{x}} \text{ s.t. } \bar{\mathbf{A}} \bar{\mathbf{x}} &= \tilde{\mathbf{b}} \\ \text{and } \bar{\mathbf{l}} &\leq \bar{\mathbf{x}} \leq \bar{\mathbf{u}}, \end{aligned} \quad (7.67)$$

where $\bar{\mathbf{A}} = [\tilde{\mathbf{A}} \ \mathbf{A}_s \ \mathbf{I}_a] \in \mathbb{R}^{m \times (n+m+a)}$, $\bar{\mathbf{x}}^T = [\mathbf{x}^T \ \mathbf{x}_s^T \ \mathbf{x}_a^T] \in \mathbb{R}^{n+m+a}$, $\bar{\mathbf{l}}^T = [\mathbf{l}^T \ \mathbf{0}_{m+a}^T]$ and $\bar{\mathbf{u}}^T = [\mathbf{u}^T \ \infty \mathbf{1}_{m+a}^T]$. The matrix \mathbf{A}_s is a diagonal matrix with entries equal to 1 or -1 depending on whether the inequality sign in (7.66) is \leq or \geq , respectively. The number a of artificial variables is defined by the number of negative entries in \mathbf{A}_s , such that the concatenation of m columns from $[\mathbf{A}_s \ \mathbf{I}_a]$ can construct the identity matrix \mathbf{I}_m . For the special case $\mathbf{b} = \tilde{\mathbf{b}}$, i.e. the entries of \mathbf{b} are non-negative, $\mathbf{A}_s = \mathbf{I}_m$. Hence, no artificial variables are needed, i.e. $a = 0$.

With the use of the simplex method to solve (7.67), the number of operations (multiplication and addition pairs) on each iteration is given by, [72, p.83],

$$3m \quad \text{or} \quad (m+1)(n+a+1) + 2m, \quad (7.68)$$

depending on whether pivoting is required or not. According to [72, p.86], in most iterations no pivoting is required and hence less computation is needed.

7.5.2 Computational Complexity of MSM for PSK Signaling

As can be seen from (7.21), there are $m = 2M + N(Q-4)$ inequalities and $n = 2N + 1$ variables. The number a of artificial variables reduces to 0, since the vector $\mathbf{b}^T = [\mathbf{0}_{2M}^T \ \cos(\frac{\pi}{Q}) \mathbf{1}_{N(Q-4)}^T]$ has only non-negative entries. Thus, the number of operations (multiplication and addition pairs) on each iteration calculates in this case to

$$6M + 3(Q-4)N, \quad (7.69)$$

or

$$2N + 4MN + 8M + 2 + (Q-4)(2N^2 + 4N). \quad (7.70)$$

For massive MIMO systems, where $N \gg M$, the number of floating-point operations for each iteration is on the order of

$$\#\text{FLOPs}_{\text{PSK}} = \mathcal{O}(4MN + (Q-4)(2N^2 + 4N)). \quad (7.71)$$

For the special case of one-bit quantization, i.e. $Q = 4$, the computational complexity of the MSM method is linear in N and M .

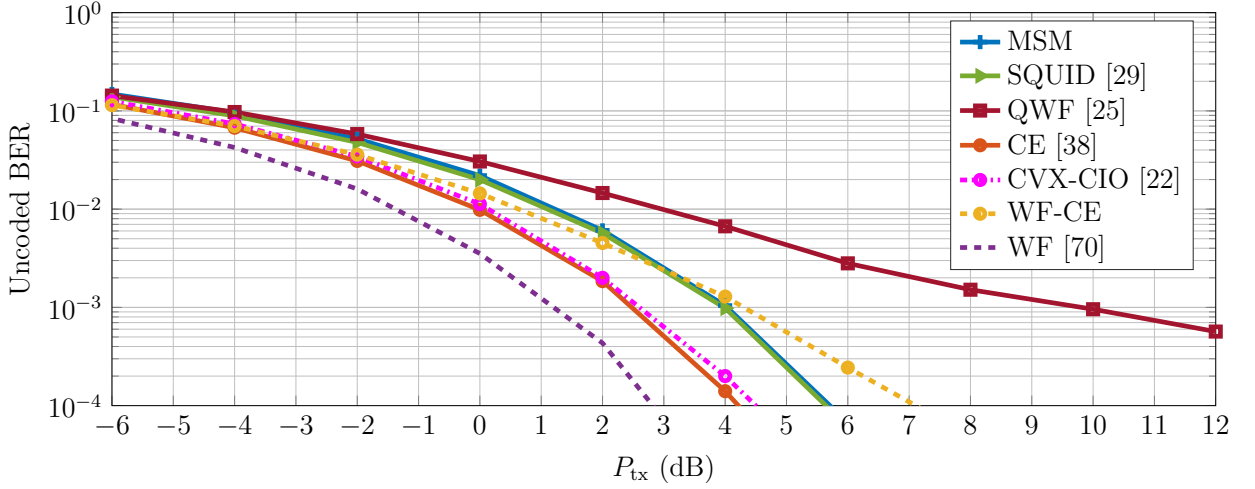


Fig. 7.6: Unencoded BER performance for a MU-MIMO system with $N = 64$ and $M = 8$ for different precoding designs and 4-PSK signaling, © 2018 IEEE.

7.5.3 Computational Complexity of MSM for QAM Signaling

From (7.50), we have $m = 4M + N(Q - 4)$ inequalities and $n = 2N + 2$ variables. The number a of artificial variables reduces to 0, since the vector $\mathbf{b}^T = [\mathbf{0}_{4M+1}^T \quad \cos(\frac{\pi}{Q}) \mathbf{1}_{N(Q-4)}^T]$ has only non-negative entries. Thus, the number of operations (multiplication and addition pairs) on each iteration calculates in this case to

$$12M + 3(Q - 4)N, \quad (7.72)$$

or

$$2N + 8MN + 20M + 3 + (Q - 4)(2N^2 + 5N). \quad (7.73)$$

For massive MIMO systems, the number of floating-point operations for each iteration is on the order of

$$\#\text{FLOPs}_{\text{QAM}} = \mathcal{O}(8MN + (Q - 4)(2N^2 + 5N)). \quad (7.74)$$

For the special case of one-bit quantization, i.e. $Q = 4$, the complexity is linear in N and M . Note that the sparsity of \mathbf{E} can be exploited by deploying the revised simplex method to reduce the number of required operations in the case of $Q > 4$ [72, p.89].

7.6 Simulation Results

For the simulations, we assume a BS with $N = 64$ antennas serving $M = 8$ single-antenna users. The channel \mathbf{H} is composed of i.i.d. Gaussian random variables with zero-mean and unit variance. The numerical results are obtained with Monte Carlo simulations of 100 independent channel realizations. The additive noise is also i.i.d. with variance one at each antenna. The performance metric is the unencoded BER averaged over the single-antenna users. For the blind estimation of the coefficients g_m we use a block length of $T = 128$.

In the first simulation set, depicted in Fig. 7.6, we assume full CSI, choose QPSK modulation and compare the uncoded BER as a function of the transmit power P_{tx} for the following precoders:

- The proposed MSM method with $Q = 4$.
- The SQUID precoder presented in [29] with $Q = 4$, where the precoder design criterion is the symbol-wise MSE between \mathbf{u} and \mathbf{s} under a quantization constraint. The latter is equivalent to the QCE constraint for the special case $Q = 4$. The SQUID precoder is a semi-definite relaxation based algorithm.
- The quantized WF precoder denoted by "QWF" from [25] with $Q = 4$. This precoder design is based on linearizing the quantizer and considering the resulting quantization noise as additive Gaussian noise.
- The CE precoder presented in [38] denoted by "CE [38]", with $Q = \infty$, where the symbol-wise MSE between \mathbf{y} and a scaled version of \mathbf{s} is minimized under the CE constraint. The scaling factor that is applied to \mathbf{s} is jointly optimized.
- The CE precoder from [22] denoted by CVX-CIO that aims at maximizing the constructive interference under the CE constraint.
- The WF precoder followed by the CE quantizer with $Q = \infty$ denoted by "WF-CE", and
- The WF precoder in the ideal case denoted by "The ideal WF" from [70], where neither quantization nor the CE constraint is applied to the transmit signal.

It can be seen that the CE constraint leads to a loss of almost 2 dB at a BER of 10^{-2} compared to the ideal WF and a loss of less than 1.5 dB when using the unquantized symbol-wise precoders proposed in [38] and in [22]. The one-bit quantization, which represents the QCE case of $Q = 4$, leads to more losses that depend on the precoder design. With the use of the linear precoder QWF a loss of more than 4 dB at a BER of 10^{-2} is noticed. However, the non-linear precoders MSM and SQUID improve the performance drastically and show a loss of slightly more than 2 dB compared to the ideal case at the cost of higher computational complexity. Nevertheless, the proposed MSM method is more efficient than SQUID as it is based on a purely linear programming formulation.

In the second simulation set, depicted in Fig. 7.7, the uncoded BER is plotted as a function of the transmit power P_{tx} using the MSM precoder for different modulation schemes, for the i.i.d. and the mmW ($N_{\text{cl}} = 2$ and $N_{\text{ray}} = 10$) channel models and for two different values of Q : $Q = 4$ and $Q = 8$. Higher values of Q are omitted since the obtained results do not differ much from the case of $Q = 8$. In addition, it is beneficial in terms of computational complexity and power consumption to keep Q as small as possible. The performance obtained by the MSM method is compared to the performance of the ideal linear system. It can be deduced from the plots that the non-linear approach improves the performance much better than the linear approaches discussed in Chapter 5. As expected, the higher the number of symbols in the modulation scheme, the higher the BER for a given P_{tx} value. However, the increase of the DAC resolution q and thus the resulting increase in Q leads to a performance improvement, which depends on the modulation scheme.

In the third simulation set, the same comparison is conducted, but with doubling the number of antennas of the quantized system. The results are plotted in Fig. 7.8. It can be concluded that the performance of the ideal system can be achieved with a QCE, which has less than the double number of antennas. When the quantization resolution is increased, less antennas are required to achieve the same performance as the one of the ideal linear system.

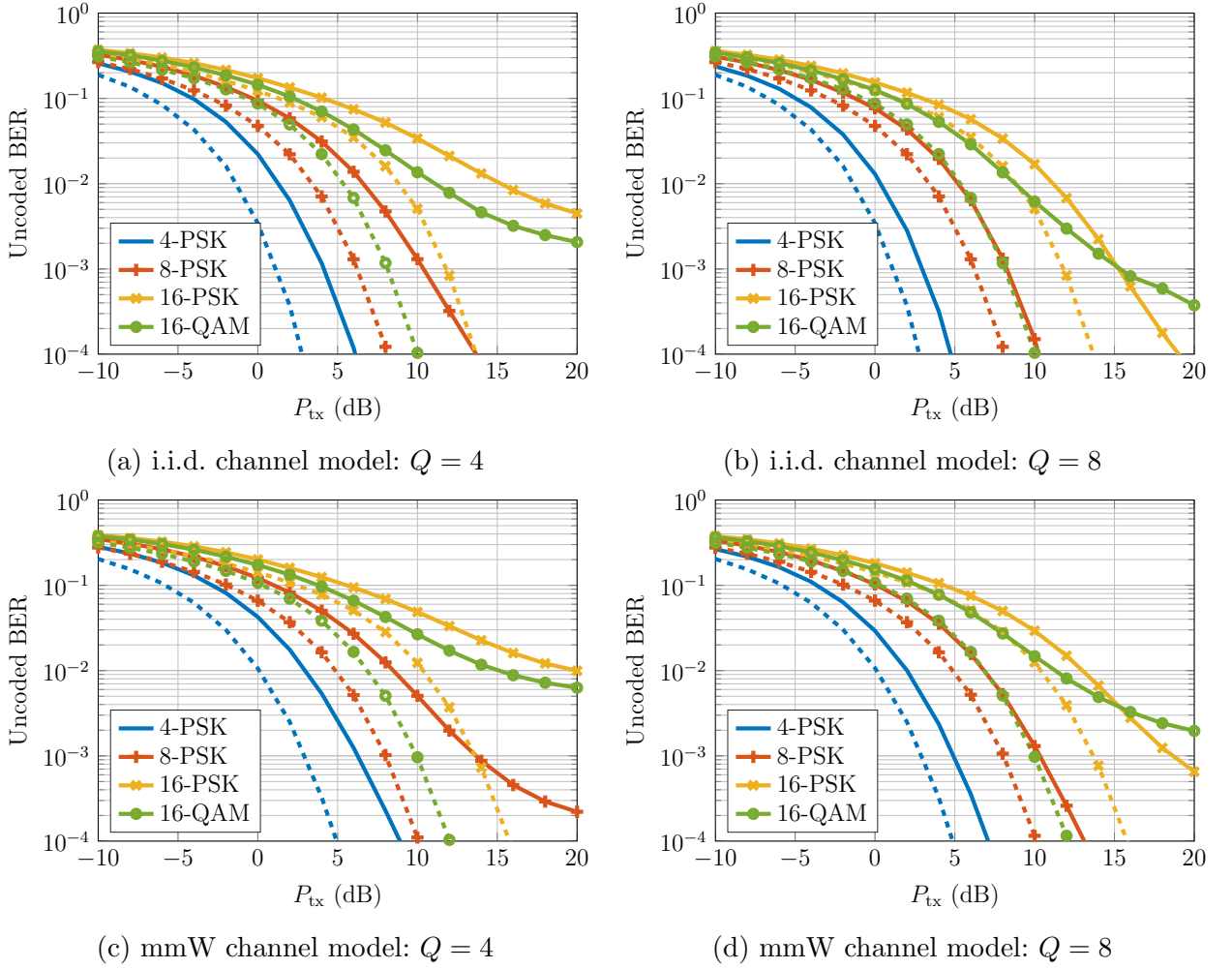


Fig. 7.7: Uncoded BER performance for a MU-MIMO system with $N = 64$ and $M = 8$ for different modulation schemes: MSM (solid lines), the ideal WF (dashed lines).

Since the optimization problem in [22] has some similarities with our proposed MSM, we compare the uncoded BER performance for both designs in Fig. 7.9. In our simulation, we pass the entries of \mathbf{x} obtained by the CVX-CIO method through the CE quantizer to get QCE signals. Additionally, we introduce the method denoted by CVX-CIO-noCE that is the same as CVX-CIO with an instantaneous power constraint instead of the CE constraint. As can be seen from the results, CVX-CIO and CVX-CIO-noCE do not perform optimally under the constraint of QCE transmit signals. However, the loss compared to MSM reduces when the quantization resolution increases.

In the last simulation set, we counted the average number of iterations required by the MSM precoder. The results are summarized in Table 7.4, where we observe that around 50 iterations are required for all modulation schemes for $Q = 4$ and more than 100 iterations for $Q > 4$.

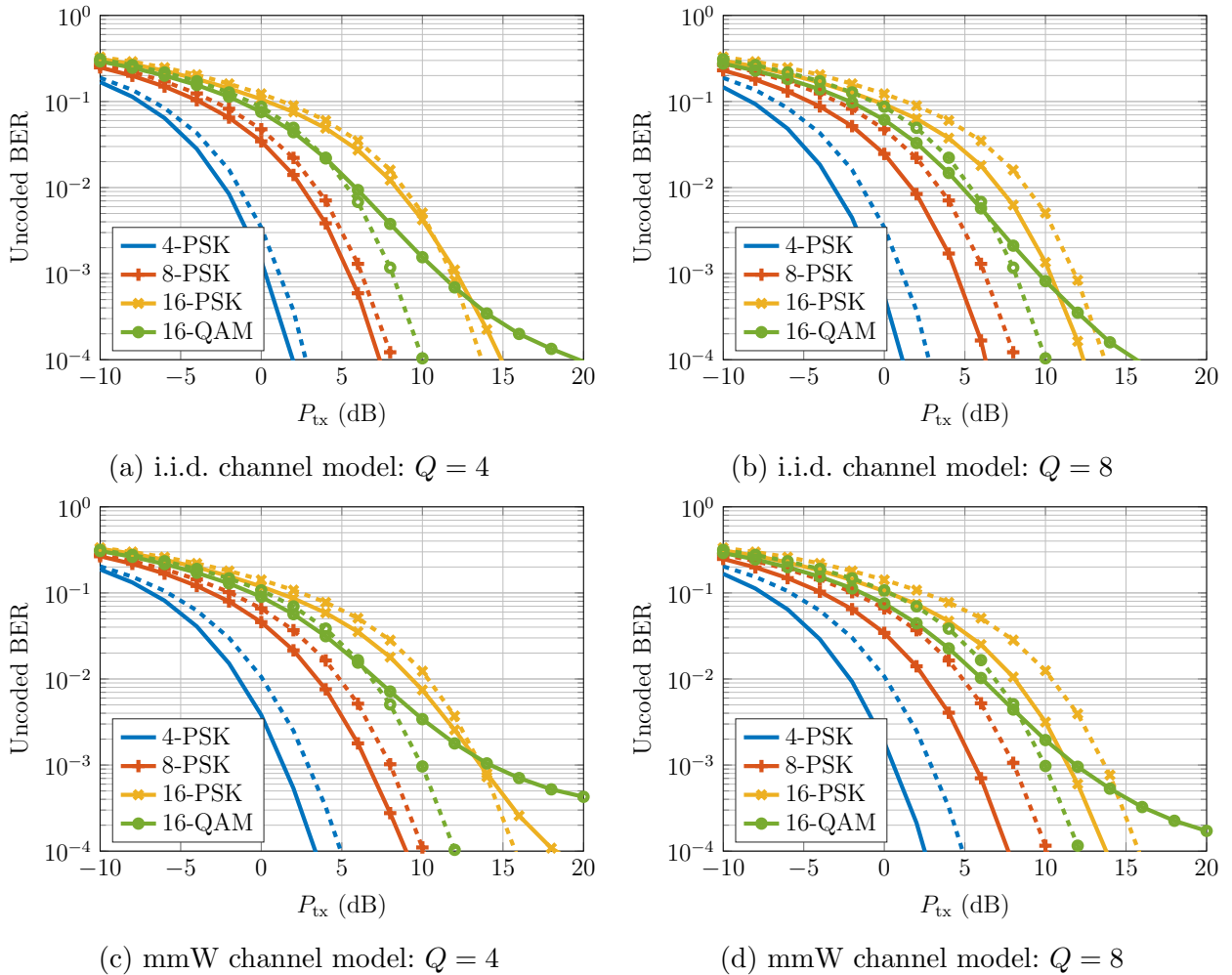


Fig. 7.8: Unencoded BER performance for a MU-MIMO system with $M = 8$ for different modulation schemes: MSM with $N = 128$ (solid lines), the ideal WF with $N = 64$ (dashed lines).

Nb. of iter	$Q = 4$	$Q = 8$	$Q = 16$
4-PSK	45.77	121.05	187.63
8-PSK	50.15	123.91	191.55
16-PSK	54.94	128.74	199.61
16-QAM	43.25	120.42	187.32
64-QAM	43.04	120.30	188.30

Table 7.4: Average number of iterations of MSM for the i.i.d. channel model with $N = 64$ and $M = 8$, © 2018 IEEE.

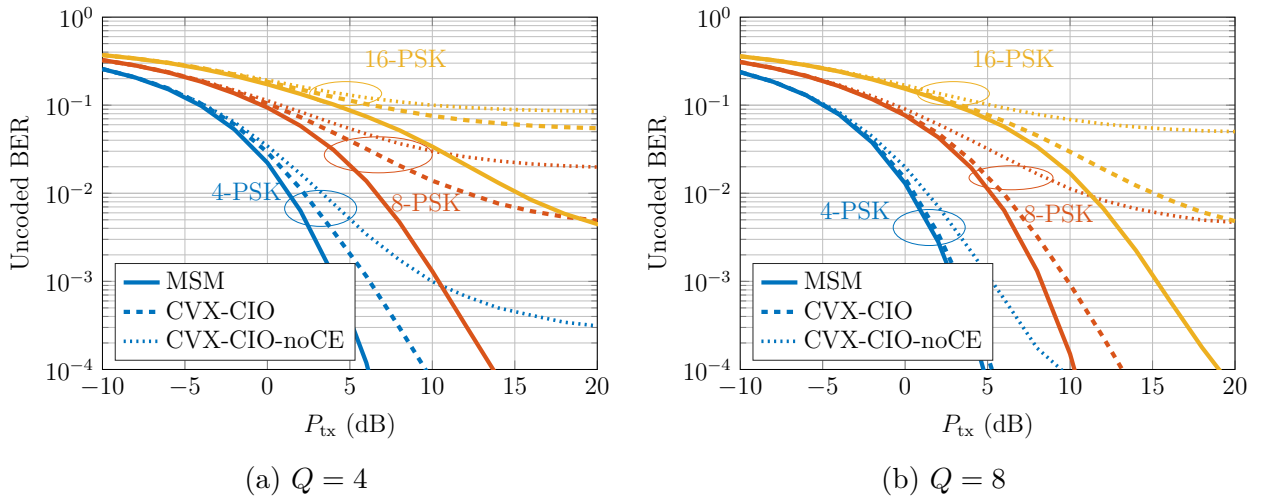


Fig. 7.9: Comparison of the uncoded BER performance between MSM, CVX-CIO from [22] and CVX-CIO-noCE for a MU-MIMO system with $N = 64$ and $M = 8$, © 2018 IEEE.

8. Frequency-Selective Channels

8.1 Input-Output Relationship

In the case of frequency-selective channels, the relaxed noiseless received signal vector is given by

$$\mathbf{y}'[t] = \mathbf{H}[t] * \mathbf{x}[t], \quad \mathbf{x}[t] \in \mathbb{X}^N. \quad (8.1)$$

To mitigate the resulting ISI, multiple time instants have to be jointly considered in the optimization problem. Therefore, we consider in general signal blocks of length B . To this end, let us define the input signal block as

$$\mathbf{s}_B[t] = [\mathbf{s}[t]^T \quad \mathbf{s}[t-1]^T \quad \cdots \quad \mathbf{s}[t-B+1]^T]^T \quad (8.2)$$

and we aim to design the optimal transmit signal block

$$\mathbf{x}_B[t] = [\mathbf{x}[t]^T \quad \mathbf{x}[t-1]^T \quad \cdots \quad \mathbf{x}[t-B+1]^T]^T. \quad (8.3)$$

The relaxed noiseless received signal block $\mathbf{y}'_B[t] = [\mathbf{y}'[t]^T \quad \mathbf{y}'[t-1]^T \quad \cdots \quad \mathbf{y}'[t-B+1]^T]^T$ is in general expressed as

$$\mathbf{y}'_B[t] = \mathbf{H}_B \mathbf{x}_B[t] + \mathbf{y}'_{\text{IBI}}[t], \quad (8.4)$$

where \mathbf{H}_B and $\mathbf{y}'_{\text{IBI}}[t]$ denote the corresponding convolution matrix and the vector containing the Inter-Block Interference (IBI) due to the channel frequency-selectivity. The exact expressions depend on the block length B and are given for different block lengths in Section 8.1.1 and Section 8.1.2.

8.1.1 Symbol-wise Processing: MSM-SP

In the case of symbol-wise processing, i.e. $B = 1$, we get

$$\mathbf{y}' = \mathbf{y}'[t] = \mathbf{H}_0 \mathbf{x}[t] + \sum_{\ell=1}^{L-1} \mathbf{H}_\ell \mathbf{x}[t-\ell]. \quad (8.5)$$

This leads to the expressions

$$\mathbf{H}_B = \mathbf{H}_0 \quad \text{and} \quad \mathbf{y}'_{\text{IBI}}[t] = \sum_{\ell=1}^{L-1} \mathbf{H}_\ell \mathbf{x}[t-\ell]. \quad (8.6)$$

8.1.2 Block-wise Processing

8.1.2.1 Without Cyclic Prefix: MSM-BP

In this section, the precoding task is conducted for a block of length B , with $B \geq L$, such that we get only one interfering block. The noiseless received block signal is given by

$$\mathbf{y}'_B[t] = \mathbf{H}_B \mathbf{x}_B[t] + \mathbf{H}_{\text{IBI}} \mathbf{x}_B[t - B], \quad (8.7)$$

where

$$\mathbf{H}_B = \sum_{\ell=0}^{L-1} \Upsilon_1(B, \ell) \otimes \mathbf{H}_\ell, \quad (8.8)$$

$$\mathbf{H}_{\text{IBI}} = \sum_{\ell=0}^{L-1} \Upsilon_2(B, \ell) \otimes \mathbf{H}_\ell \quad (8.9)$$

$$\Upsilon_1(B, \ell) = \begin{bmatrix} \mathbf{0}_{B-\ell, \ell} & \mathbf{I}_{B-\ell} \\ & \mathbf{0}_{\ell, B} \end{bmatrix} \quad (8.10)$$

$$\Upsilon_2(B, \ell) = \begin{bmatrix} \mathbf{0}_{B-\ell, B} \\ \mathbf{I}_\ell & \mathbf{0}_{\ell, B-L+\ell} \end{bmatrix}. \quad (8.11)$$

This leads to the following exact expression:

$$\mathbf{y}'_{\text{IBI}}[t] = \mathbf{H}_{\text{IBI}} \mathbf{x}_B[t - B]. \quad (8.12)$$

8.1.2.2 With Cyclic Prefix: MSM-BP-CP

To remove the IBI, we append to each block a cyclic prefix of length $L - 1$ and choose a block of length $B \geq L$. Note that due to the cyclic prefix the power per useful transmit vector $\mathbf{t}[t]$ reduces to $P_{\text{tx}} \frac{B}{B+L-1}$. The resulting noiseless received vector is given by

$$\mathbf{y}'_B[t] = \mathbf{H}_{\text{circ}} \mathbf{x}_B[t], \quad (8.13)$$

where

$$\mathbf{H}_{\text{circ}} = \sum_{\ell=0}^{L-1} (\Upsilon_1(B, \ell) + \Upsilon_2(B, \ell)) \otimes \mathbf{H}_\ell \quad (8.14)$$

is the block-circular channel matrix. Thus, we get

$$\mathbf{H}_B = \mathbf{H}_{\text{circ}} \quad \text{and} \quad \mathbf{y}'_{\text{IBI}}[t] = \mathbf{0}_{MB}. \quad (8.15)$$

8.2 Optimization Problem

The MSM precoder that was first introduced in Chapter 7, can be reformulated for frequency-selective channels as

$$\max_{\mathbf{x}_B[t]} \delta \quad (8.16)$$

$$\text{s.t. } y'_{Bm}[t] \in \text{SR}_m, \quad m = 1, \dots, MB \quad (8.17)$$

$$\text{and } \mathbf{x}_B[t] \in \mathbb{X}^{NB}. \quad (8.18)$$

Since the mathematical expressions of the SRs depend on the constellation set \mathbb{S} , we introduce the optimization problem of the MSM precoder for frequency-selective channels separately for PSK signaling and QAM signaling.

8.2.1 PSK Signaling

In analogy to the derivations introduced in

$$\begin{aligned}\mathbf{z}_B &= \text{diag}(\mathbf{s}_B[t]^*) \mathbf{y}'_B[t] \\ &= \tilde{\mathbf{H}}_B \mathbf{x}_B[t] + \tilde{\mathbf{y}}'_{\text{IBI}},\end{aligned}\quad (8.19)$$

where

$$\tilde{\mathbf{H}}_B = \text{diag}(\mathbf{s}_B[t]^*) \mathbf{H}_B \quad (8.20)$$

$$\tilde{\mathbf{y}}'_{\text{IBI}} = \text{diag}(\mathbf{s}_B[t]^*) \mathbf{y}'_{\text{IBI}}[t]. \quad (8.21)$$

When using the following real-valued representation:

$$\Re\{\tilde{\mathbf{H}}_B \mathbf{x}_B[t]\} = \underbrace{\begin{bmatrix} \Re\{\tilde{\mathbf{H}}_B\} & -\Im\{\tilde{\mathbf{H}}_B\} \end{bmatrix}}_{=\mathbf{A}_B} \underbrace{\begin{bmatrix} \Re\{\mathbf{x}_B[t]\} \\ \Im\{\mathbf{x}_B[t]\} \end{bmatrix}}_{=\bar{\mathbf{x}}_B} = \mathbf{A}_B \bar{\mathbf{x}}_B \quad (8.22)$$

$$\Im\{\tilde{\mathbf{H}}_B \mathbf{x}_B[t]\} = \underbrace{\begin{bmatrix} \Im\{\tilde{\mathbf{H}}_B\} & \Re\{\tilde{\mathbf{H}}_B\} \end{bmatrix}}_{=\mathbf{B}_B} \underbrace{\begin{bmatrix} \Re\{\mathbf{x}_B[t]\} \\ \Im\{\mathbf{x}_B[t]\} \end{bmatrix}}_{=\bar{\mathbf{x}}_B} = \mathbf{B}_B \bar{\mathbf{x}}_B, \quad (8.23)$$

and the following definitions:

$$\tilde{\mathbf{a}} = \Re\{\tilde{\mathbf{y}}'_{\text{IBI}}\} - \Im\{\tilde{\mathbf{y}}'_{\text{IBI}}\} \quad (8.24)$$

$$\tilde{\mathbf{b}} = \Re\{\tilde{\mathbf{y}}'_{\text{IBI}}\} + \Im\{\tilde{\mathbf{y}}'_{\text{IBI}}\}, \quad (8.25)$$

the constraint in (7.17) applied on \mathbf{z}_B from (8.19) can be rewritten as

$$\begin{bmatrix} \mathbf{B}_B - \tan \theta \mathbf{A}_B & \frac{1}{\cos \theta} \mathbf{1}_{MB} \\ -\mathbf{B}_B - \tan \theta \mathbf{A}_B & \frac{1}{\cos \theta} \mathbf{1}_{MB} \end{bmatrix} \begin{bmatrix} \bar{\mathbf{x}}_B \\ \delta \end{bmatrix} \leq \begin{bmatrix} \tilde{\mathbf{a}} \\ \tilde{\mathbf{b}} \end{bmatrix}. \quad (8.26)$$

Combining (8.16), (8.26) and (8.18), the optimization problem MSM for PSK signaling can be expressed for the case $P_{\text{tx}} = N$ by the real-valued linear programming problem

$$\begin{aligned}\max_{\mathbf{v}_B} & \begin{bmatrix} \mathbf{0}_{2NB}^T & 1 \end{bmatrix} \mathbf{v}_B \text{ s.t. } \begin{bmatrix} \mathbf{B}_B - \tan \theta \mathbf{A}_B & \frac{1}{\cos \theta} \mathbf{1}_{MB} \\ -\mathbf{B}_B - \tan \theta \mathbf{A}_B & \frac{1}{\cos \theta} \mathbf{1}_{MB} \\ \mathbf{E} \otimes \mathbf{I}_B & \mathbf{0}_{NB(Q-4)} \end{bmatrix} \mathbf{v}_B \leq \begin{bmatrix} \tilde{\mathbf{a}} \\ \tilde{\mathbf{b}} \\ \cos\left(\frac{\pi}{Q}\right) \mathbf{1}_{NB(Q-4)} \end{bmatrix} \\ \text{and} & \begin{bmatrix} -\cos\left(\frac{\pi}{Q}\right) \mathbf{1}_{2NB} \\ 0 \end{bmatrix} \leq \mathbf{v}_B \leq \begin{bmatrix} \cos\left(\frac{\pi}{Q}\right) \mathbf{1}_{2NB} \\ \infty \end{bmatrix},\end{aligned}\quad (8.27)$$

where $\mathbf{v}_B = [\bar{\mathbf{x}}_B^T \ \delta]^T$.

8.2.2 QAM Signaling

In analogy to the derivations introduced in Section 7.4, the noiseless received signal in the modified coordinate system, shifted by

$$\mathbf{o}_B = \alpha (\mathbf{s}_B[t] - \text{sgn}(\mathbf{s}_B[t])) \quad (8.28)$$

and rotated by the phases of the signals of interest in the shifted coordinate system, i.e.

$$\mathbf{s}_{B(\mathbf{o}_B)}[t] = \alpha \mathbf{s}_B[t] - \mathbf{o}_B, \quad (8.29)$$

is given by

$$\begin{aligned} \mathbf{z}_B &= \frac{1}{\sqrt{2}} \text{diag}(\text{sgn}(\mathbf{s}_B^*[t])) (\mathbf{y}'_B[t] - \alpha (\mathbf{s}_B[t] - \text{sgn}(\mathbf{s}_B[t]))) \\ &= \hat{\mathbf{H}}_B \mathbf{x}_B[t] + \hat{\mathbf{y}}'_{\text{IBI}} - \alpha \mathbf{c}_B, \end{aligned} \quad (8.30)$$

where

$$\hat{\mathbf{H}}_B = \frac{1}{\sqrt{2}} \text{diag}(\text{sgn}(\mathbf{s}_B^*[t])) \mathbf{H}_B \quad (8.31)$$

$$\hat{\mathbf{y}}'_{\text{IBI}} = \frac{1}{\sqrt{2}} \text{diag}(\text{sgn}(\mathbf{s}_B^*[t])) \mathbf{y}'_{\text{IBI}}[t] \quad (8.32)$$

$$\mathbf{c}_B = \frac{1}{\sqrt{2}} \text{diag}(\text{sgn}(\mathbf{s}_B^*[t])) (\mathbf{s}_B[t] - \text{sgn}(\mathbf{s}_B[t])). \quad (8.33)$$

When using the following real-valued representation:

$$\Re\{\hat{\mathbf{H}}_B \mathbf{x}_B[t]\} = \underbrace{\begin{bmatrix} \Re\{\hat{\mathbf{H}}_B\} & -\Im\{\hat{\mathbf{H}}_B\} \end{bmatrix}}_{=\mathbf{V}_B} \begin{bmatrix} \Re\{\mathbf{x}_B[t]\} \\ \Im\{\mathbf{x}_B[t]\} \end{bmatrix} = \mathbf{V}_B \bar{\mathbf{x}}_B \quad (8.34)$$

$$\Im\{\hat{\mathbf{H}}_B \mathbf{x}_B[t]\} = \underbrace{\begin{bmatrix} \Im\{\hat{\mathbf{H}}_B\} & \Re\{\hat{\mathbf{H}}_B\} \end{bmatrix}}_{=\mathbf{W}_B} \begin{bmatrix} \Re\{\mathbf{x}_B[t]\} \\ \Im\{\mathbf{x}_B[t]\} \end{bmatrix} = \mathbf{W}_B \bar{\mathbf{x}}_B, \quad (8.35)$$

and the following definitions:

$$\hat{\mathbf{a}} = \Re\{\hat{\mathbf{y}}'_{\text{IBI}}\} - \Im\{\hat{\mathbf{y}}'_{\text{IBI}}\} \quad (8.36)$$

$$\hat{\mathbf{b}} = \Re\{\hat{\mathbf{y}}'_{\text{IBI}}\} + \Im\{\hat{\mathbf{y}}'_{\text{IBI}}\}, \quad (8.37)$$

the constraint in (7.9) applied on \mathbf{z}_B from (8.30) can be rewritten as

$$\begin{bmatrix} \mathbf{W}_B - \mathbf{V}_B & \mathbf{1}_{MB} & \Re\{\mathbf{c}_B\} - \Im\{\mathbf{c}_B\} \\ -\mathbf{W}_B - \mathbf{V}_B & \mathbf{1}_{MB} & \Re\{\mathbf{c}_B\} + \Im\{\mathbf{c}_B\} \\ \mathbf{W}_B + \mathbf{V}_B & \mathbf{1}_{MB} & -\Re\{\mathbf{c}_B\} - \Im\{\mathbf{c}_B\} - \sqrt{2}\boldsymbol{\xi}_2 \\ -\mathbf{W}_B + \mathbf{V}_B & \mathbf{1}_{MB} & -\Re\{\mathbf{c}_B\} + \Im\{\mathbf{c}_B\} - \sqrt{2}\boldsymbol{\xi}_1 \end{bmatrix} \begin{bmatrix} \bar{\mathbf{x}}_B \\ \sqrt{2}\delta \\ \alpha \end{bmatrix} \leq \begin{bmatrix} \hat{\mathbf{a}} \\ \hat{\mathbf{b}} \\ -\hat{\mathbf{a}} \end{bmatrix}. \quad (8.38)$$

By adding the first line to the fourth and the second line to the third, the latter constraint recalculates to

$$\begin{bmatrix} \mathbf{W}_B - \mathbf{V}_B & \mathbf{1}_{MB} & \Re\{\mathbf{c}_B\} - \Im\{\mathbf{c}_B\} \\ -\mathbf{W}_B - \mathbf{V}_B & \mathbf{1}_{MB} & \Re\{\mathbf{c}_B\} + \Im\{\mathbf{c}_B\} \\ \mathbf{0}_{MB,2NB} & \sqrt{2} \mathbf{1}_{MB} & -\boldsymbol{\xi}_2 \\ \mathbf{0}_{MB,2NB} & \sqrt{2} \mathbf{1}_{MB} & -\boldsymbol{\xi}_1 \end{bmatrix} \begin{bmatrix} \bar{\mathbf{x}}_B \\ \sqrt{2}\delta \\ \alpha \end{bmatrix} \leq \begin{bmatrix} \hat{\mathbf{a}} \\ \hat{\mathbf{b}} \\ \mathbf{0}_{MB} \\ \mathbf{0}_{MB} \end{bmatrix}. \quad (8.39)$$

Combining (8.16), (8.39) and (8.18), the optimization problem of the MSM precoder for QAM signaling can be then for $P_{\text{tx}} = N$ expressed by

$$\begin{aligned}
 & \max_{\mathbf{v}_B} [\mathbf{0}_{2NB}^T \quad 1 \quad 0] \mathbf{v}_B \\
 & \text{s.t.} \quad \begin{bmatrix} \mathbf{W}_B - \mathbf{V}_B & \mathbf{1}_{MB} & \Re\{\mathbf{c}_B\} - \Im\{\mathbf{c}_B\} \\ -\mathbf{W}_B - \mathbf{V}_B & \mathbf{1}_{MB} & \Re\{\mathbf{c}_B\} + \Im\{\mathbf{c}_B\} \\ \mathbf{0}_{MB,2NB} & \sqrt{2} \mathbf{1}_{MB} & -\xi_2 \\ \mathbf{0}_{MB,2NB} & \sqrt{2} \mathbf{1}_{MB} & -\xi_1 \\ \mathbf{E} \otimes \mathbf{I}_B & \mathbf{0}_{NB(Q-4)} & \mathbf{0}_{NB(Q-4)} \end{bmatrix} \mathbf{v}_B \leq \begin{bmatrix} \hat{\mathbf{a}} \\ \hat{\mathbf{b}} \\ \mathbf{0}_{MB} \\ \mathbf{0}_{MB} \\ \cos\left(\frac{\pi}{Q}\right) \mathbf{1}_{NB(Q-4)} \end{bmatrix} \\
 & \text{and} \quad \begin{bmatrix} -\cos\left(\frac{\pi}{Q}\right) \mathbf{1}_{2NB} \\ 0 \\ 0 \end{bmatrix} \leq \mathbf{v}_B \leq \begin{bmatrix} \cos\left(\frac{\pi}{Q}\right) \mathbf{1}_{2NB} \\ \infty \\ \infty \end{bmatrix}, \tag{8.40}
 \end{aligned}$$

where $\mathbf{v}_B^T = [\bar{\mathbf{x}}_B^T \quad \sqrt{2}\delta \quad \alpha]$.

8.3 Computational Complexity

We recall the study of the computational complexity of linear programs in Section 7.5.1. Note that the sparsity of the linear program in (8.27) and (8.40) can be exploited by deploying the revised simplex method to reduce the computational complexity [72, p.89].

8.3.1 PSK

From (8.27), we have $m = 2MB + NB(Q - 4)$ inequalities and $n = 2NB + 1$ variables. The number a of artificial variables is defined by the number of non-negative entries in $\tilde{\mathbf{a}}$ and $\tilde{\mathbf{b}}$. We introduce the ratio ε to denote the percentage of non-negative entries in $\tilde{\mathbf{a}}$ and $\tilde{\mathbf{b}}$. It holds that $a = 2MB\varepsilon$. Note that $\varepsilon = 0$ for MSM-BP-CP due to (8.15).

According to Section 7.5.1, the number of floating-point operations for each iteration, under the assumption of massive MIMO systems, is on the order of

$$\#\text{FLOPs}_{\text{PSK}} = \mathcal{O}(4MNB^2 + (Q - 4)(2N^2B^2 + 2NMB^2\varepsilon)). \tag{8.41}$$

For the special case of one-bit quantization, i.e. $Q = 4$, the complexity is linear in N and M but quadratic in B . For $Q > 4$, the complexity becomes quadratic in N and B .

8.3.2 QAM

From (8.40), we have $m = 4MB + NB(Q - 4)$ inequalities and $n = 2NB + 2$ variables. The number a of artificial variables is defined by the number of non-negative entries in $\hat{\mathbf{a}}$ and $\hat{\mathbf{b}}$. We introduce the ratio ε to denote the percentage of non-negative entries in $\hat{\mathbf{a}}$ and $\hat{\mathbf{b}}$. It holds that $a = 2MB\varepsilon$. Note that $\varepsilon = 0$ for MSM-BP-CP due to (8.15).

According to Section 7.5.1, the number of floating-point operations for each iteration, under the assumption of massive MIMO systems, is on the order of

$$\#\text{FLOPs}_{\text{QAM}} = \mathcal{O}(8MNB^2 + (Q - 4)(2N^2B^2 + 2NMB^2\varepsilon)). \tag{8.42}$$

For the special case of one-bit quantization, i.e. $Q = 4$, the complexity is linear in N and M but quadratic in B . For $Q > 4$, the complexity becomes quadratic in N and B .

8.4 Simulation Results

In this section, we compare the performance of the different linear precoding techniques that were introduced above with the ideal WF precoder while assuming full CSI. We denote the precoders by

For the simulations, we consider the i.i.d. channel model with the exponential power delay profiles of $L = 3$ and $L = 6$ given in Table 4.1 and Table 4.2. Assuming full CSI, we compare the performance of the different precoding techniques that were introduced in this chapter with the ideal FIR-WF precoder of length L introduced in [73]. We denote the precoders by

- FIR-WF: the ideal FIR-WF, [70], where no CEQ is applied in the system.
- FIR-WF + CEQ: the ideal FIR-WF, that is followed by the optimal CEQ.
- MSM-SP: the MSM precoder based on symbol-wise processing as introduced in Section 8.1.1.
- MSM-BP: the MSM precoder based on block-wise processing as introduced in Section 8.1.2.1.
- MSM-BP-CP: the MSM precoder based on block-wise processing with the cyclic prefix as introduced in Section 8.1.2.2.

We assume a BS with $N = 64$ antennas serving $M = 8$ single-antenna users with 16-QAM signals. The numerical results are obtained with Monte Carlo simulations of 100 independent frequency-selective channel realizations. The AWGN is also i.i.d. with variance one at each antenna. The performance metric is the uncoded BER averaged over the single-antenna users. For the blind estimation of the coefficients g_m we use a block length of $T = 128$. The numerical results are plotted in Fig. 8.1a and Fig. 8.1b for $Q = 4$. It can be deduced, that the smaller the block length B , the worse the performance is. The MSM-SP performs the worst. When performing block-wise processing with $B = 4$ and $B = 8$, the gain in terms of the uncoded BER is moderate, whereas the addition of the cyclic prefix in the block-wise processing leads to a significant gain. This gain is due to the fact that the cyclic prefix removes the IBI, and comes with decreased computational complexity for the same block length and $Q > 4$. However, it comes at the cost of a decrease in the throughput, since $B/(B+L-1)$ is the percentage of the useful transmission. In every block transmission there are $L-1$ redundant (unuseful) symbols due to the cyclic prefix. Increasing the block length decreases the IBI and increases the throughput. However, it increases the computational complexity as discussed in Section 8.3. Thus, there is a trade-off between the performance, the throughput and the computational complexity.

It can be also noticed that the QCE system can achieve the same performance as the ideal linear system with the method MSM-BP-CP with an optimal value of B and by doubling the number of antennas.

In Fig. 8.2, the effect of increasing the resolution of the polar DACs is studied. We choose the precoder MSM-BP-CP with a block length of $B = 16$ and plot the uncoded BER as a function of the available power P_{tx} for different values of Q . As expected the performance improves for higher resolution. The performance difference between $Q = 8$ and $Q = 16$ is very

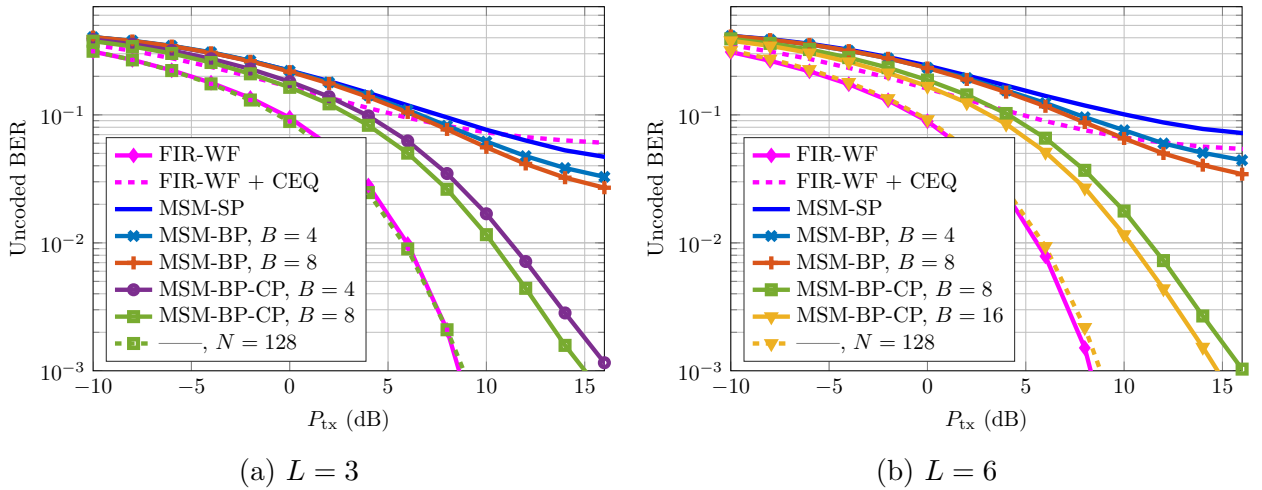


Fig. 8.1: BER performance for a MU-MIMO system with $N = 64$, $M = 8$ and $Q = 4$ with 16-QAM and the i.i.d. channel model.

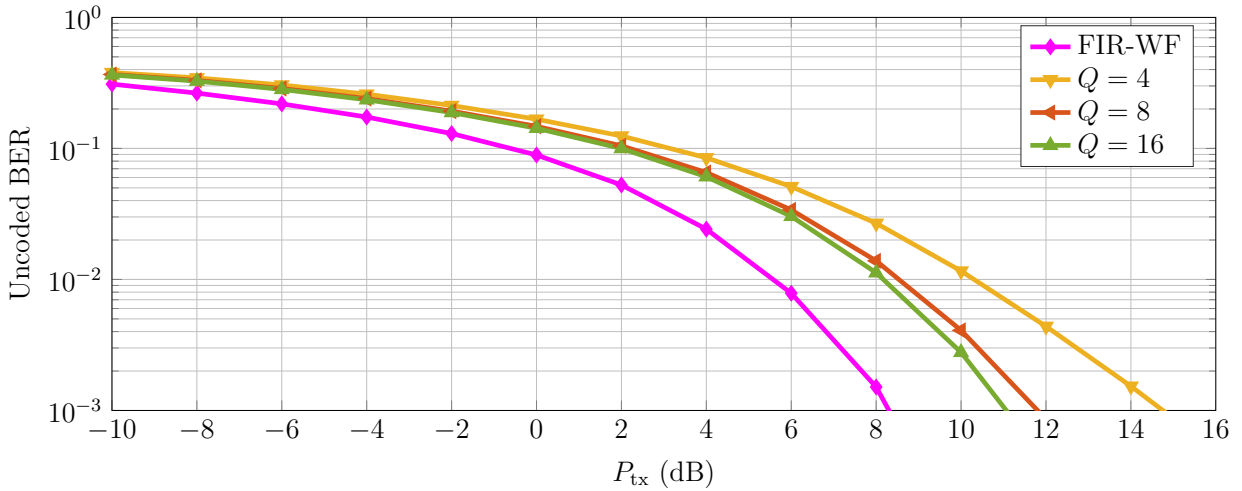


Fig. 8.2: BER performance for a MU-MIMO system with $N = 64$, $M = 8$ with 16-QAM signaling and channel exponential power delay profile with $L = 6$ for different quantization resolutions: MSM-BP-CP with $B = 16$

moderate compared to the required additional computational complexity and the resulting higher power consumption.

In summary, the proposed MSM method offers significant gains in terms of the unencoded BER for the case of flat-fading channels as well as for frequency-selective channels. It is however characterized by its high computational complexity especially for frequency-selective channels, where block-wise processing is required. This challenge can be overcome by the steady enhancement of the hardware computational power.

9. Discussion: Benefits and Challenges

Throughout this work, we are justifying the use of massive MIMO systems with CEQs by their high power efficiency. Therefore, we show in Section 9.1 a comparison between an ideal linear system and a quantized system in terms of power consumption, which proves the benefit of coarsely QCE systems. However, a strong argument against coarse quantization at the transmitter is the spectral regrowth. Therefore, we look deeper in Section 9.2 at the obtained spectrum in the presence of coarse quantization in massive MIMO systems.

9.1 Power Efficiency

As shown in the previous chapters, the same performance of an ideal linear system, i.e. linear PAs and infinite resolution DACs, can be obtained with a QCE system having at most the double number of antennas as in the ideal case. We compare the power consumed at the PAs in both cases.

A linear PA has a given power efficiency denoted by η_{LPA} . However, only a fraction of the ideal power efficiency η_{LPA} can be achieved, if the input signal at the PA shows large PAPR, which is the case for systems without CE constraint at the PA input. In average each antenna has a transmit power of P_{tx}/N , thus the dissipated power at all PAs in the ideal system is given by

$$P_{\text{DC}_i} = N \frac{P_{\text{tx}}}{N} \frac{1}{\eta_{\text{LPA}}} = \frac{P_{\text{tx}}}{\eta_{\text{LPA}}}, \quad (9.1)$$

whereas the dissipated power in a QCE system, having the same BER performance for the same transmit power P_{tx} , is given by

$$P_{\text{DC}_Q} = 2N \frac{P_{\text{tx}}}{2N} \frac{1}{\eta_{\text{NLPA}}} = \frac{P_{\text{tx}}}{\eta_{\text{NLPA}}}, \quad (9.2)$$

where η_{NLPA} denotes the power efficiency of the PA deployed in the QCE system. Since the input to the PA in this case is of CE, there is no need to use linear PAs. Switching PAs that show high power efficiency are then suitable. A recent PA that operates at the saturation region and has a power efficiency of more than 90% is the PA class M [51]. Since QCE signals have a PAPR value equal to 1, the PA in the QCE case can be operated at its highest power efficiency.

By assuming that the linear PA achieves at most 20% of power efficiency, the ratio between

the dissipated powers in the ideal case and the QCE case is given by

$$\begin{aligned} \frac{P_{\text{DC}_i}}{P_{\text{DC}_Q}} &= \frac{\eta_{\text{NLPA}}}{\eta_{\text{LPA}}} \\ &= \frac{0.9}{0.2} = 4.5. \end{aligned} \quad (9.3)$$

This means that the power dissipated at the PAs in the ideal case is more than 4 times the power dissipated in the QCE case. In this comparison, the power consumption related to other hardware components like up-converters, is ignored. However, we justify this analysis by the fact that the PA is the most power hungry device in the BS. Note that the DAC in the quantized case is of very low resolution, which first leads to reduced power consumption of the DAC itself and second simplifies the surrounding circuitry.

9.2 Strong Non-linearities vs. Spectral Shaping

The coarse quantization is usually related to the potential deterioration of the pulse shaping in the time domain and thus the increase in bandwidth and the introduction of undesired out of band radiations in the frequency domain. This is true in the case of one transmit antenna and Single-User (SU) scenario. However, does this statement hold true for MU massive MIMO systems? In the sequel, we restrict our analysis to the case of $Q = 4$ for simplicity. The analysis for the case of $Q = 8$ is left for future work.

To check whether the BS transmits signals that respect a specified spectral mask, a measuring device is put in the far field of it and measures the Power Spectral Density (PSD) of the received signal. Therefore, we assume the measuring device as being the virtual $(M+1)$ -th user and its continuous-time noiseless received signal is denoted by $y_{M+1}(t')$, which is expressed as

$$y_{M+1}(t') = \mathbf{h}_{M+1}^T \mathbf{t}'(t'), \quad (9.4)$$

where \mathbf{h}_{M+1} denotes the channel vector between the BS antennas and the measuring device. \mathbf{h}_{M+1} is a random vector that is drawn from the same distribution as the user channels. Moreover, $\mathbf{t}'(t')$ denotes the continuous-time transmit signal. The vector $\mathbf{t}'(t')$ is obtained by converting the discrete-time signal $\mathbf{t}[t]$, i.e. the CEQ output signal, to the analog world. In the following, we will differentiate between two cases

- standard QAM staggering
- offset QAM staggering.

9.2.1 Standard QAM

In this case, the continuous-time signal $\mathbf{t}'(t')$ is obtained by

$$\mathbf{t}'(t') = \sum_{t=0}^{\infty} \mathbf{t}[t] g_{\text{DAC}}(t' - tT_s), \quad (9.5)$$

where T_s represents the sample period and $g_{\text{DAC}}(t')$ is a real-valued impulse response. Since we are interested in the PSD at the measuring device, we first compute the autocorrelation

function of the signal $y_{M+1}(t')$.

$$\begin{aligned}\rho_{y_{M+1}y_{M+1}}(\tau') &= \frac{\text{E} [y_{M+1}(t')y_{M+1}^*(t' - \tau')]}{\sigma_{y_{M+1}}^2} \\ &= \frac{\mathbf{h}_{M+1}^T \text{E} [\mathbf{t}'(t')\mathbf{t}'^H(t' - \tau')] \mathbf{h}_{M+1}^*}{\mathbf{h}_{M+1}^T \text{E} [\mathbf{t}'(t')\mathbf{t}'^H(t')] \mathbf{h}_{M+1}^*}.\end{aligned}\quad (9.6)$$

The covariance matrix $\text{E} [\mathbf{t}'(t')\mathbf{t}'^H(t' - \tau')]$ is given by

$$\begin{aligned}\text{E} [\mathbf{t}'(t')\mathbf{t}'^H(t' - \tau')] &= \sum_{k=0}^{\infty} \sum_{\ell'=-\infty}^k g_{\text{DAC}}(t' - kT_s)g_{\text{DAC}}(t' - kT_s - (\tau' - \ell'T_s)) \text{E} [\mathbf{t}[k]\mathbf{t}^H[k - \ell']] \\ &= \sum_{\ell'=-\infty}^{\infty} \mathbf{C}_{\mathbf{t}\mathbf{t}}[\ell'] (g_{\text{DAC}} * g_{\text{DAC}})(\tau' - \ell'T_s) \\ &= (\mathbf{C}_{\mathbf{t}\mathbf{t}} * (g_{\text{DAC}} * g_{\text{DAC}}))(\tau').\end{aligned}\quad (9.7)$$

Therefore, we get

$$\rho_{y_{M+1}y_{M+1}}(\tau') = \frac{\mathbf{h}_{M+1}^T ((g_{\text{DAC}} * g_{\text{DAC}}) * \mathbf{C}_{\mathbf{t}\mathbf{t}})(\tau') \mathbf{h}_{M+1}^*}{\mathbf{h}_{M+1}^T ((g_{\text{DAC}} * g_{\text{DAC}}) * \mathbf{C}_{\mathbf{t}\mathbf{t}})(0) \mathbf{h}_{M+1}^*}.\quad (9.8)$$

9.2.1.1 Without Digital Pulse Shaping

If no digital pulse shaping is applied to the precoder output, then there is no time correlation for the signal \mathbf{t} ; that is

$$\mathbf{C}_{\mathbf{t}\mathbf{t}}[\tau] = \begin{cases} \mathbf{C}_{\mathbf{t}\mathbf{t}}, & \text{if } \tau = 0, \\ \mathbf{0}_{N,N}, & \text{otherwise.} \end{cases}\quad (9.9)$$

Therefore, the autocorrelation function simplifies to

$$\begin{aligned}\rho_{y_{M+1}y_{M+1}}(\tau') &= \frac{(g_{\text{DAC}} * g_{\text{DAC}})(\tau') \mathbf{h}_{M+1}^T \mathbf{C}_{\mathbf{t}\mathbf{t}} \mathbf{h}_{M+1}^*}{(g_{\text{DAC}} * g_{\text{DAC}})(0) \mathbf{h}_{M+1}^T \mathbf{C}_{\mathbf{t}\mathbf{t}} \mathbf{h}_{M+1}^*} \\ &= \frac{(g_{\text{DAC}} * g_{\text{DAC}})(\tau')}{(g_{\text{DAC}} * g_{\text{DAC}})(0)}.\end{aligned}\quad (9.10)$$

The PSD at the measuring device is then given by

$$\begin{aligned}S_{y_{M+1}y_{M+1}}(f) &= \int_{-\infty}^{\infty} \rho_{y_{M+1}y_{M+1}}(\tau') e^{-j2\pi f\tau'} d\tau' \\ &= \frac{|G_{\text{DAC}}(f)|^2}{\int_{-\infty}^{\infty} |G_{\text{DAC}}(f)|^2 df}.\end{aligned}\quad (9.11)$$

Thus, the shape of the PSD in this case is determined by the transfer function of the DAC and is independent of the number of antennas N and users M . By choosing the time impulse

response of the DAC as the rectangular function,

$$g_{\text{DAC}}(t') = \text{rect}\left(\frac{t'}{T_s}\right) = \begin{cases} 1 & , \text{if } |t'| < \frac{T_s}{2} \\ \frac{1}{2} & , \text{if } |t'| = \frac{T_s}{2} \\ 0 & , \text{otherwise,} \end{cases} \quad (9.12)$$

and thus the frequency transfer function

$$G_{\text{DAC}}(f) = \frac{\sin(\pi f/T_s)}{\pi f/T_s}, \quad (9.13)$$

we get the PSD depicted in Fig. 9.1 at the measuring device. Note that the sample period T_s is equal to the symbol period T_{sym} in the case of no digital pulse shaping.

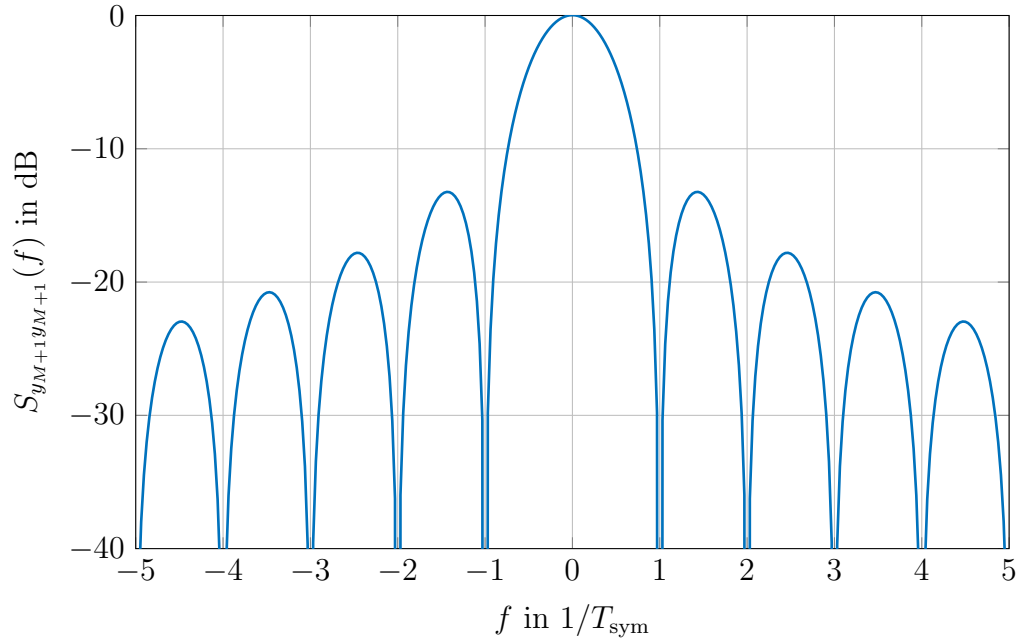


Fig. 9.1: PSD at the measuring device with the DAC as the rectangular function in standard QAM staggering and without digital pulse shaping.

9.2.1.2 With Digital Pulse Shaping

We assume that the input of the CEQ $\mathbf{x}[t]$ is given by

$$\begin{aligned} \mathbf{x}[t] &= g_{\text{PS}}[t] * \mathbf{x}'[t] \\ &= \sum_{\ell=0}^{\infty} \mathbf{x}'[\ell] g_{\text{PS}}[t - \ell], \end{aligned} \quad (9.14)$$

where

$$g_{\text{PS}}[t] = \sum_{t=-\infty}^{\infty} g_{\text{PS}}(t') \delta(t' - tT_s), T_s = \frac{T_{\text{sym}}}{2}, \quad (9.15)$$

and $\mathbf{x}'[t]$ is the output of the precoder. The Discrete Time Fourier Transform (DTFT) of $g_{\text{PS}}[t]$, leads to the transfer function $G_{\text{PS}}(e^{j2\pi f T_s})$, which is continuous and periodic. The covariance matrix at the input of the CEQ calculates to

$$\begin{aligned}\mathbf{C}_{\mathbf{xx}}[t] &= \sum_{\ell'=-\infty}^{\infty} \mathbf{C}_{\mathbf{x}'\mathbf{x}'}[\ell'] (g_{\text{PS}} * g_{\text{PS}})[t - \ell'] \\ &= \mathbf{C}_{\mathbf{x}'\mathbf{x}'} (g_{\text{PS}} * g_{\text{PS}})[t].\end{aligned}\quad (9.16)$$

Hence, the covariance matrix $\mathbf{R}_{\mathbf{xx}}[\tau]$ calculates to

$$\begin{aligned}\mathbf{R}_{\mathbf{xx}}[\tau] &= \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \mathbb{E} [\mathbf{x}[t]\mathbf{x}^{\text{H}}[t - \tau]] \text{diag}(\mathbf{C}_{\mathbf{xx}})^{-1/2} \\ &= \frac{(g_{\text{PS}} * g_{\text{PS}})[\tau]}{(g_{\text{PS}} * g_{\text{PS}})[0]} (\text{diag}(\mathbf{C}_{\mathbf{x}'\mathbf{x}'})^{-1/2} \mathbb{E} [\mathbf{x}'\mathbf{x}'^{\text{H}}] \text{diag}(\mathbf{C}_{\mathbf{x}'\mathbf{x}'})^{-1/2}) \\ &= \frac{(g_{\text{PS}} * g_{\text{PS}})[\tau]}{(g_{\text{PS}} * g_{\text{PS}})[0]} \mathbf{R}_{\mathbf{x}'\mathbf{x}'}.\end{aligned}\quad (9.17)$$

Consequently, the autocorrelation function of the signal $y_{M+1}(t')$ after applying Price's Theorem in (3.52) reads as

$$\begin{aligned}\rho_{y_{M+1}y_{M+1}}(\tau) &= \left(\mathbf{h}_{M+1}^{\text{T}} \Xi \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \left((g_{\text{DAC}} * g_{\text{DAC}}) * \arcsin \left(\Re \{ \mathbf{R}_{\mathbf{xx}} e^{-j2\Delta k\psi} \} \right) \right) (0) \Xi \mathbf{h}_{M+1}^* \right)^{-1} \\ &\quad \mathbf{h}_{M+1}^{\text{T}} \Xi \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} (g_{\text{DAC}} * g_{\text{DAC}} * \arcsin \left(\Re \{ \mathbf{R}_{\mathbf{xx}}[\tau] e^{-j2\Delta k\psi} \} \right)) (\tau) \Xi \mathbf{h}_{M+1}^* \\ &= \mathbf{h}_{M+1}^{\text{T}} \Xi \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \\ &\quad \left(g_{\text{DAC}} * g_{\text{DAC}} * \arcsin \left(\Re \left\{ \frac{(g_{\text{PS}} * g_{\text{PS}})}{(g_{\text{PS}} * g_{\text{PS}})[0]} \mathbf{R}_{\mathbf{x}'\mathbf{x}'} e^{-j2\Delta k\psi} \right\} \right) \right) (\tau) \Xi \mathbf{h}_{M+1}^* \\ &\quad \left(\mathbf{h}_{M+1}^{\text{T}} \Xi \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \right. \\ &\quad \left. \left(g_{\text{DAC}} * g_{\text{DAC}} * \arcsin \left(\Re \left\{ \frac{(g_{\text{PS}} * g_{\text{PS}})}{(g_{\text{PS}} * g_{\text{PS}})[0]} \mathbf{R}_{\mathbf{x}'\mathbf{x}'} e^{-j2\Delta k\psi} \right\} \right) \right) (0) \Xi \mathbf{h}_{M+1}^* \right)^{-1}.\end{aligned}\quad (9.18)$$

The PSD is then obtained by transforming the autocorrelation function into the frequency domain. Thus, we get

$$\begin{aligned}
S_{y_{M+1}y_{M+1}}(f) &= \int_{-\infty}^{\infty} \rho_{y_{M+1}y_{M+1}}(\tau') e^{-j2\pi f\tau'} d\tau', \\
&= \left(\mathbf{h}_{M+1}^T \Xi \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \right. \\
&\quad \left. \left(g_{\text{DAC}} * g_{\text{DAC}} * \arcsin \left(\Re \left\{ \frac{(g_{\text{PS}} * g_{\text{PS}})}{(g_{\text{PS}} * g_{\text{PS}})[0]} \mathbf{R}_{\mathbf{x}'\mathbf{x}'} e^{-j2\Delta k\psi} \right\} \right) \right) (0) \Xi \mathbf{h}_{M+1}^* \right)^{-1} \\
&\quad |G_{\text{DAC}}(f)|^2 \sum_{\Delta k=0}^{Q/2-1} e^{j(2\Delta k\psi)} \text{tr} \left(\Xi \mathbf{h}_{M+1}^* \mathbf{h}_{M+1}^T \Xi \right. \\
&\quad \left. \sum_{\tau=-\infty}^{\infty} \arcsin \left(\Re \left\{ \frac{(g_{\text{PS}} * g_{\text{PS}})[\tau]}{(g_{\text{PS}} * g_{\text{PS}})[0]} \mathbf{R}_{\mathbf{x}'\mathbf{x}'} e^{-j2\Delta k\psi} \right\} \right) e^{-j2\pi f\tau} \right). \tag{9.19}
\end{aligned}$$

In the case that the channels of the users show small correlation, we can assume that the non-diagonal entries of $\mathbf{R}_{\mathbf{x}'\mathbf{x}'}$ are small enough to use the first order Taylor approximation of the arcsin function. According to Section 3.6, we get

$$\begin{aligned}
S_{y_{M+1}y_{M+1}}(f) &\approx |G_{\text{DAC}}(f)|^2 \\
&\quad \left(\mathbf{h}_{M+1}^T \left(\alpha_q \frac{(g_{\text{DAC}} * g_{\text{DAC}} * g_{\text{PS}} * g_{\text{PS}})(0)}{(g_{\text{PS}} * g_{\text{PS}})[0]} \Xi \mathbf{R}_{\mathbf{x}'\mathbf{x}'} \Xi + \beta_q (g_{\text{DAC}} * g_{\text{DAC}})(0) \Xi^2 \right) \mathbf{h}_{M+1}^* \right)^{-1} \\
&\quad \left(\mathbf{h}_{M+1}^T \left(\alpha_q \frac{|G_{\text{PS}}(e^{j2\pi f T_s})|^2}{(g_{\text{PS}} * g_{\text{PS}})[0]} \Xi \mathbf{R}_{\mathbf{x}'\mathbf{x}'} \Xi + \beta_q \Xi^2 \right) \mathbf{h}_{M+1}^* \right) \\
&= |G_{\text{DAC}}(f)|^2 \left(\mathbf{h}_{M+1}^T \Xi \left(|G_{\text{PS}}(e^{j2\pi f T_s})|^2 \mathbf{R}_{\mathbf{x}'\mathbf{x}'} + \beta_q / \alpha_q (g_{\text{PS}} * g_{\text{PS}})[0] \mathbf{I}_N \right) \Xi \mathbf{h}_{M+1}^* \right) \\
&\quad \left(\mathbf{h}_{M+1}^T \Xi \left((g_{\text{DAC}} * g_{\text{DAC}} * g_{\text{PS}} * g_{\text{PS}})(0) \mathbf{R}_{\mathbf{x}'\mathbf{x}'} \right. \right. \\
&\quad \left. \left. + \beta_q / \alpha_q (g_{\text{PS}} * g_{\text{PS}})[0] (g_{\text{DAC}} * g_{\text{DAC}})(0) \mathbf{I}_N \right) \Xi \mathbf{h}_{M+1}^* \right)^{-1}. \tag{9.20}
\end{aligned}$$

By choosing the impulse response of the DAC as in (9.12) and the pulse shaper as the Root-Raised-Cosine (RRC) filter of roll-off factor $ro = 0.1$, we obtain the PSDs according to (9.20) depicted in Fig. 9.2 for different values of the ratio N/M , when applying linear and non-linear precoding. The PSD is measured at one of the active users. It can be observed that the spectrum is narrower for larger values of the ratio N/M . Additionally, there is a slight improvement in the spectral shape in the case of the linear precoder, which is explained by the non-equal power allocation at the antennas.

Additionally, the PSD is plotted for a random \mathbf{h}_{M+1} , where $\mathbf{h}_{M+1} \in \text{range}(\mathbf{R}_{\mathbf{x}'\mathbf{x}'})$ and for the \mathbf{h}_{M+1} being the eigenvector with the largest eigenvalue of $\mathbf{R}_{\mathbf{x}'\mathbf{x}'}$ in Fig. 9.3. For a random $\mathbf{h}_{M+1} \in \text{range}(\mathbf{R}_{\mathbf{x}'\mathbf{x}'})$, there is no improvement in the spectral shape for large ratios N/M . This is due to the fact that the signals in that direction do not overlap constructively. The precoder designs the transmit signals such that they overlap constructively at the users.

When \mathbf{h}_{M+1} is chosen as the eigenvector with the largest eigenvalue, we can see the best spectral shape that can be obtained, which improves with the ratio N/M .

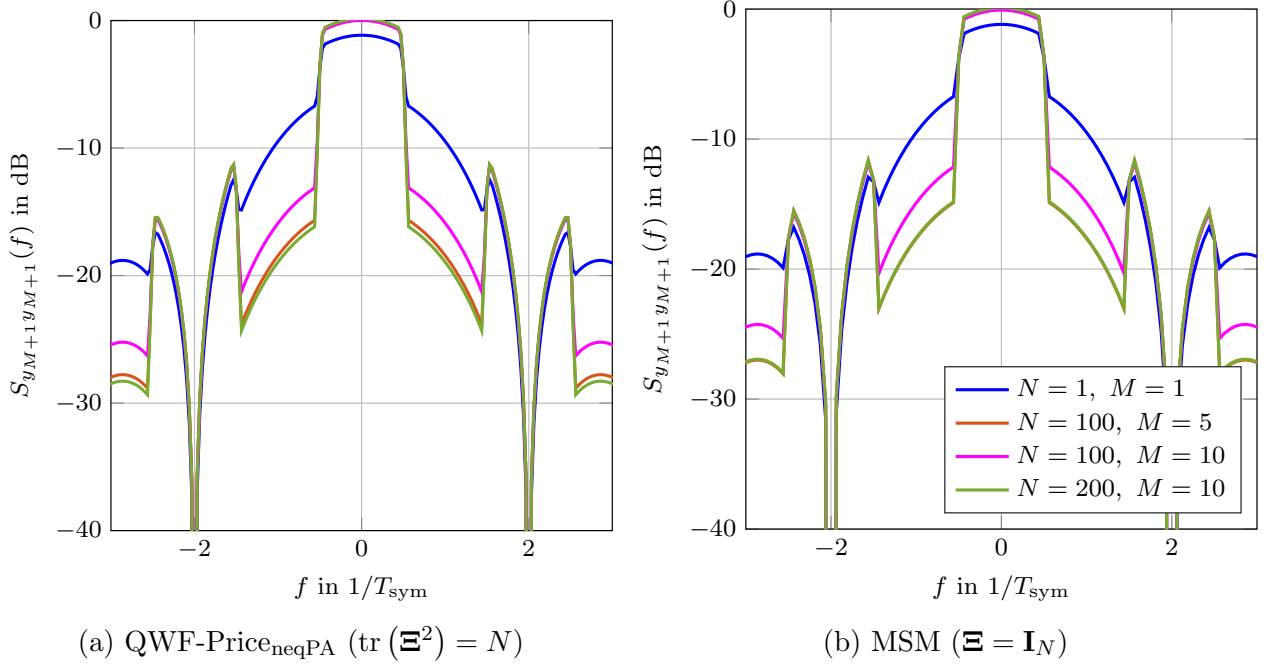


Fig. 9.2: PSD at one of the users for an i.i.d. channel realization with the DAC as the rectangular function in standard QAM staggering and with digital pulse shaping, i.e RRC filter with roll-off factor $ro = 0.1$.

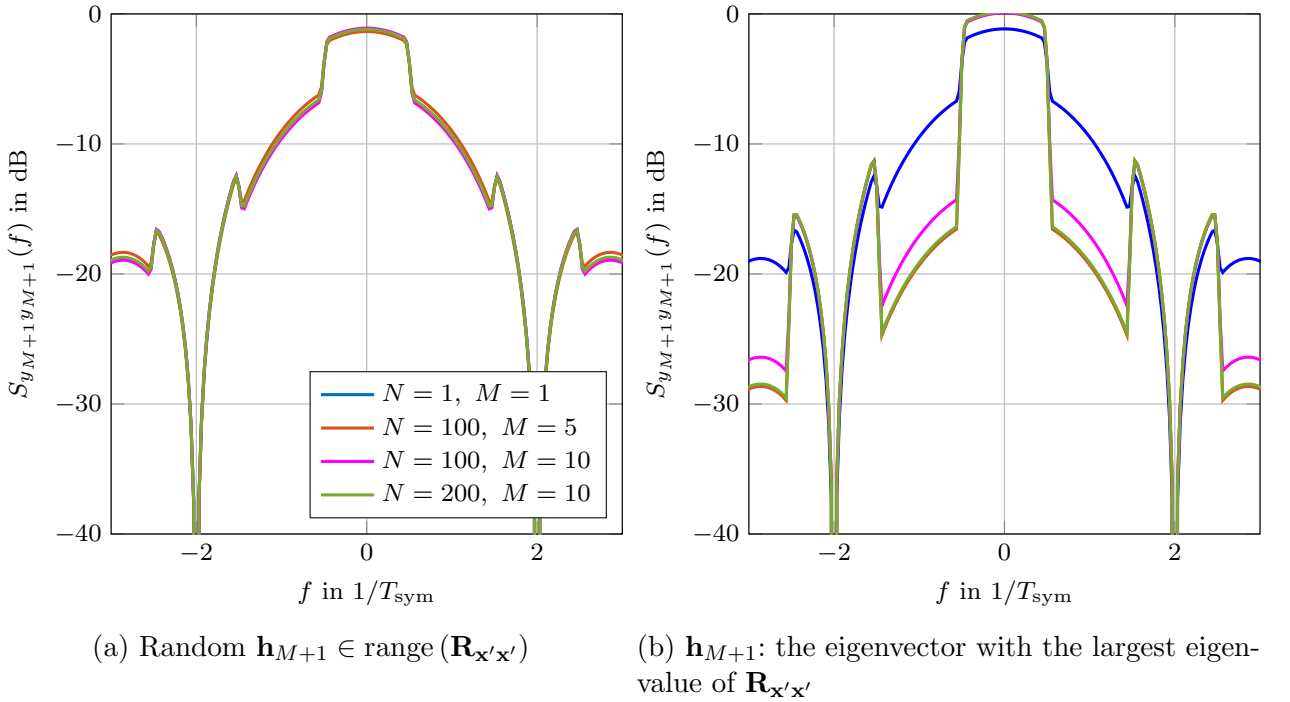


Fig. 9.3: PSD at the measuring device for an i.i.d. channel realization with the DAC as the rectangular function in standard QAM staggering, with digital pulse shaping, i.e RRC filter with roll-off factor $ro = 0.1$ and with the MSM precoder.

Asymptotic Analysis Let us look at the obtained approximate expression of the PSD in (9.20) and assume that the precoder is close to the matched filter of the channel such that

$$\mathbf{C}_{\mathbf{x}'\mathbf{x}'} \approx \mathbf{H}^H \mathbf{H}, \quad (9.21)$$

$$\mathbf{R}_{\mathbf{x}'\mathbf{x}'} \approx \text{diag}(\mathbf{H}^H \mathbf{H})^{-1/2} \mathbf{H}^H \mathbf{H} \text{diag}(\mathbf{H}^H \mathbf{H})^{-1/2}. \quad (9.22)$$

By assuming the i.i.d. channel model described in Section 4.6.1 and a large number of antennas and users with $N \gg M$, which leads to almost orthogonal user channels that span an M -dimensional subspace, we get

$$\mathbb{E} [\text{diag}(\mathbf{H}^H \mathbf{H})] = M \mathbf{I}_N, \quad (9.23)$$

$$\mathbb{E} [\mathbf{h}_{M+1}^T \mathbf{h}_{M+1}^*] = N,$$

$$0 \leq \mathbb{E} [\mathbf{h}_{M+1}^T \mathbf{H}^H \mathbf{H} \mathbf{h}_{M+1}^*] \leq N^2, \quad (9.24)$$

where the lower bound is reached, if \mathbf{h}_{M+1} lies in the null space of \mathbf{H}^* , i.e. $\mathbf{h}_{M+1} \in \text{null}(\mathbf{H}^*)$, and the right limit is reached, if \mathbf{h}_{M+1} lies in the range space of \mathbf{H}^T , i.e. $\mathbf{h}_{M+1} \in \text{range}(\mathbf{H}^T)$. We first assume that the CEQ is optimal and hence $\mathbf{\Xi} = \sqrt{\alpha_q} \text{diag}(\mathbf{C}_{\mathbf{xx}})^{1/2} = \sqrt{\alpha_q (g_{\text{PS}} * g_{\text{PS}})(0)} \text{diag}(\mathbf{C}_{\mathbf{x}'\mathbf{x}'})^{1/2}$. Therefore, if $\mathbf{h}_{M+1} \in \text{range}(\mathbf{H}^T)$, the approximation of the PSD expression is asymptotically given by

$$\begin{aligned} S_{y_{M+1}y_{M+1}}(f) &\rightarrow |G_{\text{DAC}}(f)|^2 (\alpha_q (N^2 c_1 + \beta_q / \alpha_q c_2 N M))^{-1} \\ &\quad \left(\alpha_q \left(|G_{\text{PS}}(e^{j2\pi f T_s})|^2 N^2 + \beta_q / \alpha_q c_3 N M \right) \right) \\ &\rightarrow |G_{\text{DAC}}(f)|^2 \frac{\alpha_q \frac{N}{M} |G_{\text{PS}}(e^{j2\pi f T_s})|^2 + \beta_q c_3}{\alpha_q \frac{N}{M} c_1 + \beta_q c_2}, \end{aligned} \quad (9.25)$$

where

$$c_1 = (g_{\text{PS}} * g_{\text{PS}} * g_{\text{DAC}} * g_{\text{DAC}})(0), \quad (9.26)$$

$$c_2 = (g_{\text{PS}} * g_{\text{PS}})[0] (g_{\text{DAC}} * g_{\text{DAC}})(0), \quad (9.27)$$

$$c_3 = (g_{\text{PS}} * g_{\text{PS}})[0]. \quad (9.28)$$

The obtained expression, under the i.i.d. channel assumption and large number of N and M , consists of two parts: the desired PSD of the pulse shaper that increases linearly with the ratio $\alpha_q N/M$ and a constant noise part due to the quantization errors. Indeed, from this approximation we can see that increasing the ratio N/M improves the spectral shape of the received signal.

In the extreme case that \mathbf{h}_{M+1}^T is orthogonal to all user channels, we get

$$\begin{aligned} S_{y_{M+1}y_{M+1}}(f) &\rightarrow |G_{\text{DAC}}(f)|^2 \frac{c_3}{c_2} \\ &\rightarrow \frac{|G_{\text{DAC}}(f)|^2}{\int_{-\infty}^{\infty} |G_{\text{DAC}}(f)|^2 df}. \end{aligned} \quad (9.29)$$

We see here that the digital pulse shaping vanishes and the spectrum is shaped only by the transfer function of the DAC.

Hence, this asymptotic analysis explains the behavior observed in the previous section.

9.2.2 Offset QAM

In this case, the output of the CEQ $\mathbf{t}[t]$ undergoes Offset QAM staggering. This means that the inphase and quadrature parts of the signal are alternated in a time distance of $T_s/2$. Therefore, the continuous-time signal at the transmitter $\mathbf{t}'(t')$ reads as

$$\mathbf{t}'(t') = \sum_{t=0}^{\infty} \tilde{\mathbf{t}}[2t]g_{\text{DAC}}(t' - 2tT') + \tilde{\mathbf{t}}[2t+1]g_{\text{DAC}}(t' - (2t+1)T'), \quad (9.30)$$

where

$$\tilde{\mathbf{t}}[2t] = \Re \{ \mathbf{t}[t] \}, \quad (9.31)$$

$$\tilde{\mathbf{t}}[2t+1] = \text{j} \Im \{ \mathbf{t}[t] \}. \quad (9.32)$$

We introduce two different impulse responses for the DAC. The first one is the the rectangular function introduced in (9.12)

$$g_{\text{DAC}}(t') = \text{rect} \left(\frac{t'}{2T'} \right), \quad (9.33)$$

where $T' = T_s/2$. The output signal $\mathbf{t}'(t')$ is of constant envelope and the phase follows a rectangular time shape. The maximal difference between the phase values at subsequent time instants is equal to $\pm\pi/4$.

The second time impulse response for the DAC is the cosine impulse

$$g_{\text{DAC}}(t') = g_{\text{MSK}}(t') = \frac{1}{T'} \cos \left(\frac{\pi}{2T'} t' \right), \quad -T' \leq t' \leq T', \quad (9.34)$$

which leads in the considered case, i.e. $Q = 4$, to Minimum-Shift Keying (MSK) signaling at the BS antennas. The corresponding frequency transfer function is given by

$$G_{\text{MSK}}(f) = \frac{4 \cos(2\pi f T')}{\pi 1 - (4f T')^2}. \quad (9.35)$$

Note that the output signal $\mathbf{t}'(t')$ is also of constant envelope and the phase changes linearly in time.

In order to obtain the analytical expression of the autocorrelation function $\rho_{y_{M+1}y_{M+1}}(\tau')$, we have to compute the covariance matrix $\text{E} \left[\mathbf{t}'(t') \mathbf{t}'^{\text{H}}(t' - \tau') \right]$; that is

$$\begin{aligned} \text{E} \left[\mathbf{t}'(t') \mathbf{t}'^{\text{H}}(t' - \tau') \right] &= \sum_{\ell=-\infty}^{\infty} \Re \{ \mathbf{C}_{\text{tt}}[2\ell] \} (g_{\text{DAC}} * g_{\text{DAC}})(\tau' - 2\ell T') \\ &\quad - \text{j} \sum_{\ell'=-\infty}^{\infty} \mathbf{C}_{\Re\{\mathbf{t}\}\Im\{\mathbf{t}\}}[2\ell'] (g_{\text{DAC}} * g_{\text{DAC}})(\tau' - (2\ell' - 1)T') \\ &\quad + \text{j} \sum_{\ell''=-\infty}^{\infty} \mathbf{C}_{\Im\{\mathbf{t}\}\Re\{\mathbf{t}\}}[2\ell''] (g_{\text{DAC}} * g_{\text{DAC}})(\tau' - (2\ell'' + 1)T') \\ &= (\Re \{ \mathbf{C}_{\text{tt}} \} * (g_{\text{DAC}} * g_{\text{DAC}}))(\tau') \\ &\quad - \text{j} (\mathbf{C}_{\Re\{\mathbf{t}\}\Im\{\mathbf{t}\}} * (g_{\text{DAC}} * g_{\text{DAC}}))(\tau' + T') \\ &\quad + \text{j} (\mathbf{C}_{\Im\{\mathbf{t}\}\Re\{\mathbf{t}\}} * (g_{\text{DAC}} * g_{\text{DAC}}))(\tau' - T'). \end{aligned} \quad (9.36)$$

Therefore, we get

$$\begin{aligned} \rho_{y_{M+1}y_{M+1}}(\tau') &= (\mathbf{h}_{M+1}^T ((\Re\{\mathbf{C}_{tt}\} * (g_{DAC} * g_{DAC}))(0) + j(\Im\{\mathbf{C}_{tt}\} * (g_{DAC} * g_{DAC}))(T')) \mathbf{h}_{M+1}^*)^{-1} \\ &\quad \mathbf{h}_{M+1}^T \left((\Re\{\mathbf{C}_{tt}\} * (g_{DAC} * g_{DAC}))(\tau') - j(\mathbf{C}_{\Re\{t\}\Im\{t\}} * (g_{DAC} * g_{DAC}))(\tau' + T') \right. \\ &\quad \left. + j(\mathbf{C}_{\Im\{t\}\Re\{t\}} * (g_{DAC} * g_{DAC}))(\tau' - T') \right) \mathbf{h}_{M+1}^*. \end{aligned} \quad (9.37)$$

9.2.2.1 Without Digital Pulse Shaping

In the case of no digital pulse shaping, the autocorrelation function simplifies to

$$\begin{aligned} \rho_{y_{M+1}y_{M+1}}(\tau') &= (\mathbf{h}_{M+1}^T (\Re\{\mathbf{C}_{tt}\} (g_{DAC} * g_{DAC}))(0) + j\Im\{\mathbf{C}_{tt}\} (g_{DAC} * g_{DAC})(T') \mathbf{h}_{M+1}^*)^{-1} \\ &\quad \mathbf{h}_{M+1}^T \left(\Re\{\mathbf{C}_{tt}\} (g_{DAC} * g_{DAC})(\tau') - j\mathbf{C}_{\Re\{t\}\Im\{t\}} (g_{DAC} * g_{DAC})(\tau' + T') \right. \\ &\quad \left. + j\mathbf{C}_{\Im\{t\}\Re\{t\}} (g_{DAC} * g_{DAC})(\tau' - T') \right) \mathbf{h}_{M+1}^*. \end{aligned} \quad (9.38)$$

Thus, the PSD can be computed as

$$S_{y_{M+1}y_{M+1}}(f) = \frac{|G_{DAC}(f)|^2 \mathbf{h}_{M+1}^T (\Re\{\mathbf{C}_{tt}\} + j\mathbf{C}_{\Im\{t\}\Re\{t\}} e^{-j2\pi f T'} - j\mathbf{C}_{\Re\{t\}\Im\{t\}} e^{j2\pi f T'}) \mathbf{h}_{M+1}^*}{\mathbf{h}_{M+1}^T ((g_{DAC} * g_{DAC})(0) \Re\{\mathbf{C}_{tt}\} + j(g_{DAC} * g_{DAC})(T') \Im\{\mathbf{C}_{tt}\}) \mathbf{h}_{M+1}^*}. \quad (9.39)$$

It holds that

$$\begin{aligned} \mathbf{h}_{M+1}^T \mathbf{C}_{\Re\{t\}\Im\{t\}} e^{j2\pi f T'} \mathbf{h}_{M+1}^* &= \left(\mathbf{h}_{M+1}^T \mathbf{C}_{\Re\{t\}\Im\{t\}} e^{j2\pi f T'} \mathbf{h}_{M+1}^* \right)^T \\ &= \mathbf{h}_{M+1}^H \mathbf{C}_{\Im\{t\}\Re\{t\}} e^{j2\pi f T'} \mathbf{h}_{M+1} \\ &= \left(\mathbf{h}_{M+1}^T \mathbf{C}_{\Im\{t\}\Re\{t\}} e^{-j2\pi f T'} \mathbf{h}_{M+1}^* \right)^*, \end{aligned} \quad (9.40)$$

which leads to the simplified expression

$$\begin{aligned} S_{y_{M+1}y_{M+1}}(f) &= \frac{|G_{DAC}(f)|^2 (\mathbf{h}_{M+1}^T \Re\{\mathbf{C}_{tt}\} \mathbf{h}_{M+1}^* - 2\Im\{\mathbf{h}_{M+1}^T \mathbf{C}_{\Im\{t\}\Re\{t\}} \mathbf{h}_{M+1}^* e^{-j2\pi f T'}\})}{\mathbf{h}_{M+1}^T ((g_{DAC} * g_{DAC})(0) \Re\{\mathbf{C}_{tt}\} + j(g_{DAC} * g_{DAC})(T') \Im\{\mathbf{C}_{tt}\}) \mathbf{h}_{M+1}^*} \\ &\approx |G_{DAC}(f)|^2 \left(\mathbf{h}_{M+1}^T \Xi (\alpha_q \Re\{\mathbf{R}_{xx}\} + \beta_q \mathbf{I}_N) \Xi \mathbf{h}_{M+1}^* \right. \\ &\quad \left. - 2\alpha_q \Im\left\{ \mathbf{h}_{M+1}^T \Xi \text{diag}(\mathbf{C}_{xx})^{-1/2} \mathbf{C}_{\Im\{x\}\Re\{x\}} \text{diag}(\mathbf{C}_{xx})^{-1/2} \Xi \mathbf{h}_{M+1}^* e^{-j2\pi f T'} \right\} \right) \\ &\quad \left(\mathbf{h}_{M+1}^T \Xi \left((g_{DAC} * g_{DAC})(0) (\alpha_q \Re\{\mathbf{R}_{xx}\} + \beta_q \mathbf{I}_N) \right. \right. \\ &\quad \left. \left. + j\alpha_q (g_{DAC} * g_{DAC})(T') \Im\{\mathbf{R}_{xx}\} \right) \Xi \mathbf{h}_{M+1}^* \right)^{-1}. \end{aligned} \quad (9.41)$$

In Fig. 9.4, the evaluation of the PSD expression in (9.41) is depicted for both time impulse responses for the DAC. We can observe that the spectral shape changes slightly with increasing the ratio N/M , in contrast to the case of standard QAM without pulse shaping, where no dependency on N and M was observed.

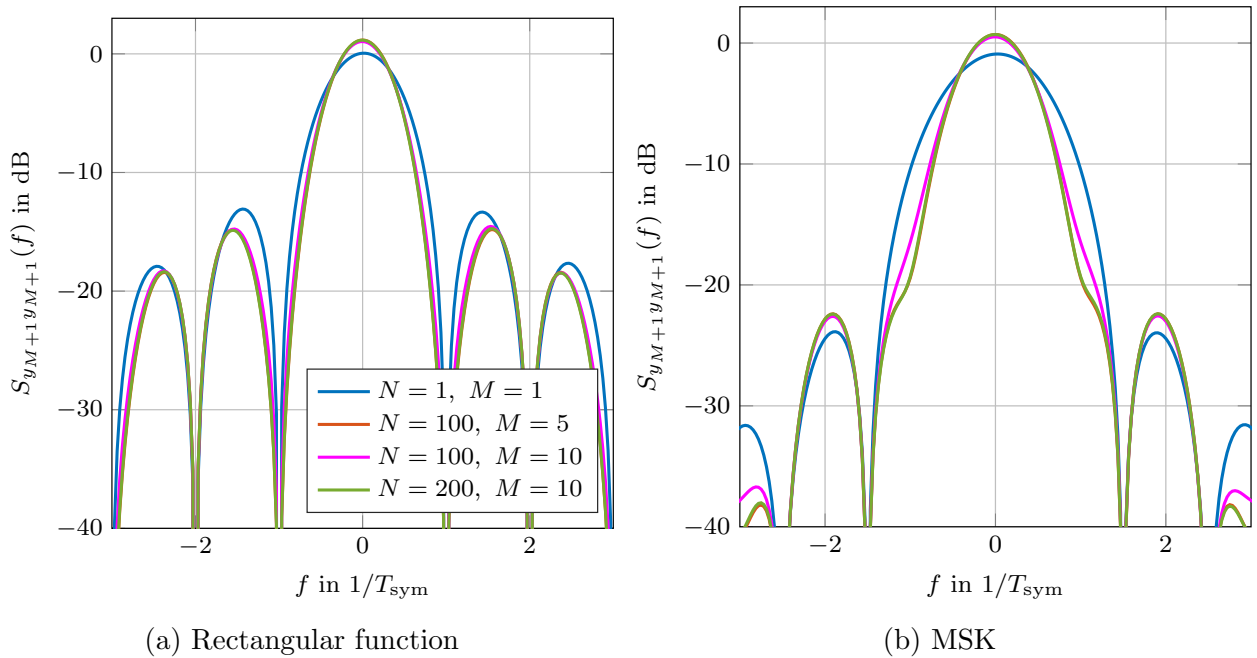


Fig. 9.4: PSD at one of the users for an i.i.d. channel realization with the MSM precoder in Offset QAM staggering without digital pulse shaping.

9.2.2.2 With Digital Pulse Shaping

We use the same identity as in (9.14). By applying the first order Taylor approximation of the arcsin function, we obtain

$$\begin{aligned}
\mathbb{E} \left[\mathbf{t}'(t') \mathbf{t}'^H(t' - \tau') \right] &\approx \sum_{\ell=-\infty}^{\infty} \mathbf{L}_{\mathcal{Q}} \Re \{ \mathbf{C}_{\mathbf{x}\mathbf{x}}[\ell] \} \mathbf{L}_{\mathcal{Q}}^H (g_{\text{DAC}} * g_{\text{DAC}}) (\tau' - 2\ell T_s) \\
&\quad + \beta_q \mathbf{\Xi}^2 (g_{\text{DAC}} * g_{\text{DAC}}) (\tau) \\
&\quad - j \sum_{\ell'=-\infty}^{\infty} \mathbf{L}_{\mathcal{Q}} \mathbf{C}_{\Re\{\mathbf{x}\}\Im\{\mathbf{x}\}}[\ell'] \mathbf{L}_{\mathcal{Q}}^H (g_{\text{DAC}} * g_{\text{DAC}}) (\tau' - (2\ell' - 1)T_s) \\
&\quad + j \sum_{\ell''=-\infty}^{\infty} \mathbf{L}_{\mathcal{Q}} \mathbf{C}_{\Im\{\mathbf{x}\}\Re\{\mathbf{x}\}}[\ell''] \mathbf{L}_{\mathcal{Q}}^H (g_{\text{DAC}} * g_{\text{DAC}}) (\tau' - (2\ell'' + 1)T_s) \\
&= \frac{1}{(g_{\text{PS}} * g_{\text{PS}})[0]} \alpha_q \mathbf{\Xi} \text{diag}(\mathbf{C}_{\mathbf{x}'\mathbf{x}'})^{-1/2} \\
&\quad \left(\Re \{ \mathbf{C}_{\mathbf{x}'\mathbf{x}'} \} (g_{\text{DAC}} * g_{\text{DAC}} * g_{\text{PS}} * g_{\text{PS}}) (\tau) \right. \\
&\quad \left. - j \mathbf{C}_{\Re\{\mathbf{x}'\}\Im\{\mathbf{x}'\}} (g_{\text{DAC}} * g_{\text{DAC}} * g_{\text{PS}} * g_{\text{PS}}) (\tau' + T_s) \right. \\
&\quad \left. + j \mathbf{C}_{\Im\{\mathbf{x}'\}\Re\{\mathbf{x}'\}} (g_{\text{DAC}} * g_{\text{DAC}} * g_{\text{PS}} * g_{\text{PS}}) (\tau' - T_s) \right) \\
&\quad \text{diag}(\mathbf{C}_{\mathbf{x}'\mathbf{x}'})^{-1/2} \mathbf{\Xi} \\
&\quad + \beta_q \mathbf{\Xi}^2 (g_{\text{DAC}} * g_{\text{DAC}}) (\tau'), \tag{9.42}
\end{aligned}$$

which leads to

$$\begin{aligned}
\rho_{y_{M+1}y_{M+1}}(\tau') &\approx \left(\mathbf{h}_{M+1}^T \mathbf{\Xi} \left(\Re \{ \mathbf{R}_{\mathbf{x}'\mathbf{x}'} \} (g_{\text{PS}} * g_{\text{PS}} * g_{\text{DAC}} * g_{\text{DAC}}) (0) \right. \right. \\
&\quad \left. \left. + \beta_q / \alpha_q (g_{\text{DAC}} * g_{\text{DAC}}) (0) (g_{\text{PS}} * g_{\text{PS}}) (0) \mathbf{I}_N \right) \right. \\
&\quad \left. + j \Im \{ \mathbf{R}_{\mathbf{x}'\mathbf{x}'} \} (g_{\text{PS}} * g_{\text{PS}}) [0] * (g_{\text{DAC}} * g_{\text{DAC}}) (T') \right) \mathbf{\Xi} \mathbf{h}_{M+1}^* \Big)^{-1} \\
&\quad \mathbf{h}_{M+1}^T \mathbf{\Xi} \left(\Re \{ \mathbf{R}_{\mathbf{x}'\mathbf{x}'} \} (g_{\text{PS}} * g_{\text{PS}} * g_{\text{DAC}} * g_{\text{DAC}}) (\tau') \right. \\
&\quad \left. - j \text{diag}(\mathbf{C}_{\mathbf{x}'\mathbf{x}'})^{-1/2} \mathbf{C}_{\Re\{\mathbf{x}'\}\Im\{\mathbf{x}'\}} \text{diag}(\mathbf{C}_{\mathbf{x}'\mathbf{x}'})^{-1/2} ((g_{\text{PS}} * g_{\text{PS}}) * (g_{\text{DAC}} * g_{\text{DAC}})) (\tau' + T') \right. \\
&\quad \left. + j \text{diag}(\mathbf{C}_{\mathbf{x}'\mathbf{x}'})^{-1/2} \mathbf{C}_{\Im\{\mathbf{x}'\}\Re\{\mathbf{x}'\}} \text{diag}(\mathbf{C}_{\mathbf{x}'\mathbf{x}'})^{-1/2} ((g_{\text{PS}} * g_{\text{PS}}) * (g_{\text{DAC}} * g_{\text{DAC}})) (\tau' - T') \right) \\
&\quad \mathbf{\Xi} \mathbf{h}_{M+1}^* \\
&\quad + \frac{\beta_q}{\alpha_q} (g_{\text{PS}} * g_{\text{PS}}) [0] (g_{\text{DAC}} * g_{\text{DAC}}) (\tau') \mathbf{h}_{M+1}^T \mathbf{\Xi}^2 \mathbf{h}_{M+1}^*. \tag{9.43}
\end{aligned}$$

Thus, the PSD can be approximated by

$$\begin{aligned}
 S_{y_{M+1}y_{M+1}}(f) \approx & |G_{\text{DAC}}(f)|^2 \left(\mathbf{h}_{M+1}^T \mathbf{\Xi} |G_{\text{PS}}(f)|^2 \left(\Re \{ \mathbf{R}_{\mathbf{x}'\mathbf{x}'} \} \right. \right. \\
 & - j \text{diag}(\mathbf{C}_{\mathbf{x}'\mathbf{x}'})^{-1/2} \mathbf{C}_{\Re\{\mathbf{x}'\}\Im\{\mathbf{x}'\}} \text{diag}(\mathbf{C}_{\mathbf{x}'\mathbf{x}'})^{-1/2} e^{j2\pi f T'} \\
 & \left. \left. + j \text{diag}(\mathbf{C}_{\mathbf{x}'\mathbf{x}'})^{-1/2} \mathbf{C}_{\Im\{\mathbf{x}'\}\Re\{\mathbf{x}'\}} \text{diag}(\mathbf{C}_{\mathbf{x}'\mathbf{x}'})^{-1/2} e^{-j2\pi f T'} \right) \mathbf{\Xi} \mathbf{h}_{M+1}^* \right. \\
 & + \frac{\beta_q}{\alpha_q} (g_{\text{PS}} * g_{\text{PS}}) [0] \mathbf{h}_{M+1}^T \mathbf{\Xi}^2 \mathbf{h}_{M+1} \\
 & \left(\mathbf{h}_{M+1}^T \mathbf{\Xi} \left(\Re \{ \mathbf{R}_{\mathbf{x}'\mathbf{x}'} \} (g_{\text{PS}} * g_{\text{PS}} * g_{\text{DAC}} * g_{\text{DAC}}) (0) \right. \right. \\
 & + \beta_q / \alpha_q (g_{\text{DAC}} * g_{\text{DAC}}) (0) (g_{\text{PS}} * g_{\text{PS}}) [0] \mathbf{I}_N \\
 & \left. \left. + j \Im \{ \mathbf{R}_{\mathbf{x}'\mathbf{x}'} \} (g_{\text{PS}} * g_{\text{PS}} * g_{\text{DAC}} * g_{\text{DAC}}) (T') \right) \mathbf{\Xi} \mathbf{h}_{M+1}^* \right)^{-1}. \quad (9.44)
 \end{aligned}$$

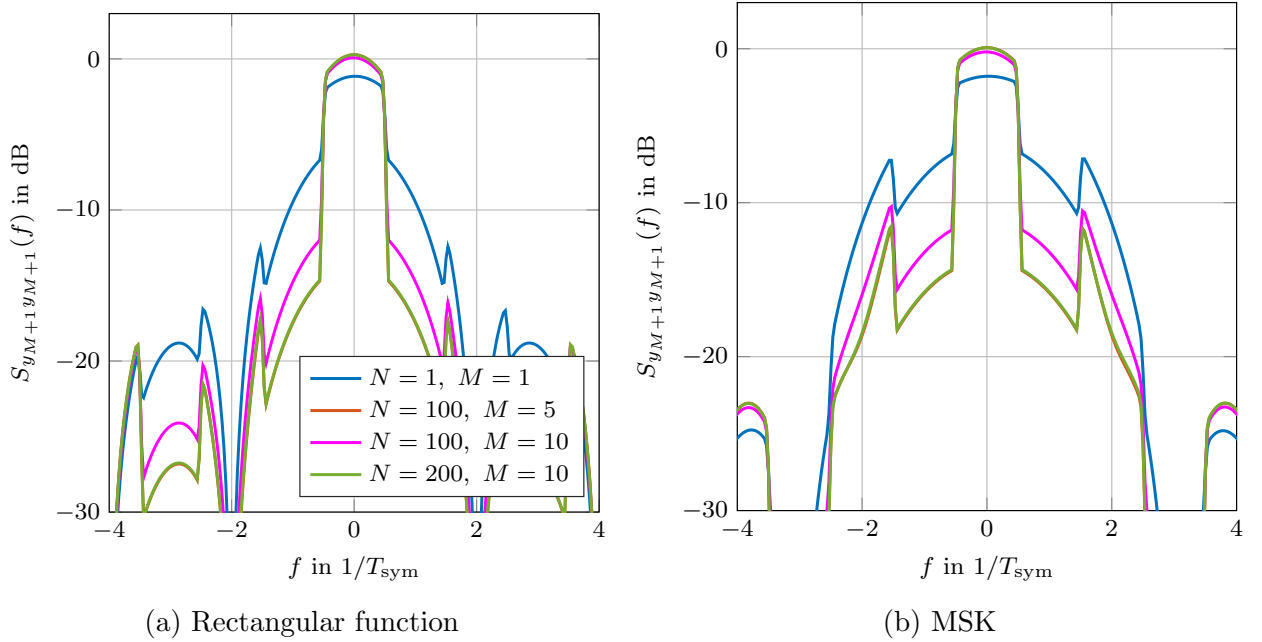


Fig. 9.5: PSD at one of the users for an i.i.d. channel realization with the MSM precoder in Offset QAM staggering with the RRC digital pulse shaping of $ro = 0.1$.

By choosing the digital pulse shaper as the RRC filter of roll-off factor $ro = 0.1$, we obtain the PSDs depicted in Fig. 9.5 for different time impulse responses for the DAC and different values of the ratio N/M . It can be again observed that the spectrum is narrower for larger values of the ratio N/M . Additionally, the spectrum with the rectangular function and in Offset-QAM staggering has less out of band radiations compared to the case of standard QAM case.

Additionally, the PSD is plotted in Fig. 9.6 for a random \mathbf{h}_{M+1} , where $\mathbf{h}_{M+1} \in \text{range}(\mathbf{R}_{\mathbf{x}'\mathbf{x}'})$.

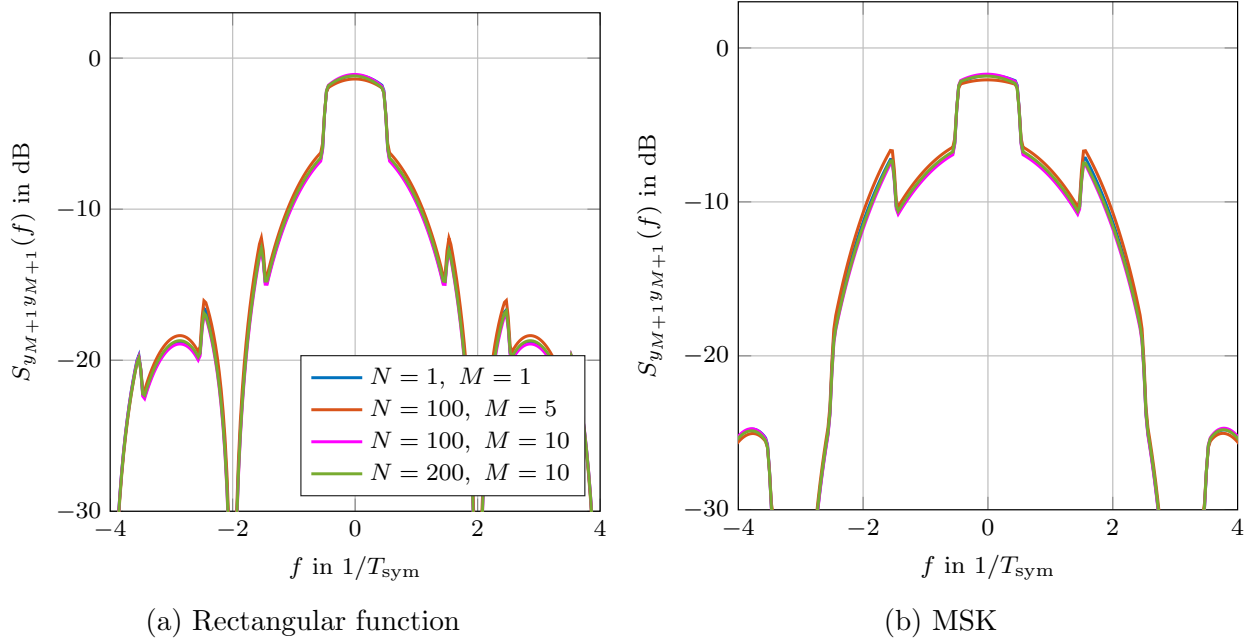


Fig. 9.6: PSD for a random $\mathbf{h}_{M+1} \in \text{range}(\mathbf{R}_{\mathbf{x}'\mathbf{x}'})$ for an i.i.d. channel realization with the MSM precoder in Offset QAM staggering with the RRC digital pulse shaping of $ro = 0.1$.

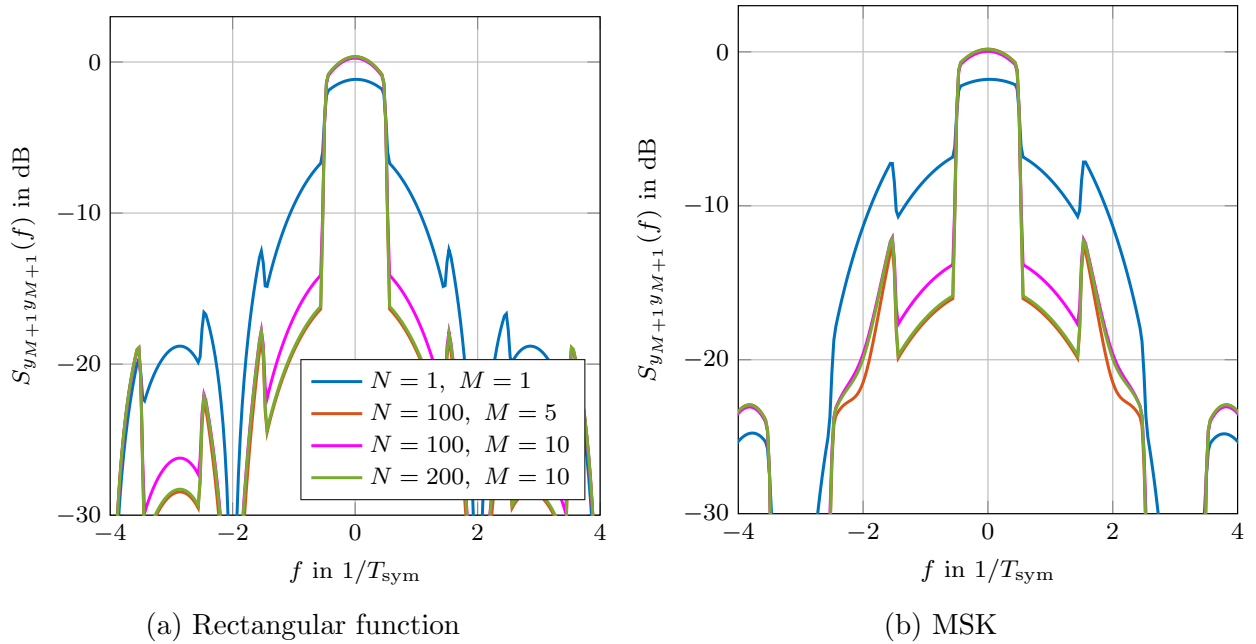


Fig. 9.7: PSD for \mathbf{h}_{M+1} being the eigenvector with the largest eigenvalue of $\mathbf{R}_{\mathbf{x}'\mathbf{x}'}$ for an i.i.d. channel realization with the MSM precoder in Offset QAM staggering with the RRC digital pulse shaping of $ro = 0.1$.

there is no improvement in the spectral shape for large ratios N/M . This is due to the fact that the signals in that direction do not overlap constructively. The precoder designs the transmit signals such that they overlap constructively at the users.

Finally, the PSD is plotted for \mathbf{h}_{M+1} being the eigenvector with the largest eigenvalue of

$\mathbf{R}_{\mathbf{x}'\mathbf{x}'}$ in Fig. 9.7. We can see the best spectral shape that can be obtained, which improves with the ratio N/M .

In summary, we have shown that the spectral shape at the measuring device, when \mathbf{h}_{M+1} is one of the user channels or one of the eigenvectors with large eigenvalues of $\mathbf{R}_{\mathbf{x}'\mathbf{x}'}$, improves by increasing the ratio N/M . Indeed, the spectral shape deterioration due to the coarse quantization diminishes in massive MIMO systems and is not that dramatic as in the SU Single-Input Single-Output (SISO) system. By using the Offset-QAM staggering combined with the rectangular or MSK impulse response, we ensure that the input of the PA is of constant envelope. Therefore, only non-linear and thus highly power efficient PAs are required in the system.

We would like to remind the reader that our analysis was restricted to the case of $Q = 4$. Our goal was to show the potential of spectral shaping despite the presence of coarse quantization in massive MIMO scenarios. The optimization problem of the PSD is left for future work. However, we would like to share our future steps with the reader.

In order to further improve the spectral properties in QCE systems, we think that the desired pulse shaping should be considered in the precoder design. Therefore, each user's signal should undergo a pulse shaping. Hence, the input to the precoder is not drawn from a discrete well-defined set \mathbb{S} anymore but from a continuous of high resolution undefined set. We suggest to quantize the resulting continuous set to obtain a well-defined discrete set \mathbb{S}' at the input of the precoder. Then, the MSM precoder should be extended to make sure that the received signals at the receiver belong to the new defined set. In other words, the new SRs should be defined that correspond to the new set \mathbb{S}' . Consequently, the received signals would follow the same pulse shape that is designed at the transmitter.

Another idea would be to optimize the pulse shaping after the CEQ. Therefore, to ensure the CE property, we can apply phase modulation for the transmission. Hence, we can consider the concepts of the Continuous Phase Modulation (CPM), [74], in particular the Gaussian Minimum-Shift Keying (GMSK) modulation [75]. This modulation type ensures a smooth change of the phase in time. This enhances the spectral efficiency of the transmit signals.

10. Conclusion

Throughout this thesis, we developed digital signal processing techniques for MU massive MIMO systems, whose transmit signals are restricted to be of CE. Additionally, the conversion from the digital into the analog world is deployed with very low resolution. The introduced CEQ models this behavior. Hence, we end up with coarsely QCE MU massive MIMO systems.

The choice of CE signaling at the BS antennas is motivated by having high power efficient PAs. In the case of massive MIMO systems, where the number of BS antennas is assumed to be very large, e.g. 100, the PA power consumption becomes a crucial concern especially at mmW frequencies. Therefore, we enhance the power efficiency of the communication systems by the CE constraint.

The choice of the coarse quantization is also motivated by having less power consuming DACs, whose number increases linearly with the number of the BS antennas and whose power consumption increases exponentially with the resolution. Other benefits of low resolution DACs are the reduced required chip area and the relaxed need for highly performing surrounding circuitry.

The combination of both properties leads to power efficient QCE MU massive MIMO systems.

However, this results in performance degradation, since beside to the MUI, the channel distortions and the AWGN, additional CEQ distortions are introduced to the system. Therefore, we proposed digital precoding techniques to mitigate all usual distortion sources and additionally the CEQ distortions. The proposed techniques are classified into two groups: linear and non-linear techniques.

The linear precoding designs were based on the MMSE criterion. The precoding matrix is optimized for every channel coherence time. For the matrix design, the CEQ was linearized based on Bussgang's theorem and the statistical properties of the quantization noise were first computed with Price's theorem and second approximated with the LCA. Furthermore, we have shown that there is no MSE duality between uplink and downlink in the case of QCE MIMO systems. Two linear precoder designs were proposed based on a virtual duality and an approximate duality. The simulation results showed that the proposed linear techniques are moderately better, in terms of the uncoded BER, than the well known linear WF followed by the CEQ. This observation made us think about non-linear precoding.

In the non-linear precoding approach, denoted by MSM, the transmit vector at each time instant is optimized. Therefore, we talk about symbol-wise precoding that depend on the desired received signal at the users and the channel realization. This approach is certainly more computationally complex than the linear approach. The design criterion is the safety

margin to the decision thresholds that should be maximized to minimize the SER. In order to consider the QCE constraint, the entries of the designed transmit vector should belong to a relaxed convex version of the QCE constraint. The problem could be formulated such that we obtained a linear programming problem. The simulation results showed significant gains compared to the linear approach. However, this gain comes at the cost of more computational complexity.

The main motivation of choosing the QCE is the power efficiency. We have shown that the PA power consumption of an ideal linear system is more than four times the PA power consumption of a QCE for the same uncoded BER performance.

To settle the doubts about the spectral regrowth combined with the coarse quantization, we analyzed the PSDs at the users in QCE MU massive MIMO systems. We have observed that the quantization widens the spectrum. However, this effect diminishes for larger ratios between the number of antennas and users, which is the case in massive MIMO systems. In this analysis, we did not propose an optimal solution for the time pulse shaping in order to optimize the PSD. However, we have shown that the problem of the spectral regrowth in the presence of quantization distortions is not that dramatic as in the case of SU SISO systems. This means that there is potential for optimization, which is left for future work.

Appendix

A1 Complex-Valued Gaussian Joint PDF

Let us assume that x_i and x_j with $i \neq j$, are complex-valued Gaussian random variables with zero means and positive variances $\sigma_{x_i}^2$ and $\sigma_{x_j}^2$; that is the PDFs are given by

$$p_{X_{i/j}}(x_{i/j}) = \frac{1}{\pi \sigma_{x_{i/j}}^2} e^{-\frac{|x_{i/j}|^2}{\sigma_{x_{i/j}}^2}}.$$

According to [55, Theorem 2.16], the random variable x_i given x_j is a complex-valued Gaussian random variable with mean value

$$\mathbb{E}[x_i|x_j] = \frac{\mathbb{E}[x_i x_j^*]}{\sigma_{x_j}^2} x_j$$

and variance

$$\mathbb{E}[|x_i - \mathbb{E}[x_i|x_j]|^2 | x_j] = \sigma_{x_i}^2 - \frac{\mathbb{E}[x_i x_j^*]}{\sigma_{x_j}^2} \mathbb{E}[x_j^* x_i].$$

After introducing the correlation coefficient factor $\rho_{x_i, x_j} = \frac{\mathbb{E}[x_i x_j^*]}{\sigma_{x_i} \sigma_{x_j}}$, it holds that

$$x_i|x_j \sim \mathcal{N}_{\mathbb{C}}\left(\rho_{x_i, x_j} \frac{\sigma_{x_i}}{\sigma_{x_j}} x_j, \sigma_{x_i}^2 \left(1 - |\rho_{x_i, x_j}|^2\right)\right).$$

Hence, the PDF of x_i given x_j is given by

$$p_{X_i|X_j}(x_i|x_j) = \frac{1}{\pi \sigma_{x_i}^2 \left(1 - |\rho_{x_i, x_j}|^2\right)} e^{-\left(\frac{|x_i - \rho_{x_i, x_j} \frac{\sigma_{x_i}}{\sigma_{x_j}} x_j|^2}{\sigma_{x_i}^2 \left(1 - |\rho_{x_i, x_j}|^2\right)}\right)}.$$

Finally, the joint PDF of x_i and x_j is obtained by

$$\begin{aligned} p_{X_i, X_j}(x_i, x_j) &= p_{X_i|X_j}(x_i|x_j) p_{X_j}(x_j) \\ &= \frac{1}{\pi^2 \left(1 - |\rho_{x_i, x_j}|^2\right) \sigma_{x_i}^2 \sigma_{x_j}^2} e^{-\frac{1}{1 - |\rho_{x_i, x_j}|^2} \left(\frac{|x_i|^2}{\sigma_{x_i}^2} + \frac{|x_j|^2}{\sigma_{x_j}^2} - \frac{2\Re\{\rho_{x_i, x_j} x_j x_i^*\}}{\sigma_{x_i} \sigma_{x_j}}\right)}. \end{aligned}$$

A2 Proof of Theorem 1

We recall the modified version of Price's theorem in [58].

Theorem 3 (Pawula's Theorem [58]). *Let $f(\mathbf{s}) = f(s_1, \dots, s_n)$ be a function of n real-valued joint Gaussian random variables. The covariance matrix and the correlation matrix of the vector \mathbf{s} are denoted by $\mathbf{C}_{\mathbf{ss}}$ and $\mathbf{R}_{\mathbf{ss}}$ of elements denoted by ρ_{s_i, s_j} , $i, j = 1, \dots, n$. The off-diagonal entries of $\mathbf{C}_{\mathbf{ss}}$ are multiplied by a perturbing term v . Then*

$$\frac{d^\ell \mathbb{E}[f(s_1, \dots, s_n)]}{d v^\ell} = \mathbb{E}_v \left[\sum_{i < j} \rho_{s_i, s_j}^\ell \sigma_{s_i}^\ell \sigma_{s_j}^\ell \frac{\partial^{2\ell} f(s_1, \dots, s_n)}{\partial s_i^\ell \partial s_j^\ell} \right], \quad s_i, s_j \in \mathbb{R}, \quad i, j = 1, \dots, n,$$

where \mathbb{E}_v denotes the expectation based on the resulting perturbed PDF.

Theorem 3 can be applied to two complex-valued Gaussian random variables $x_i = x_{iR} + j x_{iI}$ and $x_j = x_{jR} + j x_{jI}$ as follows

$$\begin{aligned} \frac{d^\ell \mathbb{E}[f(x_{iR}, x_{iI}, x_{jR}, x_{jI})]}{d v^\ell} = & \mathbb{E}_v \left[\rho_{x_{iR}, x_{iI}}^\ell \sigma_{x_{iR}}^\ell \sigma_{x_{iI}}^\ell \frac{\partial^{2\ell} f(x_{iR}, x_{iI}, x_{jR}, x_{jI})}{\partial x_{iR}^\ell \partial x_{iI}^\ell} \right. \\ & + \rho_{x_{iR}, x_{jR}}^\ell \sigma_{x_{iR}}^\ell \sigma_{x_{jR}}^\ell \frac{\partial^{2\ell} f(x_{iR}, x_{iI}, x_{jR}, x_{jI})}{\partial x_{iR}^\ell \partial x_{jR}^\ell} \\ & + \rho_{x_{iR}, x_{jI}}^\ell \sigma_{x_{iR}}^\ell \sigma_{x_{jI}}^\ell \frac{\partial^{2\ell} f(x_{iR}, x_{iI}, x_{jR}, x_{jI})}{\partial x_{iR}^\ell \partial x_{jI}^\ell} \\ & + \rho_{x_{iI}, x_{jR}}^\ell \sigma_{x_{iI}}^\ell \sigma_{x_{jR}}^\ell \frac{\partial^{2\ell} f(x_{iR}, x_{iI}, x_{jR}, x_{jI})}{\partial x_{iI}^\ell \partial x_{jR}^\ell} \\ & + \rho_{x_{iI}, x_{jI}}^\ell \sigma_{x_{iI}}^\ell \sigma_{x_{jI}}^\ell \frac{\partial^{2\ell} f(x_{iR}, x_{iI}, x_{jR}, x_{jI})}{\partial x_{iI}^\ell \partial x_{jI}^\ell} \\ & \left. + \rho_{x_{jR}, x_{jI}}^\ell \sigma_{x_{jR}}^\ell \sigma_{x_{jI}}^\ell \frac{\partial^{2\ell} f(x_{iR}, x_{iI}, x_{jR}, x_{jI})}{\partial x_{jR}^\ell \partial x_{jI}^\ell} \right]. \end{aligned}$$

While assuming circularly symmetric distributed Gaussian signals x_i and x_j , i.e.

$$\begin{aligned} \sigma_{x_{iR}} = \sigma_{x_{iI}} = \frac{1}{\sqrt{2}} \sigma_{x_i}, & \quad \sigma_{x_{jR}} = \sigma_{x_{jI}} = \frac{1}{\sqrt{2}} \sigma_{x_j}, \\ \rho_{x_{iR}, x_{jR}} = \rho_{x_{iI}, x_{jI}} = \Re\{\rho_{x_i, x_j}\}, & \quad \rho_{x_{iI}, x_{jR}} = -\rho_{x_{iR}, x_{jI}} = \Im\{\rho_{x_i, x_j}\}, \\ \rho_{x_{iR}, x_{iI}} = 0, & \quad \rho_{x_{jR}, x_{jI}} = 0, \end{aligned}$$

we get

$$\begin{aligned} \frac{d^\ell \mathbb{E}[f(x_{iR}, x_{iI}, x_{jR}, x_{jI})]}{d v^\ell} = & \frac{1}{2} \sigma_{x_i}^\ell \sigma_{x_j}^\ell \cdot \\ & \mathbb{E}_v \left[\Re\{\rho_{x_i, x_j}\}^\ell \left(\frac{\partial^{2\ell} f(x_{iR}, x_{iI}, x_{jR}, x_{jI})}{\partial x_{iR}^\ell \partial x_{jR}^\ell} + \frac{\partial^{2\ell} f(x_{iR}, x_{iI}, x_{jR}, x_{jI})}{\partial x_{iI}^\ell \partial x_{jI}^\ell} \right) \right. \\ & \left. + \Im\{\rho_{x_i, x_j}\}^\ell \left(\frac{\partial^{2\ell} f(x_{iR}, x_{iI}, x_{jR}, x_{jI})}{\partial x_{iI}^\ell \partial x_{jR}^\ell} - \frac{\partial^{2\ell} f(x_{iR}, x_{iI}, x_{jR}, x_{jI})}{\partial x_{iR}^\ell \partial x_{jI}^\ell} \right) \right]. \end{aligned} \tag{A1}$$

The following equalities

$$\begin{aligned}\frac{\partial}{\partial x_R} &= \frac{\partial}{\partial x} + \frac{\partial}{\partial x^*}, \\ \frac{\partial}{\partial x_I} &= j \left(\frac{\partial}{\partial x} - \frac{\partial}{\partial x^*} \right)\end{aligned}$$

imply that

$$\begin{aligned}\frac{\partial^2}{\partial x_{iR} \partial x_{jR}} &= \frac{\partial^2}{\partial x_i \partial x_j} + \frac{\partial^2}{\partial x_i \partial x_j^*} + \frac{\partial^2}{\partial x_i^* \partial x_j} + \frac{\partial^2}{\partial x_i^* \partial x_j^*}, \\ \frac{\partial^2}{\partial x_{iI} \partial x_{jI}} &= -\frac{\partial^2}{\partial x_i \partial x_j} + \frac{\partial^2}{\partial x_i \partial x_j^*} + \frac{\partial^2}{\partial x_i^* \partial x_j} - \frac{\partial^2}{\partial x_i^* \partial x_j^*}, \\ \frac{\partial^2}{\partial x_{iI} \partial x_{jR}} &= j \left(\frac{\partial^2}{\partial x_i \partial x_j} + \frac{\partial^2}{\partial x_i \partial x_j^*} - \frac{\partial^2}{\partial x_i^* \partial x_j} - \frac{\partial^2}{\partial x_i^* \partial x_j^*} \right), \\ \frac{\partial^2}{\partial x_{iR} \partial x_{jI}} &= j \left(\frac{\partial^2}{\partial x_i \partial x_j} - \frac{\partial^2}{\partial x_i \partial x_j^*} + \frac{\partial^2}{\partial x_i^* \partial x_j} - \frac{\partial^2}{\partial x_i^* \partial x_j^*} \right).\end{aligned}$$

Thus, (A1) simplifies to

$$\begin{aligned}\frac{d^\ell E[f(x_i, x_j)]}{dv^\ell} &= \sigma_{x_i}^\ell \sigma_{x_j}^\ell E_v \left[\Re\{\rho_{x_i, x_j}\}^\ell \left(\frac{\partial^{2\ell} f(x_i, x_j)}{\partial x_i^\ell \partial x_j^{*\ell}} + \frac{\partial^{2\ell} f(x_i, x_j)}{\partial x_i^{*\ell} \partial x_j^\ell} \right) \right. \\ &\quad \left. + j^\ell \Im\{\rho_{x_i, x_j}\}^\ell \left(\frac{\partial^{2\ell} f(x_i, x_j)}{\partial x_i^\ell \partial x_j^{*\ell}} - \frac{\partial^{2\ell} f(x_i, x_j)}{\partial x_i^{*\ell} \partial x_j^\ell} \right) \right].\end{aligned}$$

A3 Second Order Derivative in the Polar Representation

$$\begin{aligned}\frac{\partial^2 f(x_i, x_j)}{\partial x_i \partial x_j^*} &= \frac{\partial}{\partial x_i} \left(\frac{\partial f(x_i, x_j)}{\partial x_j^*} \right) \\ &= \frac{\partial}{\partial x_i} \left(\frac{\partial f(x_i, x_j)}{\partial r_j} \frac{\partial r_j}{\partial x_j^*} + \frac{\partial f(x_i, x_j)}{\partial \phi_j} \frac{\partial \phi_j}{\partial x_j^*} \right) \\ &= \frac{\partial^2 f(x_i, x_j)}{\partial r_i \partial r_j} \frac{\partial r_i}{\partial x_i} \frac{\partial r_j}{\partial x_j^*} + \frac{\partial^2 f(x_i, x_j)}{\partial \phi_i \partial r_j} \frac{\partial \phi_i}{\partial x_i} \frac{\partial r_j}{\partial x_j^*} + \frac{\partial^2 f(x_i, x_j)}{\partial r_i \partial \phi_j} \frac{\partial r_i}{\partial x_i} \frac{\partial \phi_j}{\partial x_j^*} + \frac{\partial^2 f(x_i, x_j)}{\partial \phi_i \partial \phi_j} \frac{\partial \phi_i}{\partial x_i} \frac{\partial \phi_j}{\partial x_j^*} \\ &= \begin{bmatrix} \frac{\partial r_i}{\partial x_i} & \frac{\partial \phi_i}{\partial x_i} \end{bmatrix} \begin{bmatrix} \frac{\partial^2 f(x_i, x_j)}{\partial r_i \partial r_j} & \frac{\partial^2 f(x_i, x_j)}{\partial r_i \partial \phi_j} \\ \frac{\partial^2 f(x_i, x_j)}{\partial \phi_i \partial r_j} & \frac{\partial^2 f(x_i, x_j)}{\partial \phi_i \partial \phi_j} \end{bmatrix} \begin{bmatrix} \frac{\partial r_j}{\partial x_j^*} \\ \frac{\partial \phi_j}{\partial x_j^*} \end{bmatrix}.\end{aligned}$$

Using the properties $r = \sqrt{xx^*}$ and $\phi = \arctan\left(\frac{-j(x-x^*)}{x+x^*}\right)$, we get

$$\begin{aligned}\frac{\partial^2 f(x_i, x_j)}{\partial x_i \partial x_j^*} &= \frac{1}{2} e^{-j\phi_i} \begin{bmatrix} 1 & -j \\ & r_i \end{bmatrix} \begin{bmatrix} \frac{\partial^2 f(x_i, x_j)}{\partial r_i \partial r_j} & \frac{\partial^2 f(x_i, x_j)}{\partial r_i \partial \phi_j} \\ \frac{\partial^2 f(x_i, x_j)}{\partial \phi_i \partial r_j} & \frac{\partial^2 f(x_i, x_j)}{\partial \phi_i \partial \phi_j} \end{bmatrix} \frac{1}{2} e^{j\phi_j} \begin{bmatrix} 1 \\ j \\ r_j \end{bmatrix} \\ &= \frac{1}{4} e^{-j(\phi_i - \phi_j)} \begin{bmatrix} 1 & -j \\ & r_i \end{bmatrix} \begin{bmatrix} \frac{\partial^2 f(x_i, x_j)}{\partial r_i \partial r_j} & \frac{\partial^2 f(x_i, x_j)}{\partial r_i \partial \phi_j} \\ \frac{\partial^2 f(x_i, x_j)}{\partial \phi_i \partial r_j} & \frac{\partial^2 f(x_i, x_j)}{\partial \phi_i \partial \phi_j} \end{bmatrix} \begin{bmatrix} 1 \\ j \\ r_j \end{bmatrix} \\ &= \frac{1}{4} e^{-j(\phi_i - \phi_j)} \left(\frac{\partial^2 f(x_i, x_j)}{\partial r_i \partial r_j} + \frac{1}{r_i r_j} \frac{\partial^2 f(x_i, x_j)}{\partial \phi_i \partial \phi_j} + j \frac{1}{r_j} \frac{\partial^2 f(x_i, x_j)}{\partial r_i \partial \phi_j} - j \frac{1}{r_i} \frac{\partial^2 f(x_i, x_j)}{\partial \phi_i \partial r_j} \right).\end{aligned}$$

A4 Special Integrals

The joint PDF of $x_i = r_i e^{j\phi_i}$ and $x_j = r_j e^{j\phi_j}$ for $i \neq j$ can be expressed in the polar coordinates as

$$p_{X_i, X_j, i \neq j}(r_i, r_j, \phi_i, \phi_j) = d e^{-(a^2 r_i^2 + b^2 r_j^2 - 2abc r_i r_j)},$$

where

$$\begin{aligned} a &= \frac{1}{\sqrt{1 - |\rho_{x_i, x_j}|^2} \sigma_{x_i}}, & b &= \frac{1}{\sqrt{1 - |\rho_{x_i, x_j}|^2} \sigma_{x_j}} \\ c &= \Re\{\rho_{x_i, x_j} e^{-j(\phi_i - \phi_j)}\}, & d &= \frac{1}{\pi^2 (1 - |\rho_{x_i, x_j}|^2) \sigma_{x_i}^2 \sigma_{x_j}^2}. \end{aligned}$$

In this section, two different integrals as functions of the joint PDF w.r.t. the radii r_i and r_j are computed, viz.,

- $w_1(\phi_i - \phi_j) = \int_0^\infty \int_0^\infty p_{X_i, X_j, i \neq j}(r_i, r_j, \phi_i, \phi_j) r_j dr_i dr_j$
- $w_2(\phi_i - \phi_j) = \int_0^\infty \int_0^\infty p_{X_i, X_j, i \neq j}(r_i, r_j, \phi_i, \phi_j) dr_i dr_j$

The first integral calculates to

$$\begin{aligned} w_1(\phi_i - \phi_j) &= \int_0^\infty \int_0^\infty p_{X_i, X_j, i \neq j}(r_i, r_j, \phi_i, \phi_j) r_j dr_i dr_j \\ &\stackrel{z = ar_i - bcr_j}{=} \frac{d}{a} \int_0^\infty r_j e^{-b^2(1-c^2)r_j^2} \int_{-bcr_j}^\infty e^{-z^2} dz dr_j \\ &= \sqrt{\frac{\pi}{4}} \frac{d}{a} \int_0^\infty r_j e^{-b^2(1-c^2)r_j^2} (1 + \operatorname{erf}(bcr_j)) dr_j \\ &\stackrel{x = \sqrt{2}b\sqrt{1-c^2}r_j}{=} \frac{d}{a} \frac{\pi}{\sqrt{2}b^2(1-c^2)} \frac{1}{\sqrt{2\pi}} \int_0^\infty x e^{-\frac{x^2}{2}} \frac{1}{2} \left(1 + \operatorname{erf}\left(\frac{c}{\sqrt{1-c^2}} \frac{x}{\sqrt{2}}\right) \right) dx. \end{aligned}$$

Using the equality (1,011.2) in [76], the first integral reduces to

$$\begin{aligned} w_1(\phi_i - \phi_j) &= \frac{\sqrt{1 - |\rho_{x_i, x_j}|^2}}{4\pi \sqrt{\pi} \sigma_{x_i} (1 - c)} \\ &= \frac{\sqrt{1 - |\rho_{x_i, x_j}|^2}}{4\pi \sqrt{\pi} \sigma_{x_i} (1 - \Re\{\rho_{x_i, x_j} e^{-j(\phi_i - \phi_j)}\})}. \end{aligned} \tag{A2}$$

Analogously, the second integral can be computed as

$$\begin{aligned} w_2(\phi_i - \phi_j) &= \int_0^\infty \int_0^\infty p_{X_i, X_j, i \neq j}(r_i, r_j, \phi_i, \phi_j) dr_i dr_j \\ &\stackrel{z = ar_i - bcr_j}{=} \frac{d}{a} \int_0^\infty e^{-b^2(1-c^2)r_j^2} \int_{-bcr_j}^\infty e^{-z^2} dz dr_j \\ &= \sqrt{\frac{\pi}{4}} \frac{d}{a} \int_0^\infty e^{-b^2(1-c^2)r_j^2} (1 + \operatorname{erf}(bcr_j)) dr_j \\ &\stackrel{x = \sqrt{2}b\sqrt{1-c^2}r_j}{=} \frac{d}{a} \frac{\pi}{b\sqrt{1-c^2}} \frac{1}{\sqrt{2\pi}} \int_0^\infty e^{-\frac{x^2}{2}} \frac{1}{2} \left(1 + \operatorname{erf}\left(\frac{c}{\sqrt{1-c^2}} \frac{x}{\sqrt{2}}\right) \right) dx. \end{aligned}$$

Using the equality (1,010.2) in [76], we obtain

$$\begin{aligned}
w_2(\phi_i - \phi_j) &= \frac{d}{a} \frac{\pi}{b\sqrt{1-c^2}} \left(\frac{1}{4} + \frac{1}{2\pi} \arcsin(c) \right) \\
&= \frac{1}{\pi\sigma_{x_i}\sigma_{x_j}\sqrt{1 - \Re\{\rho_{x_i,x_j} e^{-j(\phi_i-\phi_j)}\}}^2} \left(\frac{1}{4} + \frac{1}{2\pi} \arcsin(\Re\{\rho_{x_i,x_j} e^{-j(\phi_i-\phi_j)}\}) \right).
\end{aligned} \tag{A3}$$

A5 Derivative of arcsin(\mathbf{A})

Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be a quadratic matrix, whose non-diagonal elements depend on a parameter b and diagonal elements are constant. Then, the derivative of arcsin(\mathbf{A}) w.r.t. to b is given by

$$\begin{aligned}
\frac{\partial \arcsin(\mathbf{A})}{\partial b} &= \sum_{k=1}^N \sum_{l=1}^N \frac{\partial \arcsin(a_{kl})}{\partial b} \mathbf{e}_k \mathbf{e}_l^T \\
&= \sum_{k=1}^N \sum_{\substack{l=1 \\ l \neq k}}^N \frac{1}{\sqrt{1-a_{kl}^2}} \frac{\partial a_{kl}}{\partial b} \mathbf{e}_k \mathbf{e}_l^T, \quad \text{since } a_{kk} \text{ is constant, } k = 1, \dots, N \\
&= \text{nondiag} \left(\mathbf{1}_{N,N} - \text{nondiag}(\mathbf{A})^{\circ 2} \right)^{\circ -1/2} \circ \frac{\partial \mathbf{A}}{\partial b}.
\end{aligned}$$

Acronyms

ADC Analog-to-Digital Converter
AOD Angle of Departure
AWGN Additive White Gaussian Noise

BER Bit Error Ratio
BS Base Station

CE Constant Envelope
CEQ Constant Envelope Quantizer
CPM Continuous Phase Modulation
CSI Channel State Information

DAC Digital-to-Analog Converter
DTC Digital-to-Time Converter
DTFT Discrete Time Fourier Transform

FIR Finite Impulse Response

GMSK Gaussian Minimum-Shift Keying

HDTV High Definition Television

i.i.d. independent and identically distributed

IBI Inter-Block Interference
ISI Inter-Symbol Interference

KKT Karush-Kuhn-Tucker

LCA Linear Covariance Approximation

MIMO Multiple-Input Multiple-Output
MISO Multiple-Input Single-Output
MMSE Minimum Mean Squared Error

mmW Millimeter-Wave
MSE Mean Squared Error
MSK Minimum-Shift Keying
MSM Maximum Safety Margin
MU Multi-User
MUI Multi-User Interference

PA Power Amplifier
PAPR Peak-to-Average-Power Ratio
PDF Probability Density Function
PSD Power Spectral Density
PSK Phase-Shift Keying

QAM Quadrature Amplitude Modulation
QCE Quantized Constant Envelope

RF Radio Frequency
RRC Root-Raised-Cosine

SER Symbol Error Ratio
SINR Signal-to-Interference-Noise Ratio
SISO Single-Input Single-Output
SLNR Signal-to-Leakage-plus-Noise Ratio
SNR Signal-to-Noise Ratio
SR Symbol Region
SU Single-User

UHDV Ultra-High Definition Video

w.r.t. with respect to
WF Wiener Filter

List of Figures

3.1	Illustration of the CEQ output set for $q = 3$, i.e. $Q = 8$	12
3.2	Comparison of the derived closed-form expressions in (3.37) and (3.50) (solid lines) with the numerical results obtained by Monte-Carlo simulations (dashed lines with markers). In each figure, one dimension of ρ_{in} is fixed either $ \rho_{\text{in}} = 0.4$ or $\arg(\rho_{\text{in}}) = \frac{\pi}{4}$, © 2018 IEEE.	22
4.1	Downlink MU-MIMO system model.	31
4.2	Downlink system model with Bussgang decomposition.	38
4.3	Uplink system model with Bussgang decomposition.	38
5.1	Comparison between different linear precoders for 4-PSK signaling: $N = 64$, $M = 8$ with the i.i.d. channel.	53
5.2	Comparison between different linear precoders for 16-QAM signaling: $N = 128$, $M = 8$ with the i.i.d. channel.	53
5.3	Comparison between different linear precoders for 4-PSK signaling: $N = 64$, $M = 8$ with the mmW sparse channel.	54
5.4	Comparison between different linear precoders for 16-QAM signaling: $N = 128$, $M = 8$ with the mmW sparse channel.	54
6.1	Comparison between different linear precoders for 4-PSK signaling: $N = 64$, $M = 8$ with the i.i.d. channel with $L = 3$	68
6.2	Comparison between different linear precoders for 16-QAM signaling: $N = 128$, $M = 8$ with the i.i.d. channel with $L = 3$	68
6.3	Comparison between different linear precoders for 4-PSK signaling: $N = 64$, $M = 8$ with the mmW sparse channel with $L = 3$	69
6.4	Comparison between different linear precoders for 16-QAM signaling: $N = 128$, $M = 8$ with the mmW sparse channel with $L = 3$	69
7.1	Illustration of the relaxed polygon constraint for $Q=8$, © 2018 IEEE.	74
7.2	Decision regions and SRs (in red) for different constellations, © 2018 IEEE.	75
7.3	Illustration of the PSK SR in the modified coordinate system.	77
7.4	Illustration of the QAM receiver SR in the shifted and in the shifted and rotated coordinate system : $\xi_{1/2_m} \in \{2, \infty\}$	80
7.5	The noiseless received symbols at one arbitrary user m for an arbitrary i.i.d. channel realization with $N = 64$, $M = 8$ and $Q = 4$, © 2018 IEEE.	85

7.6	Uncoded BER performance for a MU-MIMO system with $N = 64$ and $M = 8$ for different precoding designs and 4-PSK signaling, © 2018 IEEE.	89
7.7	Uncoded BER performance for a MU-MIMO system with $N = 64$ and $M = 8$ for different modulation schemes: MSM (solid lines), the ideal WF (dashed lines).	91
7.8	Uncoded BER performance for a MU-MIMO system with $M = 8$ for different modulation schemes: MSM with $N = 128$ (solid lines), the ideal WF with $N = 64$ (dashed lines).	92
7.9	Comparison of the uncoded BER performance between MSM, CVX-CIO from [22] and CVX-CIO-noCE for a MU-MIMO system with $N = 64$ and $M = 8$, © 2018 IEEE.	93
8.1	BER performance for a MU-MIMO system with $N = 64$, $M = 8$ and $Q = 4$ with 16-QAM and the i.i.d. channel model.	101
8.2	BER performance for a MU-MIMO system with $N = 64$, $M = 8$ with 16-QAM signaling and channel exponential power delay profile with $L = 6$ for different quantization resolutions: MSM-BP-CP with $B = 16$	101
9.1	PSD at the measuring device with the DAC as the rectangular function in standard QAM staggering and without digital pulse shaping.	106
9.2	PSD at one of the users for an i.i.d. channel realization with the DAC as the rectangular function in standard QAM staggering and with digital pulse shaping, i.e RRC filter with roll-off factor $ro = 0.1$	109
9.3	PSD at the measuring device for an i.i.d. channel realization with the DAC as the rectangular function in standard QAM staggering, with digital pulse shaping, i.e RRC filter with roll-off factor $ro = 0.1$ and with the MSM precoder.	109
9.4	PSD at one of the users for an i.i.d. channel realization with the MSM precoder in Offset QAM staggering without digital pulse shaping.	113
9.5	PSD at one of the users for an i.i.d. channel realization with the MSM precoder in Offset QAM staggering with the RRC digital pulse shaping of $ro = 0.1$	115
9.6	PSD for a random $\mathbf{h}_{M+1} \in \text{range}(\mathbf{R}_{\mathbf{x}'\mathbf{x}'})$ for an i.i.d. channel realization with the MSM precoder in Offset QAM staggering with the RRC digital pulse shaping of $ro = 0.1$	116
9.7	PSD for \mathbf{h}_{M+1} being the eigenvector with the largest eigenvalue of $\mathbf{R}_{\mathbf{x}'\mathbf{x}'}$ for an i.i.d. channel realization with the MSM precoder in Offset QAM staggering with the RRC digital pulse shaping of $ro = 0.1$	116

List of Tables

3.1	Optimal step size for the CEQ and the corresponding distortion factor for unit variance inputs.	25
4.1	Exponential power delay profile with $L = 3$	35
4.2	Vehicular A (Case II) power delay profile with $L = 6$	35
7.1	Quantization distortion vs. B , © 2018 IEEE.	84
7.2	$\Delta\alpha$ and δ_{nom} vs. N and M : single α , © 2018 IEEE.	87
7.3	$\Delta\alpha$ and δ_{nom} vs. N and M : M α 's, © 2018 IEEE.	87
7.4	Average number of iterations of MSM for the i.i.d. channel model with $N = 64$ and $M = 8$, © 2018 IEEE.	92

Bibliography

- [1] A. Osseiran, F. Boccardi, V. Braun, K. Kusume, P. Marsch, M. Maternia, O. Queseth, M. Schellmann, H. Schotten, H. Taoka, H. Tullberg, M. A. Uusitalo, B. Timus, and M. Fallgren, “Scenarios for 5G Mobile and Wireless Communications: The vision of the METIS project,” *IEEE Communications Magazine*, vol. 52, no. 5, pp. 26–35, 2014.
- [2] T. L. Marzetta, “Noncooperative Cellular Wireless with Unlimited Numbers of Base Station Antennas,” *IEEE Transactions on Wireless Communications*, vol. 9, no. 11, pp. 3590–3600, 2010.
- [3] F. Rusek, D. Persson, B. K. Lau, E. G. Larsson, T. L. Marzetta, and F. Tufvesson, “Scaling Up MIMO: Opportunities and Challenges with Very Large Arrays,” *IEEE Signal Processing Magazine*, vol. 30, no. 1, pp. 40–60, 2013.
- [4] J. Hoydis, S. ten Brink, and M. Debbah, “Massive MIMO in the UL/DL of Cellular Networks: How Many Antennas Do We Need?” *IEEE Journal on Selected Areas in Communications*, vol. 31, no. 2, pp. 160–171, 2013.
- [5] H. Q. Ngo, E. G. Larsson, and T. L. Marzetta, “Energy and Spectral Efficiency of Very Large Multiuser MIMO Systems,” *IEEE Transactions on Communications*, vol. 61, no. 4, pp. 1436–1449, 2013.
- [6] L. Lu, G. Y. Li, A. L. Swindlehurst, A. Ashikhmin, and R. Zhang, “An Overview of Massive MIMO: Benefits and Challenges,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 8, no. 5, pp. 742–758, 2014.
- [7] A. L. Swindlehurst, E. Ayanoglu, P. Heydari, and F. Capolino, “Millimeter-Wave Massive MIMO: The Next Wireless Revolution?” *IEEE Communications Magazine*, vol. 52, no. 9, pp. 56–62, 2014.
- [8] Y. Niu, Y. Li, D. Jin, L. Su, and A. V. Vasilakos, “A Survey of Millimeter Wave Communications (mmWave) for 5G: Opportunities and Challenges,” *Wireless Networks*, vol. 21, no. 8, pp. 2657–2676, 2015.
- [9] T. S. Rappaport, *Millimeter Wave Wireless Communications*. Upper Saddle River N.J.: Prentice Hall, 2015.
- [10] “EARTH PROJECT INFISO-ICT-247733 EARTH Deliverable D2.3,” January 31, 2012.
- [11] O. Blume, D. Zeller, and U. Barth, “Approaches to Energy Efficient Wireless Access Networks,” in *Proc. 4th International Symposium on Communications, Control and Signal Processing (ISCCSP)*, 2010.
- [12] P. Varahram, S. Mohammady, B. M. Ali, and N. Sulaiman, *Power Efficiency in Broadband Wireless Communications*, 1st ed. Boca Raton: CRC Press, 2014.
- [13] S. Cui, A. J. Goldsmith, and A. Bahai, “Energy-Constrained Modulation Optimization,” *IEEE Transactions on Wireless Communications*, vol. 4, no. 5, pp. 2349–2360, 2005.

-
- [14] S. K. Mohammed and E. G. Larsson, "Single-User Beamforming in Large-Scale MISO Systems with Per-Antenna Constant-Envelope Constraints: The Doughnut Channel," *IEEE Transactions on Wireless Communications*, vol. 11, no. 11, pp. 3992–4005, 2012.
- [15] —, "Per-Antenna Constant Envelope Precoding for Large Multi-User MIMO Systems," *IEEE Transactions on Communications*, vol. 61, no. 3, pp. 1059–1071, 2013.
- [16] —, "Constant-Envelope Multi-User Precoding for Frequency-Selective Massive MIMO Systems," *IEEE Wireless Communications Letters*, vol. 2, no. 5, pp. 547–550, 2013.
- [17] C. Mollen and E. G. Larsson, "Multiuser MIMO Precoding with Per-Antenna Continuous-Time Constant-Envelope Constraints," in *2015 IEEE 16th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*. IEEE, 2015, pp. 261–265.
- [18] J.-C. Chen, "Low-Complexity Constant Envelope Precoding Using Finite Resolution Phase Shifters for Multiuser MIMO Systems With Large Antenna Arrays," *IEEE Transactions on Vehicular Technology*, p. 1, 2018.
- [19] J. Pan and W.-K. Ma, "Constant Envelope Precoding for Single-User Large-Scale MISO Channels: Efficient Precoding and Optimal Designs," *IEEE Journal of Selected Topics in Signal Processing*, vol. 8, no. 5, pp. 982–995, 2014.
- [20] S. Zhang, R. Zhang, and T. J. Lim, "MISO Multicasting With Constant Envelope Precoding," *IEEE Wireless Communications Letters*, vol. 5, no. 6, pp. 588–591, 2016.
- [21] F. Liu, C. Masouros, P. V. Amadori, and H. Sun, "An Efficient Manifold Algorithm for Constructive Interference Based Constant Envelope Precoding," *IEEE Signal Processing Letters*, vol. 24, no. 10, pp. 1542–1546, 2017.
- [22] P. V. Amadori and C. Masouros, "Constant Envelope Precoding by Interference Exploitation in Phase Shift Keying-Modulated Multiuser Transmission," *IEEE Transactions on Wireless Communications*, vol. 16, no. 1, pp. 538–550, 2017.
- [23] H. Shen, W. Xu, A. Lee Swindlehurst, and C. Zhao, "Transmitter Optimization for Per-Antenna Power Constrained Multi-Antenna Downlinks: An SLNR Maximization Methodology," *IEEE Transactions on Signal Processing*, vol. 64, no. 10, pp. 2712–2725, 2016.
- [24] S. Zhang, R. Zhang, and T. J. Lim, "Constant Envelope Precoding for MIMO Systems," *IEEE Transactions on Communications*, vol. 66, no. 1, pp. 149–162, 2018.
- [25] A. Mezghani, R. Ghiat, and J. A. Nossek, "Transmit Processing with Low Resolution D/A-Converters," in *Proc. 16th IEEE International Conference on Electronics, Circuits and Systems - (ICECS)*, 2009, pp. 683–686.
- [26] O. B. Usman, H. Jedda, A. Mezghani, and J. A. Nossek, "MMSE Precoder for Massive MIMO Using 1-Bit Quantization," in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2016, pp. 3381–3385.
- [27] R. Price, "A Useful Theorem for Nonlinear Devices Having Gaussian Inputs," *IEEE Transactions on Information Theory*, vol. 4, no. 2, pp. 69–72, 1958.
- [28] H. Jedda, J. A. Nossek, and A. Mezghani, "Minimum BER Precoding in 1-Bit Massive MIMO Systems," in *Proc. IEEE Sensor Array and Multichannel Signal Processing Workshop (SAM)*, 2016.
- [29] S. Jacobsson, G. Durisi, M. Coldrey, T. Goldstein, and C. Studer, "Nonlinear 1-Bit Precoding for Massive MU-MIMO with Higher-Order Modulation," in *Proc. 50th Asilomar Conference on Signals, Systems and Computers*, 2016, pp. 763–767.

-
- [30] —, “Quantized Precoding for Massive MU-MIMO,” *IEEE Transactions on Communications*, vol. 65, no. 11, pp. 4670–4684, 2017.
- [31] H. Jedda, A. Mezghani, J. A. Nossek, and A. L. Swindlehurst, “Massive MIMO Downlink 1-Bit Precoding with Linear Programming for PSK Signaling,” in *Proc. 18th IEEE International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, 2017.
- [32] O. Castañeda, T. Goldstein, and C. Studer, “POKEMON: A Non-Linear Beamforming Algorithm for 1-Bit Massive MIMO,” in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 3464–3468.
- [33] A. Swindlehurst, A. Saxena, A. Mezghani, and I. Fijalkow, “Minimum Probability-of-Error Perturbation Precoding for the One-Bit Massive MIMO Downlink,” in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 6483–6487.
- [34] M. Shao, Q. Li, and W.-K. Ma, “One-Bit Massive MIMO Precoding via a Minimum Symbol-Error Probability Design,” *arXiv:1803.04787*, Mar. 2018.
- [35] O. Castañeda, S. Jacobsson, G. Durisi, M. Coldrey, T. Goldstein, and C. Studer, “1-bit Massive MU-MIMO Precoding in VLSI,” *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, vol. 7, no. 4, pp. 508–522, 2017.
- [36] H. Jedda, A. Mezghani, J. A. Nossek, and A. L. Swindlehurst, “Massive MIMO Downlink 1-Bit Precoding for Frequency Selective Channels,” in *2017 IEEE 7th International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*. IEEE, 2017.
- [37] L. T. N. Landau and R. C. de Lamare, “Branch-and-Bound Precoding for Multiuser MIMO Systems with 1-Bit Quantization,” *IEEE Wireless Communications Letters*, 2017.
- [38] A. Noll, H. Jedda, and J. A. Nossek, “PSK Precoding in Multi-User MISO Systems,” in *Proc. 21st International ITG Workshop on Smart Antennas (WSA)*. Berlin: VDE Verlag GMBH, 2017.
- [39] S. Jacobsson, O. Castaneda, C. Jeon, G. Durisi, and C. Studer, “Nonlinear Precoding for Phase-Quantized Constant-Envelope Massive MU-MIMO-OFDM,” in *2018 25th International Conference on Telecommunications (ICT)*. IEEE, 2018, pp. 367–372.
- [40] A. Nedelcu, F. Steiner, M. Staudacher, G. Kramer, W. Zirwas, R. S. Ganesan, P. Baracca, and S. Wesemann, “Quantized Precoding for Multi-Antenna Downlink Channels with MAGIQ,” in *Proc. 22nd International ITG Workshop on Smart Antennas (WSA)*. Berlin: VDE Verlag GMBH, 2018.
- [41] M. Kazemi, H. Aghaeinia, and T. M. Duman, “Discrete-Phase Constant Envelope Precoding for Massive MIMO Systems,” *IEEE Transactions on Communications*, vol. 65, no. 5, pp. 2011–2021, 2017.
- [42] C. Masouros, T. Ratnarajah, M. Sellathurai, C. B. Papadias, and A. K. Shukla, “Known Interference in the Cellular Downlink: A Performance Limiting Factor or a Source of Green Signal Power?” *IEEE Communications Magazine*, vol. 51, no. 10, pp. 162–171, 2013.
- [43] C. Masouros and G. Zheng, “Exploiting Known Interference as Green Signal Power for Downlink Beamforming Optimization,” *IEEE Transactions on Signal Processing*, vol. 63, no. 14, pp. 3628–3640, 2015.

-
- [44] H. Jedda, A. Mezghani, A. L. Swindlehurst, and J. A. Nossek, "Quantized Constant Envelope Precoding with PSK and QAM Signaling," *IEEE Transactions on Wireless Communications*, vol. 17, no. 12, pp. 8022–8034, © 2018 IEEE.
- [45] H. Jedda and J. A. Nossek, "Quantized Constant Envelope Precoding for Frequency Selective Channels," in *2018 IEEE Statistical Signal Processing Workshop (SSP)*. IEEE, 2018, pp. 213–217.
- [46] H. Jedda, M. M. Ayub, J. Munir, A. Mezghani, and J. A. Nossek, "Power- and Spectral Efficient Communication System Design Using 1-Bit Quantization," in *Proc. International Symposium on Wireless Communication Systems (ISWCS)*, 2015, pp. 296–300.
- [47] H. Jedda, A. Mezghani, and J. A. Nossek, "Spectral Shaping with Low Resolution Signals," in *Proc. 49th Asilomar Conference on Signals, Systems and Computers*, 2015, pp. 1437–1441.
- [48] C. Mollen, E. G. Larsson, and T. Eriksson, "Waveforms for the Massive MIMO Downlink: Amplifier Efficiency, Distortion, and Performance," *IEEE Transactions on Communications*, vol. 64, no. 12, pp. 5050–5063, 2016.
- [49] H. Jedda and J. A. Nossek, "On the Statistical Properties of Constant Envelope Quantizers," *IEEE Wireless Communications Letters*, vol. 7, no. 6, pp. 1006–1009, © 2018 IEEE.
- [50] S. C. Cripps, *RF Power Amplifiers for Wireless Communications*, 2nd ed., ser. Artech House microwave library. Boston, Mass. and London: Artech House, 2006.
- [51] B. Lehmeier, A. Mezghani, and J. A. Nossek, "Electronic Amplifier for Amplifying an Input Signal," International Patent WO 2018/024 756 A1, 2018.
- [52] B. Razavi, *Principles of Data Conversion System Design*. New York: IEEE Press, 1995.
- [53] J. Groe, "Polar Transmitters for Wireless Communications," *IEEE Communications Magazine*, vol. 45, no. 9, pp. 58–63, 2007.
- [54] M. Abolfadl Ibrahim, "The Polar Transmitter: Analysis and Algorithms," Doctoral Thesis, Universität Stuttgart, 2015. [Online]. Available: <http://elib.uni-stuttgart.de/handle/11682/3647>
- [55] H. H. Andersen, M. Højbjerg, D. Sørensen, and P. S. Eriksen, "The Multivariate Complex Normal Distribution," in *Linear and Graphical Models*, ser. Lecture Notes in Statistics, P. Diggle, S. Fienberg, K. Krickeberg, I. Olkin, N. Wermuth, H. H. Andersen, M. Højbjerg, D. Sørensen, and P. S. Eriksen, Eds. New York, NY: Springer New York, 1995, vol. 101, pp. 15–37.
- [56] E. McMahon, "An Extension of Price's theorem (Corresp.)," *IEEE Transactions on Information Theory*, vol. 10, no. 2, p. 168, 1964.
- [57] A. Papoulis, "Comments on 'An Extension of Price's Theorem' by McMahon, E. L.," *IEEE Transactions on Information Theory*, vol. 11, no. 1, p. 154, 1965.
- [58] R. Pawula, "A Modified Version of Price's Theorem," *IEEE Transactions on Information Theory*, vol. 13, no. 2, pp. 285–288, 1967.
- [59] W. McGee, "Circularly Complex Gaussian Noise—A price Theorem and a Mehler Expansion (Corresp.)," *IEEE Transactions on Information Theory*, vol. 15, no. 2, pp. 317–319, 1969.
- [60] A. van den Bos, "Price's Theorem for Complex Variates," *IEEE Transactions on Information Theory*, vol. 42, no. 1, pp. 286–287, 1996.

-
- [61] J. J. Busgang, "Crosscorrelation Functions of Amplitude-Distorted Gaussian Signals." [Online]. Available: <http://hdl.handle.net/1721.1/4847>
- [62] A. Saleh and R. Valenzuela, "A Statistical Model for Indoor Multipath Propagation," *IEEE Journal on Selected Areas in Communications*, vol. 5, no. 2, pp. 128–137, 1987.
- [63] L. M. Correia and P. Smulders, "Characterisation of Propagation in 60 GHz Radio Channels," *Electronics & Communication Engineering Journal*, vol. 9, no. 2, pp. 73–80, 1997.
- [64] A. M. Sayeed, "Deconstructing Multiantenna Fading Channels," *IEEE Transactions on Signal Processing*, vol. 50, no. 10, pp. 2563–2579, 2002.
- [65] H. Xu, V. Kukshya, and T. S. Rappaport, "Spatial and Temporal Characteristics of 60-GHz Indoor Channels," *IEEE Journal on Selected Areas in Communications*, vol. 20, no. 3, pp. 620–630, 2002.
- [66] O. E. Ayach, S. Rajagopal, S. Abu-Surra, Z. Pi, and R. W. Heath, "Spatially Sparse Precoding in Millimeter Wave MIMO Systems," *IEEE Transactions on Wireless Communications*, vol. 13, no. 3, pp. 1499–1513, 2014.
- [67] A. Forenza, D. J. Love, and R. W. Heath, "Simplified Spatial Correlation Models for Clustered MIMO Channels With Different Array Configurations," *IEEE Transactions on Vehicular Technology*, vol. 56, no. 4, pp. 1924–1934, 2007.
- [68] W. Feller, *An Introduction to Probability Theory and Its Applications*, 3rd ed., ser. Wiley series in probability and mathematical statistics. New York: Wiley, 1970.
- [69] A. Kakkavas, "Precoding Techniques for Coarsely Quantized MIMO Systems," Master's thesis, Technical University of Munich, 2015.
- [70] M. Joham, W. Utschick, and J. A. Nossek, "Linear Transmit Processing in MIMO Communications Systems," *IEEE Transactions on Signal Processing*, vol. 53, no. 8, pp. 2700–2712, 2005.
- [71] S. P. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge: Cambridge University Press, 2004.
- [72] G. B. Dantzig and M. N. Thapa, *Linear Programming*, ser. Springer series in operations research. New York and London: Springer, 1997.
- [73] M. Joham, *Optimization of Linear and Nonlinear Transmit Signal Processing*. Aachen: Shaker, 2004.
- [74] T. Aulin and C. Sundberg, "Continuous Phase Modulation—Part I: Full Response Signaling," *IEEE Transactions on Communications*, vol. 29, no. 3, pp. 196–209, 1981.
- [75] K. Murota and K. Hirade, "GMSK Modulation for Digital Mobile Radio Telephony," *IEEE Transactions on Communications*, vol. 29, no. 7, pp. 1044–1050, 1981.
- [76] D. B. Owen, "A table of Normal Integrals," *Communications in Statistics - Simulation and Computation*, vol. 9, no. 4, pp. 389–419, 2007.