

Minimum Union Bound Symbol Error Probability Precoding for PSK Modulation and Phase Quantization

Erico S. P. Lopes and Lukas T. N. Landau
Centre for Telecommunications Studies
Pontifical Catholic University of Rio de Janeiro
Rio de Janeiro, Brazil 22453-900
erico;lukas.landau@cetuc.puc-rio.br

Amine Mezghani
Department of Electrical and Computer Engineering
University of Manitoba
Winnipeg, Canada
amine.mezghani@umanitoba.ca

Abstract—This study formulates a novel precoding criterion for multiuser MIMO systems based on the minimization of the symbol error probability at the users. Unlike previous formulations that require the symbols to belong to a QPSK modulation the proposed criterion allows the utilization of PSK modulation in general. Based on the proposed minimum symbol error probability criterion a discrete programming problem is derived. Using a sophisticated branch-and-bound method the proposed precoding problem is optimally solved. Numerical results show that the proposed precoding method outperforms the state-of-the-art techniques for all examined SNR values in terms of symbol-error-rate.

Index Terms—Energy Efficiency, Precoding, Constant Envelope, Low-Resolution Quantization, MIMO systems, Symbol Error Probability.

I. INTRODUCTION

Multiuser multiple-input multiple-output (MU-MIMO) systems are considered as a promising physical-layer technique and are expected to be vital for the future of wireless communications networks [1]. However, due to the high number of radio frequency front ends (RFFE) the energy consumption of the radio frequency chains imposes a challenge for this kind of technology [2].

This energy efficiency (EE) challenge led to the development of different studies that analyzed the circuit of RFFEs to dissect which are the most consuming elements, e.g., [3], [4]. These works conclude that the power amplifiers (PAs) and data converters are two of the most consuming elements in the RFFE. With this, many of the recent studies consider adopting features to minimize the power consumption of these elements. In most cases, to increase the PA's efficiency the adoption of constant envelope (CE) signaling is considered, and, to decrease the power consumption of the data converters, low-resolution in amplitude is utilized. The main drawback of adopting these features is the error-rate performance degradation they yield.

To mitigate the performance degradation CE low-resolution precoding has received increasing attention from the wireless

communications community. Linear approaches, e.g., [5]–[7], have benefit from a relatively low computational complexity. However, to achieve higher reliability nonlinear symbol-level-precoding methods have been presented based on different design criteria. Two of the most popular criteria used for precoding are the minimum mean squared error (MMSE) used in [8]–[14] and the maximum minimum distance to the decision threshold (MMDDT) which is utilized in [11], [15]–[21]. Aside from MMSE and MMDDT, the work from [22] proposed the direct optimization of the symbol error probability (SEP) for the 1-bit quantization case.

Following the direction from the work in [22] the present study considers the minimization of the SEP as the design criterion and focuses on the development of precoding techniques for a MU-MIMO downlink system with PSK modulation and phase quantization.

The novel formulation is developed based on the minimization of the union bound SEP (MUBSEP) which then allows for the utilization of PSK modulation in general. Based on this formulation an optimization problem is devised which is optimally solved using a sophisticated branch-and-bound (B&B) method.

Numerical results confirm that the proposed B&B algorithm based on the MUBSEP formulation outperforms other state-of-the-art methods in terms of symbol-error-rate (SER) for medium and high SNR values.

The remainder of this paper is organized as follows: Section II describes the system model, whereas Section III exposes the derivation of the formulation utilized for precoding. Section IV presents the design of the proposed B&B precoder. Section V presents and discusses numerical results, while Section VI gives the conclusions.

Regarding the notation, bold lower case and upper case letters indicate vectors and matrices, respectively. Non-bold letters express scalars. The operators $(\cdot)^*$ and $(\cdot)^T$ denote complex conjugation and transposition, respectively. Real and imaginary part operator, as well as the functions $\operatorname{erfc}(\cdot)$, $\operatorname{erf}(\cdot)$ and $\log(\cdot)$, are also applied to vectors and matrices, e.g., $\operatorname{Re}\{\mathbf{x}\} = [\operatorname{Re}\{\mathbf{x}_1\}, \dots, \operatorname{Re}\{\mathbf{x}_M\}]^T$. The operator

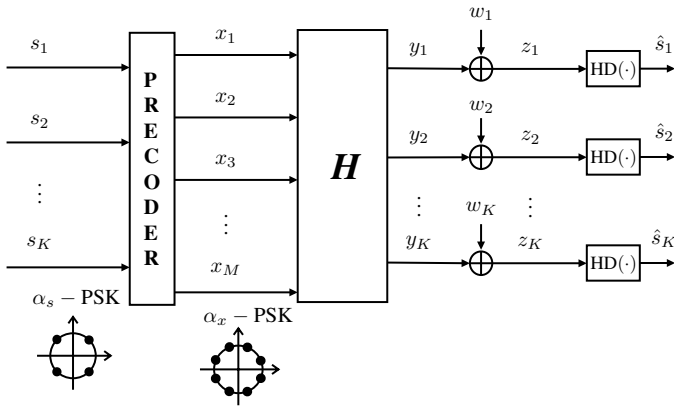


Fig. 1: Multiuser MIMO downlink with discrete precoding

$R(\cdot)$ converts a complex-valued vector into the equivalent real-valued notation. For a given column vector $\mathbf{a} \in \mathbb{C}^M$ the equivalent real-valued vector $\mathbf{a}_r = R(\mathbf{a})$ reads as

$$\mathbf{a}_r = [\text{Re}\{\mathbf{a}_1\} \text{Im}\{\mathbf{a}_1\} \cdots \text{Re}\{\mathbf{a}_M\} \text{Im}\{\mathbf{a}_M\}]^T. \quad (1)$$

The inverse operation is denoted as $C(\cdot)$ which converts equivalent real-valued notation into complex-valued. Finally, for the given vectors \mathbf{a} and \mathbf{b} , $P(\mathbf{a} = \mathbf{b})$ denotes the probability of the event $\mathbf{a} = \mathbf{b}$.

II. SYSTEM MODEL

The system model, illustrated in Fig. 1, consists of a single-cell MU-MIMO scenario where the BS has perfect channel state information (CSI) and is equipped with M transmit antennas serving K single antenna users.

In this study, a symbol level transmission is considered where s_k represents the data symbol for the k -th user. Each symbol s_k is considered to belong to the set \mathcal{S} that represents all possible symbols of a α_s -PSK modulation and is given by

$$\mathcal{S} = \left\{ s : s = e^{\frac{j\pi(2i+1)}{\alpha_s}}, \text{ for } i = 1, \dots, \alpha_s \right\}. \quad (2)$$

The symbols of all users are described in a stacked vector notation as $\mathbf{s} = [s_1, \dots, s_K]^T \in \mathcal{S}^K$. It is considered that different users' symbols are independent and that $P(s_k = s_i) = 1/\alpha_s, \forall i \in \{1, \dots, \alpha_s\}$. Based on \mathbf{s} the precoder computes the transmit vector $\mathbf{x} = [x_1, \dots, x_M]^T$. The entries of \mathbf{x} are constrained to the set \mathcal{X} which is given by

$$\mathcal{X} = \left\{ \mathbf{x} : \mathbf{x} = \sqrt{\frac{P_{\text{tx}}}{M}} e^{\frac{j\pi(2i+1)}{\alpha_x}}, \text{ for } i = 1, \dots, \alpha_x \right\}. \quad (3)$$

The vector \mathbf{x} is transmitted over a frequency flat fading channel described by the matrix $\mathbf{H} \in \mathbb{C}^{K \times M}$ such that the received signal corresponding to the k -th user reads as

$$z_k = y_k + w_k = \mathbf{h}_k \mathbf{x} + w_k, \quad (4)$$

where y_k is the noiseless received signal from the k -th user, \mathbf{h}_k is the k -th row of the channel matrix \mathbf{H} and the complex random variable $w_k \sim \mathcal{CN}(0, \sigma_w^2)$ represents additive white

Gaussian noise (AWGN). Using stacked vector notation (4) can be extended to

$$\mathbf{z} = \mathbf{y} + \mathbf{w} = \mathbf{H} \mathbf{x} + \mathbf{w}, \quad (5)$$

where $\mathbf{z} = [z_1 \dots z_K]^T$, $\mathbf{y} = [y_1 \dots y_K]^T$ and $\mathbf{w} = [w_1 \dots w_K]^T$. Each received symbol z_k is, then, hard detected based on which decision region it belongs to, meaning that z_k is detected as s_i if $z_k \in \mathcal{S}_i$. In the case of PSK modulation, the decision regions are circle sectors with infinite radius and angle of 2θ , where θ is given by $\theta = \pi/\alpha_s$. With this, the estimated symbol from the k -th user is given by $\hat{s}_k = \text{HD}(z_k)$, where $\text{HD}(\cdot)$ represents the hard detection operation. Finally, the detected symbol vector is written as $\hat{\mathbf{s}} = [\hat{s}_1, \dots, \hat{s}_K]$.

III. MUBSEP PRECODING FORMULATION

The probability of detecting the data vector \mathbf{s} conditioned on the transmit vector \mathbf{x} can be computed based on the probabilities of detection of the individual users as

$$P(\hat{\mathbf{s}} = \mathbf{s} | \mathbf{x}) = \prod_{k=1}^K P(\hat{s}_k = s_k | \mathbf{x}). \quad (6)$$

To simplify the notation we denote $P(\hat{\mathbf{s}} = \mathbf{s} | \mathbf{x})$ as $P(\hat{\mathbf{s}} | \mathbf{x})$ and $P(\hat{s}_k = s_k | \mathbf{x})$ as $P(\hat{s}_k | \mathbf{x})$. With this, (6) is rewritten as

$$P(\hat{\mathbf{s}} | \mathbf{x}) = \prod_{k=1}^K P(\hat{s}_k | \mathbf{x}). \quad (7)$$

As stated before, the detector decides for s_k when the received symbol z_k belongs to \mathcal{S}_k . Thus, the individual user probabilities are given by

$$P(\hat{s}_k | \mathbf{x}) = P(z_k \in \mathcal{S}_k | \mathbf{x}) = \frac{1}{\pi \sigma_w^2} \int_{\mathcal{S}_k} e^{-\frac{|t - y_k|^2}{\sigma_w^2}} dt. \quad (8)$$

The integral from (8) has tabled solutions for $\alpha_s \in \{2, 4\}$. Yet, for $\alpha_s \notin \{2, 4\}$, (8) requires the utilization of Monte Carlo methods, which are not suitable for symbol-level-precoding algorithms due to their relatively high computation complexity. Thus, to achieve a general design, the maximization of a lower bound on the probability of correct detection is considered.

The union bound states that for any finite or countable set of events, the probability that at least one of the events happens is smaller or equal than the sum of the probabilities of the individual events [23], meaning

$$P\left(\bigcup_i A_i\right) \leq \sum_i P(A_i). \quad (9)$$

with A_i representing an event. With this, the error probability for the k -th user, $P_e(\hat{s}_k | \mathbf{x})$, can be bounded by

$$\begin{aligned} P_e(\hat{s}_k | \mathbf{x}) &= P(z_k \in \mathcal{Z}_1 \cup \mathcal{Z}_2 | \mathbf{x}) \\ &\leq P(z_k \in \mathcal{Z}_1 | \mathbf{x}) + P(z_k \in \mathcal{Z}_2 | \mathbf{x}), \end{aligned} \quad (10)$$

where the sets \mathcal{Z}_1 and \mathcal{Z}_2 are depicted in Fig. 2. The individual

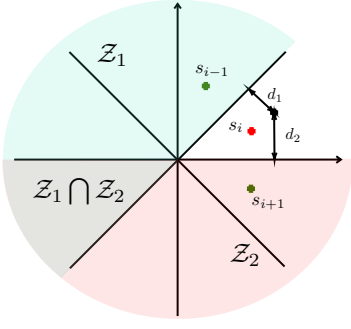


Fig. 2: Representation of the union bound

probabilities can be computed based on the minimum distances to the decision thresholds (MDDTs), $d_{1,k}$ and $d_{2,k}$, as

$$\begin{aligned} \mathbb{P}(z_k \in \mathcal{Z}_1 | \mathbf{x}) &= \int_{d_{1,k}}^{\infty} \frac{1}{\sqrt{\pi\sigma_w^2}} e^{-\frac{t^2}{\sigma_w^2}} dt = \frac{1}{2} \operatorname{erfc} \left(\frac{d_{1,k}}{\sigma_w} \right) \\ \mathbb{P}(z_k \in \mathcal{Z}_2 | \mathbf{x}) &= \int_{d_{2,k}}^{\infty} \frac{1}{\sqrt{\pi\sigma_w^2}} e^{-\frac{t^2}{\sigma_w^2}} dt = \frac{1}{2} \operatorname{erfc} \left(\frac{d_{2,k}}{\sigma_w} \right). \end{aligned}$$

The MDDTs are computed, similarly to in [16] and [20], by applying a rotation of $\arg\{s_k^*\} = -\phi_{s_k}$ to the coordinate system such that the symbol of interest is placed on the real axis. This is done by multiplying both s_k and y_k by $e^{-j\phi_{s_k}}$ which results in $e^{-j\phi_{s_k}} s_k = 1$ and $\omega_k = e^{-j\phi_{s_k}} y_k$. Based on the rotated coordinate system the MDDTs are computed as

$$d_{1,k}(\mathbf{x}) = \operatorname{Re}\{s_k^* \mathbf{h}_k \mathbf{x}\} \sin \theta - \operatorname{Im}\{s_k^* \mathbf{h}_k \mathbf{x}\} \cos \theta \quad (11)$$

$$d_{2,k}(\mathbf{x}) = \operatorname{Re}\{s_k^* \mathbf{h}_k \mathbf{x}\} \sin \theta + \operatorname{Im}\{s_k^* \mathbf{h}_k \mathbf{x}\} \cos \theta. \quad (12)$$

Using (10) to (12) one can construct a bound on the probability of correct detection of the k -th user as

$$\begin{aligned} \mathbb{P}(\hat{s}_k | \mathbf{x}) &= 1 - \mathbb{P}_e(\hat{s}_k | \mathbf{x}) \\ &\geq 1 - (\mathbb{P}(z_k \in \mathcal{Z}_1 | \mathbf{x}) + \mathbb{P}(z_k \in \mathcal{Z}_2 | \mathbf{x})) \\ &= 1 - \frac{1}{2} \operatorname{erfc} \left(\frac{d_{1,k}(\mathbf{x})}{\sigma_w} \right) - \frac{1}{2} \operatorname{erfc} \left(\frac{d_{2,k}(\mathbf{x})}{\sigma_w} \right) \\ &= \frac{1}{2} \operatorname{erf} \left(\frac{d_{1,k}(\mathbf{x})}{\sigma_w} \right) + \frac{1}{2} \operatorname{erf} \left(\frac{d_{2,k}(\mathbf{x})}{\sigma_w} \right). \quad (13) \end{aligned}$$

The probability of correct detection can then be bounded by the union bound probability, meaning $\mathbb{P}(\hat{s}_k | \mathbf{x}) \geq \mathbb{P}_{\text{ub}}(\hat{s}_k | \mathbf{x})$, with

$$\mathbb{P}_{\text{ub}}(\hat{s}_k | \mathbf{x}) = \frac{1}{2^K} \prod_{k=1}^K \left(\operatorname{erf} \left(\frac{d_{1,k}(\mathbf{x})}{\sigma_w} \right) + \operatorname{erf} \left(\frac{d_{2,k}(\mathbf{x})}{\sigma_w} \right) \right).$$

Based on $\mathbb{P}_{\text{ub}}(\hat{s}_k | \mathbf{x})$ an optimization problem can be cast as

$$\max_{\mathbf{x} \in \mathcal{X}^M} \prod_{k=1}^K \left(\operatorname{erf} \left(\frac{d_{1,k}(\mathbf{x})}{\sigma_w} \right) + \operatorname{erf} \left(\frac{d_{2,k}(\mathbf{x})}{\sigma_w} \right) \right). \quad (14)$$

Since $\log(\cdot)$ is a monotonically increasing function, applying it to the objective from (14) yields an equivalent problem. With this, the proposed MUBSEP optimization problem for an α_s -PSK modulation reads as

$$\min_{\mathbf{x} \in \mathcal{X}^M} - \sum_{k=1}^K \log \left(\operatorname{erf} \left(\frac{d_{1,k}(\mathbf{x})}{\sigma_w} \right) + \operatorname{erf} \left(\frac{d_{2,k}(\mathbf{x})}{\sigma_w} \right) \right). \quad (15)$$

An equivalent real-valued formulation of (15) can be cast as

$$\begin{aligned} \min_{\mathbf{x}_r} & - \sum_{k=1}^K \log \left(\operatorname{erf}(\mathbf{u}_{1,k}^T \mathbf{x}_r) + \operatorname{erf}(\mathbf{u}_{2,k}^T \mathbf{x}_r) \right) \quad (16) \\ \text{s.t.} & \quad \mathbf{x}_{r,2m-1} + j\mathbf{x}_{r,2m} \in \mathcal{X} \quad \text{for } m = 1, \dots, M. \end{aligned}$$

where $\mathbf{x}_r = R(\mathbf{x})$, $\mathbf{u}_{1,k} = (\mathbf{h}_{R,\theta,k}^{s*} - \mathbf{h}_{1,\theta,k}^{s*})^T$ and $\mathbf{u}_{2,k} = (\mathbf{h}_{R,\theta,k}^{s*} + \mathbf{h}_{1,\theta,k}^{s*})^T$ with $\mathbf{h}_{R,\theta,k}^{s*}$ and $\mathbf{h}_{1,\theta,k}^{s*}$ being the k -th rows of matrices $\mathbf{H}_{R,\theta}^{s*}$ and $\mathbf{H}_{1,\theta}^{s*}$. The matrices $\mathbf{H}_{R,\theta}^{s*}$ and $\mathbf{H}_{1,\theta}^{s*}$ are given by $\mathbf{H}_{R,\theta}^{s*} = \frac{\sin(\theta)}{\sigma_w} \mathbf{H}_R^{s*}$, and $\mathbf{H}_{1,\theta}^{s*} = \frac{\cos(\theta)}{\sigma_w} \mathbf{H}_1^{s*}$, with

$$\begin{aligned} \mathbf{H}_R^{s*} &= \begin{bmatrix} \operatorname{Re}\{h_{11}^{s*}\} & -\operatorname{Im}\{h_{11}^{s*}\} & \cdots & \operatorname{Re}\{h_{1M}^{s*}\} & -\operatorname{Im}\{h_{1M}^{s*}\} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \operatorname{Re}\{h_{K1}^{s*}\} & -\operatorname{Im}\{h_{K1}^{s*}\} & \cdots & \operatorname{Re}\{h_{KM}^{s*}\} & -\operatorname{Im}\{h_{KM}^{s*}\} \end{bmatrix} \\ \mathbf{H}_1^{s*} &= \begin{bmatrix} \operatorname{Im}\{h_{11}^{s*}\} & \operatorname{Re}\{h_{11}^{s*}\} & \cdots & \operatorname{Im}\{h_{1M}^{s*}\} & \operatorname{Re}\{h_{1M}^{s*}\} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \operatorname{Im}\{h_{K1}^{s*}\} & \operatorname{Re}\{h_{K1}^{s*}\} & \cdots & \operatorname{Re}\{h_{KM}^{s*}\} & \operatorname{Im}\{h_{KM}^{s*}\} \end{bmatrix}, \end{aligned}$$

where h_{ij}^{s*} is the element of the i -th row and j -th column of the matrix $\mathbf{H}^{s*} = \operatorname{diag}\{s^*\} \mathbf{H}$.

A. Conditions for convexity of the MUBSEP objective

Considering the real-valued formulation described in (16) the MUBSEP objective can be cast as

$$g(\mathbf{x}_r) = - \sum_{k=1}^K \log \left(\operatorname{erf}(\mathbf{u}_{1,k}^T \mathbf{x}_r) + \operatorname{erf}(\mathbf{u}_{2,k}^T \mathbf{x}_r) \right). \quad (17)$$

Convexity can be proven by evaluating the conditions under which the Hessian is positive semi-definite (PSD) [24]. To this end, the Hessian is calculated in what follows. Taking the derivative of $g(\mathbf{x}_r)$ with respect to \mathbf{x}_r yields

$$\begin{aligned} \frac{\partial g(\mathbf{x}_r)}{\partial \mathbf{x}_r} &= - \sum_{k=1}^K \frac{\partial}{\partial \mathbf{x}_r} \log \left(\operatorname{erf}(\mathbf{u}_{1,k}^T \mathbf{x}_r) + \operatorname{erf}(\mathbf{u}_{2,k}^T \mathbf{x}_r) \right) \\ &= - \sum_{k=1}^K \frac{\frac{\partial}{\partial \mathbf{x}_r} \left(\operatorname{erf}(\mathbf{u}_{1,k}^T \mathbf{x}_r) + \operatorname{erf}(\mathbf{u}_{2,k}^T \mathbf{x}_r) \right)}{\operatorname{erf}(\mathbf{u}_{1,k}^T \mathbf{x}_r) + \operatorname{erf}(\mathbf{u}_{2,k}^T \mathbf{x}_r)} \\ &= - \sum_{k=1}^K \frac{\frac{2}{\sqrt{\pi}} e^{-(\mathbf{u}_{1,k}^T \mathbf{x}_r)^2} \mathbf{u}_{1,k} + \frac{2}{\sqrt{\pi}} e^{-(\mathbf{u}_{2,k}^T \mathbf{x}_r)^2} \mathbf{u}_{2,k}}{\operatorname{erf}(\mathbf{u}_{1,k}^T \mathbf{x}_r) + \operatorname{erf}(\mathbf{u}_{2,k}^T \mathbf{x}_r)}. \end{aligned}$$

Note that $\frac{\partial g(\mathbf{x}_r)}{\partial \mathbf{x}_r}$ can be written in the form

$$\frac{\partial g(\mathbf{x}_r)}{\partial \mathbf{x}_r} = - \sum_{k=1}^K \frac{\mathbf{m}_k(\mathbf{x}_r)}{q_k(\mathbf{x}_r)}$$

where

$$\begin{aligned} \mathbf{m}_k(\mathbf{x}_r) &= \frac{2}{\sqrt{\pi}} e^{-(\mathbf{u}_{1,k}^T \mathbf{x}_r)^2} \mathbf{u}_{1,k} + \frac{2}{\sqrt{\pi}} e^{-(\mathbf{u}_{2,k}^T \mathbf{x}_r)^2} \mathbf{u}_{2,k} \\ q_k(\mathbf{x}_r) &= \operatorname{erf}(\mathbf{u}_{1,k}^T \mathbf{x}_r) + \operatorname{erf}(\mathbf{u}_{2,k}^T \mathbf{x}_r). \end{aligned}$$

The Hessian is, then, given by

$$\frac{\partial^2 g(\mathbf{x}_r)}{\partial \mathbf{x}_r \partial \mathbf{x}_r^T} = - \sum_{k=1}^K \frac{\frac{\partial \mathbf{m}_k(\mathbf{x}_r)}{\partial \mathbf{x}_r^T} q_k(\mathbf{x}_r) - \mathbf{m}_k(\mathbf{x}_r) \frac{\partial q_k(\mathbf{x}_r)}{\partial \mathbf{x}_r^T}}{(q_k(\mathbf{x}_r))^2}.$$

The terms $\frac{\partial \mathbf{m}_k(\mathbf{x}_r)}{\partial \mathbf{x}_r^T}$ and $\frac{\partial q_k(\mathbf{x}_r)}{\partial \mathbf{x}_r^T}$ are calculated as

$$\begin{aligned} \frac{\partial \mathbf{m}_k(\mathbf{x}_r)}{\partial \mathbf{x}_r^T} &= -(\Psi_{1,k} + \Psi_{2,k}) \\ \frac{\partial q_k(\mathbf{x}_r)}{\partial \mathbf{x}_r^T} &= \mathbf{m}_k^T(\mathbf{x}_r), \end{aligned}$$

where $\Psi_{1,k}$ and $\Psi_{2,k}$ read as

$$\begin{aligned} \Psi_{1,k} &= \frac{4}{\sqrt{\pi}} e^{-(\mathbf{u}_{1,k}^T \mathbf{x}_r)^2} \mathbf{u}_{1,k} \mathbf{u}_{1,k}^T \mathbf{x}_r \mathbf{u}_{1,k}^T, \\ \Psi_{2,k} &= \frac{4}{\sqrt{\pi}} e^{-(\mathbf{u}_{2,k}^T \mathbf{x}_r)^2} \mathbf{u}_{2,k} \mathbf{u}_{2,k}^T \mathbf{x}_r \mathbf{u}_{2,k}^T. \end{aligned}$$

The Hessian then reads as

$$\frac{\partial^2 g(\mathbf{x}_r)}{\partial \mathbf{x}_r \partial \mathbf{x}_r^T} = \sum_{k=1}^K \frac{(\Psi_{1,k} + \Psi_{2,k}) q_k(\mathbf{x}_r) + \mathbf{m}_k(\mathbf{x}_r) \mathbf{m}_k^T(\mathbf{x}_r)}{(q_k(\mathbf{x}_r))^2}.$$

A sufficient condition for PSD is $(\Psi_{1,k} + \Psi_{2,k}) q_k(\mathbf{x}_r) \succeq \mathbf{0} \forall k \in \{1, \dots, K\}$. With this, positive semi-definiteness is guaranteed for $\mathbf{u}_{1,k}^T \mathbf{x}_r \geq 0 \forall k \in \{1, \dots, K\}$ and $\mathbf{u}_{2,k}^T \mathbf{x}_r \geq 0 \forall k \in \{1, \dots, K\}$. Finally, the condition for convexity of the MUBSEP objective function can be cast in a stacked manner for all k as

$$\begin{bmatrix} \mathbf{H}_{R,\theta}^{s*} - \mathbf{H}_{1,\theta}^{s*} \\ \mathbf{H}_{R,\theta}^{s*} + \mathbf{H}_{1,\theta}^{s*} \end{bmatrix} \mathbf{x}_r \succeq \mathbf{0}. \quad (18)$$

From these results, convexity can be guaranteed by restricting the original feasible set using the condition from (18). With this, a convex MUBSEP optimization problem can be formulated as

$$\begin{aligned} \min_{\mathbf{x}_r} & - \sum_{k=1}^K \log(\operatorname{erf}(\mathbf{u}_{1,k}^T \mathbf{x}_r) + \operatorname{erf}(\mathbf{u}_{2,k}^T \mathbf{x}_r)) \\ \text{s.t.} & \begin{bmatrix} \mathbf{H}_{R,\theta}^{s*} - \mathbf{H}_{1,\theta}^{s*} \\ \mathbf{H}_{R,\theta}^{s*} + \mathbf{H}_{1,\theta}^{s*} \end{bmatrix} \mathbf{x}_r \succeq \mathbf{0}, \\ & x_{r,2m-1} + jx_{r,2m} \in \mathcal{X} \text{ for } m = 1, \dots, M. \end{aligned} \quad (19)$$

Note that, the optimal solution from (15) is not necessarily the same as the optimal solution from (19). However, different solutions are only possible if the constraint described in (18) is active. This would imply that for at least one user, the optimal solution of (15) yields a noiseless received symbol y_i in the incorrect decision region. This leads to a SEP, for this user, greater than half. This is, in general, not a relevant case, since future wireless communications systems will be designed to provide high reliability and avoid this kind of scenario.

IV. OPTIMAL MUBSEP PRECODING DESIGN VIA BRANCH-AND-BOUND

In this section, we propose a B&B algorithm that solves optimally the discrete programming problem (DPP) described in (19). In what follows, the main aspects of the proposed B&B algorithm are exposed.

A. Initialization via convex hull relaxation

The first part of the proposed MUBSEP B&B method is an initialization step where a feasible solution is computed by solving a relaxed MUBSEP problem. To this end, the discrete feasible set \mathcal{X}^M is relaxed to its convex hull \mathcal{P} . Note that \mathcal{P} is a polyhedron and thus can be described as the solution set of a finite number of linear equalities and inequalities [24]. Similarly as done in [13], [14] and [16] the relaxed feasible set is described in real-valued notation using the inequality $\mathbf{R} [\mathbf{x}_r^T, \mathbf{1}]^T \preceq \mathbf{0}$, where $\mathbf{R} = [\mathbf{A}, -\mathbf{b}]$ and

$$\begin{aligned} \mathbf{A} &= [(\mathbf{I}_M \otimes \beta_1)^T, (\mathbf{I}_M \otimes \beta_2)^T, \dots, (\mathbf{I}_M \otimes \beta_{\alpha_x})^T]^T, \\ \beta_i &= \left[\cos\left(\frac{2\pi i}{\alpha_x}\right), -\sin\left(\frac{2\pi i}{\alpha_x}\right) \right], \quad i \in \{1, \dots, \alpha_x\}, \\ \mathbf{b} &= \sqrt{\frac{P_{\text{TX}}}{M}} \cos\left(\frac{\pi}{\alpha_x}\right) \mathbf{1}_{M\alpha_x}. \end{aligned} \quad (20)$$

With this, one can readily write the relaxed MUBSEP problem by substituting \mathcal{X}^M by \mathcal{P} in (19), which then yields

$$\begin{aligned} \min_{\mathbf{x}_r} & - \sum_{k=1}^K \log(\operatorname{erf}(\mathbf{u}_{1,k}^T \mathbf{x}_r) + \operatorname{erf}(\mathbf{u}_{2,k}^T \mathbf{x}_r)) \\ \text{s.t.} & \begin{bmatrix} \mathbf{H}_{R,\theta}^{s*} - \mathbf{H}_{1,\theta}^{s*} \\ \mathbf{H}_{R,\theta}^{s*} + \mathbf{H}_{1,\theta}^{s*} \end{bmatrix} \mathbf{x}_r \succeq \mathbf{0}, \quad \mathbf{R} [\mathbf{x}_r^T, \mathbf{1}]^T \preceq \mathbf{0}. \end{aligned} \quad (21)$$

Replacing \mathcal{X}^M by \mathcal{P} yields a convex problem since (21) is the minimization of a convex objective under a convex set.

The solution of the relaxed problem is termed $\mathbf{x}_{r,\text{lb}}$ which can be written in complex form as $\mathbf{x}_{\text{lb}} = C(\mathbf{x}_{r,\text{lb}})$. Note that, $\mathbf{x}_{\text{lb}} \in \mathcal{P}$ can also belong to the original feasible set \mathcal{X}^M as $\mathcal{P} \cap \mathcal{X}^M \neq \emptyset$. If this is the case, \mathbf{x}_{lb} is also the optimal solution from (19) and no further processing is required. However, if $\mathbf{x}_{\text{lb}} \notin \mathcal{X}^M$ the solution \mathbf{x}_{lb} must be projected to \mathcal{X}^M .

The projection method utilized in this study is uniform quantization (UQ). When using this approach the projected vector is given by $\mathbf{x}_{\text{ub}} = Q(\mathbf{x}_{\text{lb}})$, where $Q(\cdot)$ represents the quantization operation. The quantization criterion utilized is based on the elementwise Euclidean distance. By this approach, the p -th entry of \mathbf{x}_{ub} , denoted as $x_{\text{ub},p}$, is computed as $x_{\text{ub},p} = \arg \min_{i \in \{1, \dots, \alpha_x\}} |x_{\text{lb},p} - x_i|^2$, where $x_{\text{lb},p}$ denotes the p -th entry of \mathbf{x}_{lb} and x_i the i -th element of \mathcal{X} . Note that, the quantized vector \mathbf{x}_{ub} attains the low-resolution constraints, meaning $\mathbf{x}_{\text{ub}} \in \mathcal{X}^M$ and thus could be used for transmission.

As mentioned before if $\mathbf{x}_{\text{lb}} \in \mathcal{X}^M$ the algorithm returns \mathbf{x}_{lb} and no further processing is required. However, if $\mathbf{x}_{\text{lb}} \notin \mathcal{X}^M$, the initial smallest known upper bound \check{g} and its corresponding vector $\check{\mathbf{x}}$ are stored as $\check{g} = g(\mathbf{x}_{\text{ub}})$ and $\check{\mathbf{x}} = \mathbf{x}_{\text{ub}}$, with the objective $g(\mathbf{x})$ given by

$$g(\mathbf{x}) = -\frac{1}{K} \log(\operatorname{erf}(\alpha_1(\mathbf{x})) + \operatorname{erf}(\alpha_2(\mathbf{x}))), \quad (22)$$

where

$$\begin{aligned} \alpha_1(\mathbf{x}) &= \operatorname{Re}\{\mathbf{H}^s \mathbf{x}\} - \operatorname{Im}\{\mathbf{H}^c \mathbf{x}\}, \\ \alpha_2(\mathbf{x}) &= \operatorname{Re}\{\mathbf{H}^s \mathbf{x}\} + \operatorname{Im}\{\mathbf{H}^c \mathbf{x}\}, \\ \mathbf{H}^s &= \frac{\sin(\theta)}{\sigma_w} \operatorname{diag}(s^*) \mathbf{H}, \quad \mathbf{H}^c = \frac{\cos(\theta)}{\sigma_w} \operatorname{diag}(s^*) \mathbf{H}. \end{aligned}$$

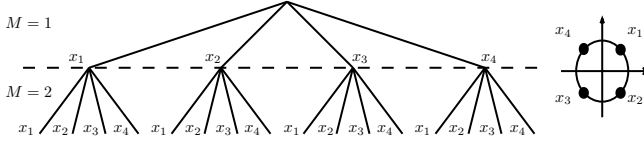


Fig. 3: Tree representation of the set \mathcal{X}^M for a system with $M = 2$ BS antennas and QPSK precoding modulation ($\alpha_x = 4$)

B. Assembling the Tree

A B&B algorithm is a tree search-based method where the tree represents the feasible set which, in this study, is \mathcal{X}^M . For the construction of the tree, it is considered that the p -th BS antenna represents the p -th layer and each possibility for a subvector $\mathbf{f} \in \mathcal{X}^p$ represents one branch. With this, the tree has M layers with α_x^M branches in the last layer. An example of a tree representation of the feasible set is shown in Fig. 3 for the case of $M = 2$ BS antennas and QPSK signaling.

C. Subproblem formulation

In a B&B algorithm, a DPP is solved by considering partially fixed subvectors and computing upper and lower bounds to evaluate if the fixed subvector is part of the optimal solution.

As mentioned in subsection IV-B, the branches of the tree represent a subvector $\mathbf{f} \in \mathcal{X}^p$ for the p -th layer. With this, the subproblems are derived by fixing, for each branch, the corresponding subvector \mathbf{f} and optimizing the remaining subvector $\mathbf{v} \in \mathcal{X}^{M-p}$, with the total vector given by $\mathbf{x} = [\mathbf{f}^T, \mathbf{v}^T]^T$. Yet, such that the optimization problems are real-valued, \mathbf{x}_r is considered instead of \mathbf{x} , with this

$$\mathbf{x}_r = [\mathbf{f}_r^T, \mathbf{v}_r^T]^T, \quad (23)$$

where for the p -th layer, the length of the fixed subvector \mathbf{f}_r is $2p$ and, consequently, the length of the subvector \mathbf{v}_r is $2(M-p)$. The subproblems are then derived based on the formulation presented by (21) for the relaxed problem.

The MUBSEP subproblems are written considering the minimization of the real-valued version of objective function shown in (19) for a given \mathbf{f}_r . To this end, the matrices \mathbf{H}_R^{s*} and \mathbf{H}_I^{s*} are divided as

$$\mathbf{H}_R^{s*} = [\mathbf{F}_{R,\theta}^{s*}, \mathbf{V}_{R,\theta}^{s*}] \quad \mathbf{H}_I^{s*} = [\mathbf{F}_{I,\theta}^{s*}, \mathbf{V}_{I,\theta}^{s*}]. \quad (24)$$

With this, the convex subproblem associated with the fixed subvector \mathbf{f}_r is given by

$$\begin{aligned} & \min_{\mathbf{v}_r} -\mathbf{1}_K^T (\log(\operatorname{erf}(\lambda_1(\mathbf{v}_r)) + \operatorname{erf}(\lambda_2(\mathbf{v}_r)))) \\ & \text{s.t.} \quad \begin{bmatrix} \mathbf{V}_{R,\theta}^{s*} - \mathbf{V}_{I,\theta}^{s*} \\ \mathbf{V}_{R,\theta}^{s*} + \mathbf{V}_{I,\theta}^{s*} \end{bmatrix} \mathbf{v}_r \succeq - \begin{bmatrix} \mathbf{F}_{R,\theta}^{s*} - \mathbf{F}_{I,\theta}^{s*} \\ \mathbf{F}_{R,\theta}^{s*} + \mathbf{F}_{I,\theta}^{s*} \end{bmatrix} \mathbf{f}_r, \\ & \quad \mathbf{R}' [\mathbf{v}_r^T, \mathbf{1}]^T \preceq \mathbf{0}. \end{aligned} \quad (25)$$

where

$$\lambda_1(\mathbf{v}_r) = \mathbf{F}_{R,\theta}^{s*} \mathbf{f}_r + \mathbf{V}_{R,\theta}^{s*} \mathbf{v}_r - \mathbf{F}_{I,\theta}^{s*} \mathbf{f}_r + \mathbf{V}_{I,\theta}^{s*} \mathbf{v}_r \quad (26)$$

$$\lambda_2(\mathbf{v}_r) = \mathbf{F}_{R,\theta}^{s*} \mathbf{f}_r + \mathbf{V}_{R,\theta}^{s*} \mathbf{v}_r + \mathbf{F}_{I,\theta}^{s*} \mathbf{f}_r + \mathbf{V}_{I,\theta}^{s*} \mathbf{v}_r. \quad (27)$$

and \mathbf{R}' represents the last $2(M-p)$ columns of \mathbf{R} .

D. Tree Search Process

For the tree search process, breadth-first search is considered. The process starts by setting the layer value $p = 1$ and, accordingly, solving the subproblems which yields the solution $\mathbf{v}_{r,lb|\mathbf{f}}$. The vector $\mathbf{x}_{lb|\mathbf{f}} = [\mathbf{f}^T, C(\mathbf{v}_{r,lb|\mathbf{f}})^T]^T$ is, then, constructed and the value of $g(\mathbf{x}_{lb|\mathbf{f}})$ is computed and stored. The solution subvector $\mathbf{v}_{lb|\mathbf{f}}$ is quantized to \mathcal{X}^{M-p} which yields $\mathbf{v}_{ub|\mathbf{f}} = Q(\mathbf{v}_{lb|\mathbf{f}})$. With this, one can construct $\mathbf{x}_{ub|\mathbf{f}} = [\mathbf{f}^T, \mathbf{v}_{ub|\mathbf{f}}^T]^T$. Note that $\mathbf{x}_{ub|\mathbf{f}} \in \mathcal{X}^M$ is an upper bound solution meaning $g(\mathbf{x}_{opt}) \leq g(\mathbf{x}_{ub|\mathbf{f}})$, with \mathbf{x}_{opt} being the optimal solution.

To evaluate if $g(\mathbf{x}_{ub|\mathbf{f}})$ is the smallest known upper bound the condition $g(\mathbf{x}_{ub|\mathbf{f}}) < \check{g}$ is checked. If true, the smallest known upper bound and its corresponding value of \mathbf{x} are updated as $\check{g} = g(\mathbf{x}_{ub|\mathbf{f}})$ and $\check{\mathbf{x}} = \mathbf{x}_{ub|\mathbf{f}}$.

After all possible valid branches in one layer are evaluated, i.e., all valid values of \mathbf{f} were fixed and its conditioned upper and lower bounds computed, the lower bounds are evaluated against \check{g} . If $\check{g} < g(\mathbf{x}_{lb|\mathbf{f}})$ then conditioned subvector \mathbf{f} cannot be a subvector of the optimal solution \mathbf{x}_{opt} and \mathbf{f} and all its evolutions can be excluded from the search process. In the context of tree search, this means pruning the branch of \mathbf{f} from the tree.

After pruning, the set of valid \mathbf{f} subvectors is updated and the algorithm repeats this process in the next layer. In the last layer, it is expected that only a few valid candidate solutions remain. With this, they are all evaluated against $\check{\mathbf{x}}$ and the optimal value is determined by the vector that yields the minimum value of the objective function. The steps of the MUBSEP B&B algorithm are detailed in Algorithm 1.

Algorithm 1 Proposed MUBSEP B&B Precoding Algorithm

Inputs: $\mathbf{H}, \mathbf{s}, \mathcal{X}$ **Output:** \mathbf{x}_{opt}
Solve (21) to get $\mathbf{x}_{r,lb}$ and compute $\mathbf{x}_{ub} = Q(C(\mathbf{x}_{r,lb}))$
If $\mathbf{x}_{ub} == C(\mathbf{x}_{r,lb})$
 return $\mathbf{x}_{opt} = \mathbf{x}_{ub}$
end if
Define $\check{\mathbf{x}} = \mathbf{x}_{ub}$ and compute $\check{g} = g(\check{\mathbf{x}})$ using (22)
Define the first level ($p = 1$) of the tree by $\mathcal{G}_p := \mathcal{X}$
for $p = 1 : M - 1$ **do**
 Partition \mathcal{G}_p in $\mathbf{f}_1, \dots, \mathbf{f}_{|\mathcal{G}_p|}$
 for $i = 1 : |\mathcal{G}_p|$ **do**
 Conditioned on $\mathbf{f}_{r,i} = R(\mathbf{f}_i)$ solve (25) to get $\mathbf{v}_{r,lb|\mathbf{f}_i}$
 Construct $\mathbf{x}_{lb,i} = [\mathbf{f}_i^T, C(\mathbf{v}_{r,lb|\mathbf{f}_i})^T]^T$
 Determine the lower bound $g_{lb,i} = g(\mathbf{x}_{lb,i})$
 Compute the upper bound solution $\mathbf{x}_{ub,i} = Q(\mathbf{x}_{lb,i})$
 Compute the upper bound $g_{ub,i} = g(\mathbf{x}_{ub,i})$ with (22)
 Update $\check{g} = \min(\check{g}, g_{ub,i})$ and update $\check{\mathbf{x}}$ accordingly
 end for
 Construct a reduced set by comparing conditioned lower bounds with the global upper bound \check{g}
 $\mathcal{G}'_p := \{\mathbf{x}_{lb,i} \mid g_{lb,i} < \check{g}, i = 1, \dots, |\mathcal{G}_p|\}$
 Define the set for the next level in the tree: $\mathcal{G}_{p+1} := \mathcal{G}'_p \times \mathcal{X}$
end for
The global solution is $\mathbf{x}_{opt} = \operatorname{argmin}_{\mathbf{x} \in \{\mathcal{G}_M \cup \{\check{\mathbf{x}}\}\}} g(\mathbf{x})$

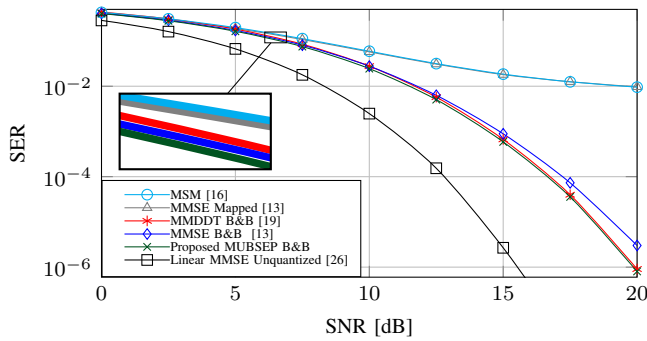


Fig. 4: SER versus SNR for $K = 3$, $M = 14$, $\alpha_s = 8$, $\alpha_x = 4$

V. NUMERICAL RESULTS

In this section, the proposed MUBSEP B&B algorithm is evaluated in terms of SER. To this end, the SNR is defined as $\text{SNR} = \|\mathbf{x}\|_2^2 / \sigma_w^2$. In terms of the channel coefficients Rayleigh fading is considered [25] as done in [5], [15], [16].

The proposed method is evaluated against the following state-of-the-art approaches: 1-The 1-bit MMDDT B&B precoder [19]; 2-The MMSE B&B precoder [13]; 3- The MMSE Mapped precoder [13]; 4- The MSM precoder [16] and 5- The unquantized Linear MMSE precoder [26].

The considered MIMO scenario has a BS with $M = 14$ antennas which serve $K = 3$ users with user symbols drawn for an 8-PSK modulation and transmit symbols drawn from a QPSK modulation, meaning that $\alpha_s = 8$ and $\alpha_x = 4$. The results are depicted in Fig. 4

Fig. 4 shows that the proposed MUBSEP B&B method outperforms, in terms of SER, the MMSE B&B method for the intermediate and high-SNR regimes while presenting similar SER for low-SNR. Moreover, the proposed MUBSEP B&B approach outperforms all other methods for all examined SNR values.

VI. CONCLUSION

In this study, the novel MUBSEP criterion is formulated. Based on the novel formulation an optimal low-resolution precoding method was proposed using the branch-and-bound algorithm. Numerical results show that the proposed precoding method outperforms the state-of-the-art techniques for medium and high SNR values.

REFERENCES

- [1] L. U. Khan, I. Yaqoob, M. Imran, Z. Han, and C. S. Hong, "6G Wireless Systems: A Vision, Architectural Elements, and Future Directions," *IEEE Access*, vol. 8, pp. 147 029–147 044, 2020.
- [2] F. Rusek, D. Persson, B. K. Lau, E. G. Larsson, T. L. Marzetta, O. Edfors, and F. Tufvesson, "Scaling Up MIMO: Opportunities and Challenges with Very Large Arrays," *IEEE Signal Process. Mag.*, vol. 30, no. 1, 2013.
- [3] A. Mezghani and J. A. Nossek, "Power efficiency in communication systems from a circuit perspective," in *2011 IEEE International Symposium of Circuits and Systems (ISCAS)*, 2011, pp. 1896–1899.
- [4] —, "Modeling and minimization of transceiver power consumption in wireless networks," in *2011 International ITG Workshop on Smart Antennas (WSA 2011)*, 2011.

- [5] S. K. Mohammed and E. G. Larsson, "Per-Antenna Constant Envelope Precoding for Large Multi-User MIMO Systems," *IEEE Trans. Commun.*, vol. 61, no. 3, pp. 1059–1071, March 2013.
- [6] S. Jacobsson, G. Durisi, M. Coldrey, T. Goldstein, and C. Studer, "Quantized Precoding for Massive MU-MIMO," *IEEE Trans. Commun.*, vol. 65, no. 11, pp. 4670–4684, 2017.
- [7] A. Mezghani, D. Plabst, L. A. Swindlehurst, I. Fijalkow, and J. A. Nossek, "Sparse Linear Precoders For Mitigating Nonlinearities In Massive MIMO," in *2021 IEEE Statistical Signal Processing Workshop (SSP)*, 2021, pp. 391–395.
- [8] S. Jacobsson, G. Durisi, M. Coldrey, T. Goldstein, and C. Studer, "Nonlinear 1-bit precoding for massive MU-MIMO with higher-order modulation," in *2016 50th Asilomar Conference on Signals, Systems and Computers*, 2016, pp. 763–767.
- [9] S. Jacobsson, O. Castañeda, C. Jeon, G. Durisi, and C. Studer, "Non-linear precoding for phase-quantized constant-envelope massive MU-MIMO-OFDM," in *2018 25th International Conference on Telecommunications (ICT)*, 2018.
- [10] A. Noll, H. Jedda, and J. Nossek, "PSK Precoding in Multi-User MISO Systems," in *WSA 2017; 21th International ITG Workshop on Smart Antennas*, Berlin, Germany, 2017, pp. 1–7.
- [11] A. Li, F. Liu, C. Masouros, Y. Li, and B. Vucetic, "Interference Exploitation 1-Bit Massive MIMO Precoding: A Partial Branch-and-Bound Solution With Near-Optimal Performance," *IEEE Trans. Wireless Commun.*, vol. 19, no. 5, pp. 3474–3489, 2020.
- [12] S. Jacobsson, W. Xu, G. Durisi, and C. Studer, "MSE-optimal 1-bit precoding for multiuser MIMO via branch and bound," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, Calgary, Alberta, Canada, April 2018, pp. 3589–3593.
- [13] E. S. P. Lopes and L. T. N. Landau, "Optimal and Suboptimal MMSE Precoding for Multiuser MIMO Systems Using Constant Envelope Signals with Phase Quantization at the Transmitter and PSK Modulation," in *WSA 2020; 24th International ITG Workshop on Smart Antennas*, Hamburg, Germany, 2020.
- [14] E. S. P. Lopes and L. T. N. Landau, "Discrete MMSE Precoding for Multiuser MIMO Systems with PSK Modulation," *IEEE Trans. Wireless Commun.*, 2022, accepted.
- [15] P. V. Amadori and C. Masouros, "Constant envelope precoding by interference exploitation in phase shift keying-modulated multiuser transmission," *IEEE Trans. Commun.*, Jan 2017.
- [16] H. Jedda, A. Mezghani, A. L. Swindlehurst, and J. A. Nossek, "Quantized constant envelope precoding with PSK and QAM signaling," *IEEE Trans. Wireless Commun.*, vol. 17, no. 12, pp. 8022–8034, Dec 2018.
- [17] F. Askerbeyli, W. Xu, and J. A. Nossek, "1-Bit Precoding for Massive MIMO Downlink with Linear Programming and a Greedy Algorithm Extension," in *2021 IEEE 93rd Vehicular Technology Conference (VTC2021-Spring)*, 2021, pp. 1–5.
- [18] G.-J. Park and S.-N. Hong, "Construction of 1-Bit Transmit-Signal Vectors for Downlink MU-MISO Systems With PSK Signaling," *IEEE Trans. Veh. Technol.*, vol. 68, no. 8, pp. 8270–8274, 2019.
- [19] L. T. N. Landau and R. C. de Lamare, "Branch-and-bound precoding for multiuser MIMO systems with 1-bit quantization," *IEEE Wireless Commun. Lett.*, vol. 6, no. 6, pp. 770–773, Dec 2017.
- [20] E. S. P. Lopes and L. T. N. Landau, "Optimal Precoding for Multiuser MIMO Systems With Phase Quantization and PSK Modulation via Branch-and-Bound," *IEEE Wireless Commun. Lett.*, 2020.
- [21] F. Liu, C. Masouros, P. V. Amadori, and H. Sun, "An Efficient Manifold Algorithm for Constructive Interference Based Constant Envelope Precoding," *IEEE Signal Process. Lett.*, 2017.
- [22] A. Mezghani and R. W. Heath, "Massive MIMO Precoding and Spectral Shaping with Low Resolution Phase-only DACs and Active Constellation Extension," *IEEE Trans. Wireless Commun.*, 2022.
- [23] G. Boole, *The Mathematical Analysis of Logic*. Philosophical Library, 1847.
- [24] S. Boyd and L. Vandenberghe, *Convex Optimization*. New York, NY, USA: Cambridge University Press, 2004.
- [25] H. Yang and T. L. Marzetta, "Performance of conjugate and zero-forcing beamforming in large-scale antenna systems," *IEEE J. Sel. Areas Commun.*, vol. 31, no. 2, pp. 172–179, 2013.
- [26] M. Joham, W. Utschick, and J. A. Nossek, "Linear transmit processing in MIMO communications systems," *IEEE Trans. Signal Process.*, vol. 53, no. 8, pp. 2700–2712, Aug 2005.