

# An Efficient and Unified Framework for Downlink Linear Precoding with QoS Constraints

Ruiding Hou, *Graduate Student Member, IEEE*, Jiaheng Wang, *Senior Member, IEEE*, Rui Zhou, *Member, IEEE*, Daniel P. Palomar, *Fellow, IEEE*, Xiqi Gao, *Fellow, IEEE*, Björn Ottersten, *Fellow, IEEE*

**Abstract**—Precoding techniques, particularly linear precoding, are widely employed in multiple-input multiple-output (MIMO) systems. Although well-studied in the literature, linear precoding design still faces two fundamental challenges: high computational complexity and the lack of a general design approach. This paper presents an efficient and unified framework for linear precoding design in downlink multiuser systems that accommodates diverse criteria, such as weighted sum rate (WSR) maximization and weighted symbol error rate (WSER) minimization, while ensuring quality of service (QoS) requirements. The proposed framework achieves an order-of-magnitude reduction in per-iteration computational complexity compared to existing methods. In particular, by accurately characterizing the feasible signal-to-interference-plus-noise ratio (SINR) region, we transform the high-dimensional precoding design problem into a more manageable, low-dimensional SINR allocation problem. We propose an efficient SINR-based precoding (SBP) framework that employs a water-filling solution at each iteration, without the need for matrix inversion. The proposed framework can be extended to broadcast and interference channels with multi-antenna users under pre-fixed receivers. Simulation results demonstrate that our method achieves near-optimal performance while significantly reducing computational complexity compared to existing methods, such as the weighted minimum mean square error (WMMSE) method.

**Index Terms**—Downlink linear precoding, QoS requirements, SINR-based precoding, efficient algorithms.

This work was supported in part by the Science and Technology Major Project of Jiangsu under Grant BG2025038, the Natural Science Foundation on Frontier Leading Technology Basic Research Project of Jiangsu under Grants BK20212001 and BK20222001, the Key Technologies R&D Program of Jiangsu (Prospective and Key Technologies for Industry) under Grants BE2022068 and BE2022068-3, the National Natural Science Foundation of China under Grants 62331024, U22B2006, and 62201362, the Taihu Lake Innovation Fund for the School of Future Technology of Southeast University, the Fundamental Research Funds for the Central Universities under Grants 2242022K60002 and 2242022K60001, and the Shenzhen Science and Technology Program under Grant JCYJ20250604191208011. The work of Daniel P. Palomar was supported by the Hong Kong GRF 16206123 research grant. The work of Björn Ottersten was supported by the Luxembourg National Research Fund (FNR), grant reference INTER/MOBILITY/2023/IS/18014377/MCR. (*Corresponding author: Jiaheng Wang.*)

Ruiding Hou, Jiaheng Wang and Xiqi Gao are with the National Mobile Communications Research Laboratory, Southeast University, Nanjing 210096, China, and also with the Purple Mountain Laboratory, Nanjing 210023, China. (e-mail: rdhou@seu.edu.cn, jhwang@seu.edu.cn, xqgao@seu.edu.cn).

Rui Zhou is with Shenzhen Research Institute of Big Data, Chinese University of Hong Kong-Shenzhen, Shenzhen 518172, China (e-mail: rui.zhou@sribd.cn).

Daniel P. Palomar is with the Hong Kong University of Science and Technology (HKUST), Clear Water Bay, Kowloon, Hong Kong (e-mail: palomar@ust.hk).

Björn Ottersten is with the Interdisciplinary Centre for Security, Reliability and Trust, University of Luxembourg, 1855 Luxembourg City, Luxembourg (e-mail: bjorn.ottersten@uni.lu).

## I. INTRODUCTION

MULTIPLE-Input multiple-output (MIMO) transmission is a key enabler for the current 5G and future 6G mobile systems [1], where precoding technologies are employed to improve data rates and maintain quality of service (QoS) [2]. Linear precoding, by linearly pre-processing the transmitted signal according to the channel state information (CSI), can effectively suppress inter-user interference [3] and has attracted significant attention in practical systems due to its low implementation complexity [4].

Linear precoding is typically customized by a range of specific criteria, such as maximizing the weighted sum rate (WSR) [5]–[12], reducing the weighted mean square error (WMSE) [13]–[15], or minimizing the weighted symbol error rate (WSER) [16], [17], yielding different optimization problems. Among various linear precoding design problems, only two are globally solved, i.e., minimizing transmit power under QoS requirements [2], [18], [19], and maximizing the worst signal-to-interference-plus-noise ratio (SINR) [2], [20], which can be transformed into a convex problem. In most cases, unfortunately, linear precoding problems are non-convex and even NP-hard [21]. Therefore, instead of searching for the globally optimal solution with prohibitive complexity, efficient algorithms that can provide locally optimal or satisfactory performance with low computational complexity are often preferred in practice. Particularly, in [5]–[8], the successive convex approximation (SCA) methods were proposed to iteratively approximate the WSR maximization problem with a series of simpler convex problems. Alternatively, the weighted minimum mean square error (WMMSE) method [9]–[11] has emerged as a widely adopted method for WSR maximization, benefiting from its closed-form update in each iteration. Furthermore, algorithms combining the WMMSE method with the alternating direction method of multipliers (ADMM) were proposed in [12], [22] to solve the QoS-constrained WSR maximization problem. Meanwhile, alternative optimization methods were proposed for WMSE minimization [13], [14] and WSER minimization [17]. In [23], the bilevel optimization method was introduced to maximize the resource efficiency. However, existing precoding schemes, especially linear precoding methods, still have several principal limitations.

First, current approaches suffer from high computational costs, especially in systems with a large number of users and antennas. In particular, the SCA methods [5]–[8] require solving a series of convex problems through standard numerical algorithms, e.g., interior-point methods. Similarly, the

bilevel optimization methods [23] have to repeatedly solve its subproblems via uplink-downlink duality (UDD) or fixed-point iteration (FPI) methods [2], [18], [19]. Despite their satisfactory performance, the computational complexity of these methods is prohibitive due to the absence of closed-form iterative steps. The WMMSE methods [9]–[11], though possessing closed-form updates, require repeatedly computing matrix inversions, each with a computational complexity at least scaling cubically with the number of users [11], hindering their implementation in practical systems. Some studies simplify the problem by fixing the direction of precoding matrix, e.g., zero-forcing (ZF) precoding [3], which, however, results in performance degradation.

Second, another limitation of existing linear precoding methods is their lack of generality and insufficient consideration of QoS constraints. Specifically, most algorithms are tailored to maximize a specific system utility, i.e., particular performance metrics that the system aims to optimize. For example, the WMMSE algorithm is designed for WSR maximization, relying on a specific logarithmic objective function [9]. Similarly, SCA-based methods are inherently problem-specific [5]–[8], [17], often requiring case-by-case derivations and yielding solutions that vary significantly even for the same objective [5]–[8]. Furthermore, many of these methods fail to explicitly incorporate QoS constraints, which are crucial for ensuring reliable communication [2]. Conversely, precoding designs that do account for QoS typically rely on complex, iterative algorithms with prohibitively high computational overhead [12], [22]–[25]. For instance, when user QoS constraints are included, the WMMSE algorithm no longer has a closed-form update, and each iteration requires solving a high-dimensional optimization problem [12], [22]. This critical gap highlights the need for a more efficient and unified precoding framework that can accommodate diverse performance criteria while simultaneously satisfying practical QoS requirements.

This paper presents an efficient and unified framework for downlink linear precoding design to maximize system utility under sum power and QoS constraints. The proposed framework can be applied to various system utilities, and its computational complexity is an order-of-magnitude lower than that of existing methods in each iteration.

The main contributions of this paper are outlined as follows.

- We consider a unified linear precoding design in downlink multiuser systems under a variety of commonly used system utility functions, such as WSR maximization and WSER minimization, and formulate the precoding design problem as a general utility maximization problem with QoS constraints. Then, we reformulate this precoding design problem into an equivalent low-dimensional SINR allocation problem and approximate it by introducing an alternative feasible SINR region, which possesses a closed-form expression. We analytically prove that this feasible SINR region is asymptotically tight.
- We propose a novel SINR-based precoding (SBP) design framework, consisting of two main stages: solving the approximate SINR allocation problem, and constructing the precoding matrix from the optimized SINR values. Our approach efficiently solves the SINR allocation prob-

lem by using a water-filling solution in each iteration, which dramatically reduces per-iteration complexity by eliminating the need for matrix inversion and converges to a Karush-Kuhn-Tucker (KKT) stationary point. For the second stage, we develop an efficient algorithm with lower complexity and a faster convergence rate. Overall, our SBP framework provides a complete and highly efficient solution for the linear precoding design.

- Furthermore, the proposed framework is extended to the linear precoding design for interference channels as well as broadcast and interference channels with multi-antenna users under pre-fixed linear or successive interference cancellation (SIC) receivers.

The remainder of this paper is organized as follows. Section II introduces the system model and problem formulation. Section III focuses on the transformation to an SINR allocation problem and the derivation of an analytical or approximate expression of the feasible SINR region. Section IV introduces the proposed framework and provides its convergence and complexity analysis. Section V extends the proposed framework to other applicable scenarios. Numerical results are provided in Section VI, and Section VII concludes the paper.

*Notation:* Bold lower and upper case letters denote vectors and matrices. The superscripts  $(\cdot)^*$ ,  $(\cdot)^T$ , and  $(\cdot)^H$  are the conjugate, transpose, and conjugate transpose.  $\mathbf{I}$ ,  $\mathbf{1}_{n \times k}$ , and  $\mathbf{0}_{n \times k}$  are the identity, all-ones, and all-zeros matrices. For a matrix  $\mathbf{A}$ ,  $[\mathbf{A}]_{mn}$ ,  $\mathbf{A}^{-1}$ ,  $\text{diag}(\mathbf{A})$ , and  $\|\mathbf{A}\|_F$  denote its  $(m, n)$ -th element, inverse, a vector of its diagonal elements, and its Frobenius norm. For a vector  $\mathbf{a}$ ,  $\text{Diag}(\mathbf{a})$  is a diagonal matrix with  $\mathbf{a}$  on its diagonal. The operators  $\odot$ ,  $\mathbb{E}[\cdot]$ , and  $|\cdot|$  are the Hadamard product, expectation, and absolute value.  $\|\cdot\|_2$  and  $\|\cdot\|_\infty$  are the Euclidean and infinity norms. Finally,  $\mathbf{a} \geq \mathbf{b}$  denotes elementwise inequality and  $\mathbf{A} \succeq \mathbf{0}$  means  $\mathbf{A}$  is positive semi-definite.

## II. PROBLEM STATEMENT

### A. System Model

Consider a downlink multiuser system comprising a BS equipped with  $N$  antennas, serving  $K$  single-antenna users simultaneously. It is assumed that all users share the same time-frequency resources. The transmitted signal is

$$\mathbf{x} = \sum_{k=1}^K \mathbf{w}_k s_k, \quad (1)$$

where  $\mathbf{w}_k \in \mathbb{C}^{N \times 1}$  denotes the linear precoder of user  $k$  and  $s_k$  is the transmitted data of user  $k$ , which satisfies  $\mathbb{E}[s_k] = 0$ ,  $\mathbb{E}[|s_k|^2] = 1$ , and  $\mathbb{E}[s_j^* s_k] = 0, \forall j \neq k$ . The received signal of user  $k$  in the broadcast channel model [26]–[28] is

$$y_k = \mathbf{h}_k^H \mathbf{x} + n_k = \mathbf{h}_k^H \mathbf{w}_k s_k + \sum_{j \neq k}^K \mathbf{h}_k^H \mathbf{w}_j s_j + n_k, \quad (2)$$

where  $\mathbf{h}_k \in \mathbb{C}^{N \times 1}$  is the channel matrix from BS to user  $k$  and  $n_k \sim \mathcal{CN}(0, \sigma_k^2)$  is the additive noise at user  $k$ , with  $\sigma_k^2$  denoting its variance. Given precoding matrix  $\mathbf{W} \triangleq [\mathbf{w}_1, \dots, \mathbf{w}_K] \in \mathbb{C}^{N \times K}$ , the SINR of user  $k$  is

$$\gamma_k(\mathbf{W}) = \frac{|\mathbf{h}_k^H \mathbf{w}_k|^2}{\sum_{j \neq k}^K |\mathbf{h}_k^H \mathbf{w}_j|^2 + \sigma_k^2}, \quad \forall k. \quad (3)$$

Without loss of generality, we assume  $\sigma_k^2 = 1, \forall k$ .<sup>1</sup>

### B. Problem Formulation

We consider a linear precoding design framework incorporating both QoS and power constraints, formulated as follows:

$$\begin{aligned} & \underset{\mathbf{W}}{\text{maximize}} && f(\gamma(\mathbf{W})) \triangleq \sum_{k=1}^K f_k(\gamma_k(\mathbf{W})) \\ & \text{subject to} && \gamma(\mathbf{W}) \geq \gamma_{\text{th}}, \\ & && \sum_{k=1}^K \|\mathbf{w}_k\|_2^2 \leq P_{\text{sum}}, \end{aligned} \quad (4)$$

where  $\gamma(\mathbf{W}) \triangleq [\gamma_1(\mathbf{W}), \dots, \gamma_K(\mathbf{W})]^T$  and  $\gamma_{\text{th}} \triangleq [\gamma_{1,\text{th}}, \dots, \gamma_{K,\text{th}}]^T$ . Each function  $f_k(\rho_k)$  is continuously differentiable and monotonically increasing with respect to  $\rho_k$ .

This framework includes a variety of commonly used system utility functions. Specific examples are detailed below:

- **WSR:** Assuming  $s_k$  is a Gaussian signal and invoking the Shannon's capacity formula, the achievable rate of user  $k$  is  $\log(1 + \rho_k)$ . Then, the objective function for WSR maximization problem is  $\sum_{k=1}^K \omega_k \log(1 + \rho_k)$ , where  $\omega_k \geq 0$  is the weight of user  $k$  to prioritize the users.
- **WMSE:** Assuming that each user employs the linear minimum mean square error (LMMSE) equalizer, the mean square error (MSE) of user  $k$  is  $(1 + \rho_k)^{-1}$  [29]. Then the objective function for WMSE minimization problem is  $-\sum_{k=1}^K \omega_k (1 + \rho_k)^{-1}$ .
- **WSER:** For a given modulation, e.g., quadrature amplitude modulation (QAM), the symbol error rate (SER) of user  $k$  can be expressed as  $\text{SER}_k(\rho_k) = a_k Q(\sqrt{b_k \rho_k})$ , where  $a_k, b_k$  depend on its signal constellation [30] and the  $Q$  function is defined as  $Q(x) = (1/\sqrt{2\pi}) \int_x^\infty e^{-z^2/2} dz$ . In this case, the objective function is  $-\sum_{k=1}^K \omega_k \text{SER}_k(\rho_k)$ .

In addition to the system utility functions listed above, our framework also includes other criteria, e.g., maximizing the WSR and minimizing the weighted block error rate in finite-block-length systems.

Note that problem (4) is usually difficult to solve due to the non-convexity of  $\gamma_k(\mathbf{W})$  and is an NP-hard problem in several cases [21]. Besides, since the precoding matrix  $\mathbf{W}$  possesses a real dimension of  $2NK$ , solving problem (4) directly entails significant computational complexity when  $NK$  is large.

## III. ANALYSIS OF FEASIBLE SINR REGION

### A. Equivalent Transformation

Enlightened by previous work [19], [23], a one-to-one correspondence can be established between the precoding matrix  $\mathbf{W} \in \mathbb{C}^{N \times K}$  and the SINR vector  $\boldsymbol{\rho} = (\rho_1, \dots, \rho_K)^T \in \mathbb{R}^K$ . Specifically, the optimal precoding matrix  $\mathbf{W}^*$  to problem (4) is identified as the one with the lowest Frobenius norm among all the precoding matrices that achieve the same SINR values  $\gamma(\mathbf{W}^*)$ . Consequently, once the optimal SINR values  $\boldsymbol{\rho}^* = \gamma(\mathbf{W}^*)$  for problem (4) are determined, the corresponding  $\mathbf{W}^*$  can be constructed. This motivates reformulating the precoding design problem (4) as an SINR allocation problem.

<sup>1</sup>For the case where  $\sigma_k \neq 1$ , we can set an equivalent channel  $\tilde{\mathbf{h}}_k = \mathbf{h}_k/\sigma_k$  and  $\tilde{\sigma}_k = 1$ .

The transformation is particularly advantageous because the objective function is concave in SINR values and its optimization variable  $\boldsymbol{\rho}$  only has a dimensionality of  $K$ . Based on this principle, problem (4) is equivalently transformed into the following more tractable optimization problem:

$$\begin{aligned} & \underset{\boldsymbol{\rho}}{\text{maximize}} && f(\boldsymbol{\rho}) \\ & \text{subject to} && \boldsymbol{\rho} \in \mathcal{R}, \end{aligned} \quad (5)$$

where the feasible SINR region  $\mathcal{R}$  is the set of all possible SINR vectors  $\gamma(\mathbf{W})$  that can be realized by some feasible precoding matrix  $\mathbf{W}$ , given as:

$$\mathcal{R} = \left\{ \gamma(\mathbf{W}) \mid \gamma(\mathbf{W}) \geq \gamma_{\text{th}}, \sum_{k=1}^K \|\mathbf{w}_k\|_2^2 \leq P_{\text{sum}} \right\}. \quad (6)$$

The primary challenge in solving problem (5) lies in the non-convexity of  $\mathcal{R}$  [31]. Additionally, since  $\mathcal{R}$  generally lacks an analytical expression, previous work [23] has to iteratively update  $\boldsymbol{\rho}$  and compute  $\mathbf{W}$  for each updated  $\boldsymbol{\rho}$ , leading to substantial computational complexity. If  $\mathcal{R}$  can be expressed or approximated analytically, problem (5) is expected to be significantly simplified.

### B. Approximation of Feasible SINR Region

Notice that  $\boldsymbol{\rho} \in \mathcal{R}$  ensures that the minimum required power across all precoding matrices satisfying the QoS constraints remains within the available power budget  $P_{\text{sum}}$ . Let  $P(\boldsymbol{\rho})$  be the optimal objective value of the following minimum power (MP) problem under given user SINRs  $\boldsymbol{\rho}$ :

$$\begin{aligned} & \underset{\{\mathbf{w}_k\}_{k=1}^K}{\text{minimize}} && \sum_{k=1}^K \|\mathbf{w}_k\|_2^2 \\ & \text{subject to} && 1 + \sum_{j \neq k} |\mathbf{h}_k^H \mathbf{w}_j|^2 \leq \rho_k^{-1} |\mathbf{h}_k^H \mathbf{w}_k|^2, \quad \forall k, \end{aligned} \quad (7)$$

which can be transformed into a convex problem [2], [18], [19]. Then, the next result can further simplify the feasible SINR region  $\mathcal{R}$  defined in (6).

**Lemma 1.** *The feasible SINR region  $\mathcal{R}$ , as specified in (6), can be equivalently represented as:*

$$\mathcal{R} = \{\boldsymbol{\rho} \mid \boldsymbol{\rho} \geq \gamma_{\text{th}}, P(\boldsymbol{\rho}) \leq P_{\text{sum}}\}. \quad (8)$$

*Proof.* See Appendix A. □

As shown in Lemma 1,  $\mathcal{R}$  can be derived from the formulation of  $P(\boldsymbol{\rho})$ . To establish a rigorous characterization of the analytical expression of  $P(\boldsymbol{\rho})$ , we consider the dual problem of problem (7) as follows:

$$\begin{aligned} & \underset{\{\lambda_k\}_{k=1}^K}{\text{maximize}} && \sum_{k=1}^K \lambda_k \\ & \text{subject to} && \mathbf{I} + \sum_{j \neq k} \lambda_j \mathbf{h}_j \mathbf{h}_j^H \succeq \rho_k^{-1} \lambda_k \mathbf{h}_k \mathbf{h}_k^H, \quad \forall k. \end{aligned} \quad (9)$$

Let the optimal solution to problems (7) and (9) under given user SINRs  $\boldsymbol{\rho}$  be  $\mathbf{W}^*(\boldsymbol{\rho}) \triangleq [\mathbf{w}_1^*(\boldsymbol{\rho}), \dots, \mathbf{w}_K^*(\boldsymbol{\rho})]$  and  $\boldsymbol{\lambda}^*(\boldsymbol{\rho}) \triangleq [\lambda_1^*(\boldsymbol{\rho}), \dots, \lambda_K^*(\boldsymbol{\rho})]^T$ , respectively. As the strong duality of problem (7) holds [2], problems (7) and (9) share the same optimal objective value, i.e.,  $\sum_{k=1}^K \|\mathbf{w}_k^*(\boldsymbol{\rho})\|_2^2 = \sum_{k=1}^K \lambda_k^*(\boldsymbol{\rho}) = P(\boldsymbol{\rho})$ . The subsequent result presents the relationship between  $\boldsymbol{\lambda}^*(\boldsymbol{\rho})$  and  $\mathbf{W}^*(\boldsymbol{\rho})$ .

**Proposition 1.** *The optimal solution  $\lambda^*(\rho)$  to problem (9) satisfies*

$$\Lambda^*(\rho) \odot (\Lambda^*(\rho) + \bar{\mathbf{H}}^{-1})^{-1} = (\mathbf{I} + \Gamma^{-1})^{-1}, \quad (10)$$

where  $\mathbf{H} \triangleq [\mathbf{h}_1, \dots, \mathbf{h}_K] \in \mathbb{C}^{N \times K}$ ,  $\bar{\mathbf{H}} \triangleq \mathbf{H}^H \mathbf{H}$ ,  $\Lambda^*(\rho) \triangleq \text{Diag}(\lambda^*(\rho))$ , and  $\Gamma \triangleq \text{Diag}(\rho)$ . Moreover, the optimal solution  $\mathbf{W}^*(\rho)$  to problem (7) is given by

$$\mathbf{W}^*(\rho) = \mathbf{H} \bar{\mathbf{H}}^{-1} (\bar{\mathbf{H}}^{-1} + \Lambda^*(\rho))^{-1} \text{Diag}(\mathbf{p})^{\frac{1}{2}}, \quad (11)$$

where  $\mathbf{p} \triangleq [(-\mathbf{1}_{K \times K} + \mathbf{I} + \Gamma^{-1}) \odot \mathbf{G}]^{-1} \mathbf{1}_{K \times 1}$  and  $\mathbf{G} \triangleq (\Lambda^*(\rho) + \bar{\mathbf{H}}^{-1})^{-1} \odot [(\Lambda^*(\rho) + \bar{\mathbf{H}}^{-1})^{-1}]^*$ .

*Proof.* See Appendix B.  $\square$

Proposition 1 indicates that, one can directly derive the optimal  $\mathbf{W}^*(\rho)$  to problem (7) from the optimal  $\lambda^*(\rho)$  to problem (9) given  $\text{rank}(\mathbf{H}) = K$ .<sup>2</sup> Conventional approaches for determining  $\lambda^*(\rho)$  usually require several iterations, each with computational complexity of  $O(N^3)$ , such as FPI [2]. To further reduce computational complexity, we aim to discover whether there is an analytical expression of  $\lambda^*(\rho)$  or a well-approximated  $\tilde{\lambda}(\rho) \in \mathbb{R}^K$ , thereby expediting the computation of  $\lambda^*(\rho)$ . For the simpler case  $K = 2$ , the closed-form expression of  $\lambda^*(\rho)$  is presented in Lemma 2.

**Lemma 2.** *When  $K = 2$ , the optimal solution to problem (9) is expressible in a closed-form as shown in (12) and (13), where  $h_{ij} = [\bar{\mathbf{H}}^{-1}]_{ij}$ .*

*Proof.* See Appendix C.  $\square$

Notably, in the proof of Lemma 2, we derive the closed-form solution to problem (9) by solving a univariate second-order polynomial equation. Consequently, from the Abel–Ruffini theorem [32], it is reasonable to conclude that no closed-form solution exists for problem (9) when  $K > 4$ , as there is no solution in radicals for general polynomial equations of degree greater than four. Therefore, we next seek to identify a suitable approximation  $\tilde{\lambda}(\rho)$  for  $\lambda^*(\rho)$ .

**Theorem 1.** *An approximation of  $\lambda^*(\rho)$  is*

$$[\tilde{\lambda}(\rho)]_k \triangleq \alpha_k \rho_k \left( 1 - \sum_{j \neq k} \frac{\beta_{kj}}{1 + \rho_j} \right), \quad \forall k, \quad (14)$$

where  $\alpha_k = [\bar{\mathbf{H}}^{-1}]_{kk}$  and  $\beta_{kj} = |[\bar{\mathbf{H}}^{-1}]_{kj}|^2 \alpha_k^{-1} \alpha_j^{-1}$ . When  $\rho \rightarrow \infty$ ,<sup>3</sup>  $\|\lambda^*(\rho) - \tilde{\lambda}(\rho)\|_\infty \rightarrow 0$ .

<sup>2</sup>In practice, when user channels exhibit strong spatial coherence, frequency or time division is typically used instead of space division multiplexing.

<sup>3</sup>We define  $\rho \rightarrow \infty$  as  $\min_k \{\rho_k\} \rightarrow +\infty$  with  $\max_{k,j} \{\rho_k/\rho_j\} \leq \kappa < +\infty$  for some constant  $\kappa$ .

*Proof.* See Appendix D.  $\square$

Unlike previous work [33], which approximates  $\lambda^*(\rho)$  in the asymptotic regime where  $N, K \rightarrow \infty$  with fixed  $c = K/N$ , Theorem 1 states that  $\lambda^*(\rho)$  can be approximated by  $\tilde{\lambda}(\rho)$  under the mild condition that  $\max_k [\sum_{j \neq k} \frac{\beta_{kj}}{1 + \rho_j}] < 1$  to guarantee  $\tilde{\lambda}(\rho) > \mathbf{0}$ . When  $\rho$  is large, this approximation is asymptotically tight. Since  $P(\rho) = \mathbf{1}^T \lambda^*(\rho) \approx \mathbf{1}^T \tilde{\lambda}(\rho)$ , an approximation of  $P(\rho)$  follows directly from Theorem 1.

**Corollary 1.** *An approximation of  $P(\rho)$  is*

$$\tilde{P}(\rho) \triangleq \sum_{k=1}^K \alpha_k \rho_k \left( 1 - \sum_{j \neq k} \frac{\beta_{kj}}{1 + \rho_j} \right), \quad (15)$$

and when  $\rho \rightarrow \infty$ ,  $|P(\rho) - \tilde{P}(\rho)| \rightarrow 0$ .

Given the expression of  $\tilde{P}(\rho)$ , the approximated feasible SINR region is

$$\begin{aligned} \tilde{\mathcal{R}} &\triangleq \left\{ \rho \mid \rho \geq \gamma_{\text{th}}, \tilde{P}(\rho) \leq P_{\text{sum}} \right\} \\ &= \left\{ \rho \mid \rho \geq \gamma_{\text{th}}, \sum_{k=1}^K \alpha_k \rho_k \left( 1 - \sum_{j \neq k} \frac{\beta_{kj}}{1 + \rho_j} \right) \leq P_{\text{sum}} \right\}, \end{aligned} \quad (16)$$

where  $\gamma_{\text{th}}$  needs to satisfy  $\max_k [\sum_{j \neq k} \frac{\beta_{kj}}{1 + \gamma_{j,\text{th}}}] < 1$ .

*Remark 1:* Consider the well-known linear ZF precoding design problem to minimize the average transmitted power:

$$\begin{aligned} &\text{minimize}_{\{\mathbf{w}_k\}_{k=1}^K} \sum_{k=1}^K \|\mathbf{w}_k\|_2^2 \\ &\text{subject to} \quad 1 - \rho_k^{-1} |\mathbf{h}_k^H \mathbf{w}_k|^2 \leq 0, \quad \forall k, \\ &\quad \quad \quad \mathbf{h}_j^H \mathbf{w}_k = 0, \quad \forall k \neq j. \end{aligned} \quad (17)$$

The optimal objective value of problem (17) is  $P^{\text{ZF}}(\rho) = \sum_{k=1}^K \alpha_k \rho_k$  [3]. Interestingly, Corollary 1 can be expressed as  $\tilde{P}(\rho) = P^{\text{ZF}}(\rho) - \sum_{k=1}^K \sum_{j \neq k} \frac{\alpha_k \beta_{kj} \rho_k}{1 + \rho_j}$ , implying that the performance gap between ZF precoding and the optimal precoding for power minimization is non-negligible unless  $\beta_{kj} = 0, \forall k \neq j$ , which is equivalent to the case where the columns of  $\mathbf{H}$  are orthogonal. The feasible SINR region with ZF precoding is  $\mathcal{R}^{\text{ZF}} = \left\{ \rho \mid \rho \geq \gamma_{\text{th}}, \sum_{k=1}^K \alpha_k \rho_k \leq P_{\text{sum}} \right\}$ .

## IV. SINR-BASED PRECODING

### A. Primary Steps of SBP

The SBP framework for solving problem (4) consists of two primary steps:

**Step 1:** Determining the optimal solution  $\rho^*$  to the approximate SINR allocation problem:

$$\begin{aligned} &\text{maximize}_{\rho} \quad f(\rho) \\ &\text{subject to} \quad \rho \in \tilde{\mathcal{R}}. \end{aligned} \quad (18)$$

$$\lambda_1^*(\rho) = \frac{h_{11}}{2} (\rho_1 - 1) + \frac{|h_{21}|^2 \rho_2 - \rho_1}{2h_{22} \rho_2 + 1} + \frac{\rho_1 + 1}{2h_{22}} \left( h_{11}^2 h_{22}^2 - \frac{2h_{11}h_{22}|h_{21}|^2(\rho_1 + \rho_2)}{(\rho_1 + 1)(\rho_2 + 1)} + \frac{|h_{21}|^4(\rho_2 - \rho_1)^2}{(\rho_1 + 1)^2(\rho_2 + 1)^2} \right)^{\frac{1}{2}}, \quad (12)$$

$$\lambda_2^*(\rho) = \frac{h_{22}}{2} (\rho_2 - 1) + \frac{|h_{12}|^2 \rho_1 - \rho_2}{2h_{11} \rho_1 + 1} + \frac{\rho_2 + 1}{2h_{11}} \left( h_{11}^2 h_{22}^2 - \frac{2h_{11}h_{22}|h_{12}|^2(\rho_1 + \rho_2)}{(\rho_1 + 1)(\rho_2 + 1)} + \frac{|h_{12}|^4(\rho_1 - \rho_2)^2}{(\rho_1 + 1)^2(\rho_2 + 1)^2} \right)^{\frac{1}{2}}. \quad (13)$$

**Step 2:** Constructing the precoding matrix  $\mathbf{W}^{\text{SBP}}$  from  $\boldsymbol{\rho}^*$ .

The core of the SBP framework lies in determining  $\boldsymbol{\rho}^*$  in **Step 1**. However, solving problem (18) is still challenging since  $\hat{\mathcal{R}}$  is not guaranteed to be a convex set. To address this issue, we first introduce the following problem:

$$\begin{aligned} & \underset{\boldsymbol{\rho}, \mathbf{Z}}{\text{maximize}} && f(\boldsymbol{\rho}) \\ & \text{subject to} && \boldsymbol{\rho} \in \hat{\mathcal{R}}(\mathbf{Z}), \end{aligned} \quad (19)$$

where  $\mathbf{Z} = (z_{kj})_{K \times K} \in \mathbb{R}^{K \times K}$  is an auxiliary variable and

$$\hat{\mathcal{R}}(\mathbf{Z}) \triangleq \left\{ \boldsymbol{\rho} \mid \rho \geq \gamma_{\text{th}}, \hat{P}(\boldsymbol{\rho}, \mathbf{Z}) \leq P_{\text{sum}} \right\}, \quad (20)$$

with  $\hat{P}(\boldsymbol{\rho}, \mathbf{Z})$  defined in (21). The equivalence between problems (18) and (19) is given below.

**Proposition 2.** *Let  $(\boldsymbol{\rho}^*, \mathbf{Z}^*)$  be the optimal solution to problem (19). Then,  $\boldsymbol{\rho}^*$  is the optimal solution to problem (18). Moreover, if  $(\bar{\boldsymbol{\rho}}, \bar{\mathbf{Z}})$  is a KKT-stationary point of problem (19), then  $\bar{\boldsymbol{\rho}}$  is a KKT-stationary point of problem (18) when  $\bar{z}_{kj} = \sqrt{\bar{\rho}_k} / (\bar{\rho}_j + 1), \forall j \neq k$ .*

*Proof.* See Appendix E.  $\square$

Building on Proposition 2, we implement the following iterative procedure to solve problem (18):

$$z_{kj}^{(n)} = \sqrt{\rho_k^{(n)}} / (\rho_j^{(n)} + 1), \quad \forall k \neq j, \quad (23)$$

$$\boldsymbol{\rho}^{(n+1)} = \underset{\boldsymbol{\rho}}{\text{argmax}} f(\boldsymbol{\rho}) \text{ s.t. } \boldsymbol{\rho} \in \hat{\mathcal{R}}(\mathbf{Z}^{(n)}). \quad (24)$$

By cyclically updating  $\mathbf{Z}^{(n)}$  and  $\boldsymbol{\rho}^{(n)}$  according to (23) and (24),  $\boldsymbol{\rho}^{(n)}$  will converge to a KKT-stationary point of problem (18), as established in Theorem 4. Next, we explore efficient methods for problem (24).

### B. Optimal Solution to Problem (24)

Problem (24) can be represented as the following problem:

$$\begin{aligned} & \underset{\boldsymbol{\rho} \geq \gamma_{\text{th}}}{\text{maximize}} && \sum_{k=1}^K f_k(\rho_k) \\ & \text{subject to} && \sum_{k=1}^K p_k(\rho_k) \leq P_{\text{sum}}, \end{aligned} \quad (25)$$

where  $p_k(\rho_k) = A_k(\mathbf{Z}^{(n)})\rho_k - 2B_k(\mathbf{Z}^{(n)})\sqrt{\rho_k} + C_k(\mathbf{Z}^{(n)})$  is continuously differentiable and strictly convex with respect to  $\rho_k$ . For simplicity, we denote  $A_k(\mathbf{Z}^{(n)})$ ,  $B_k(\mathbf{Z}^{(n)})$ , and  $C_k(\mathbf{Z}^{(n)})$  as  $A_k$ ,  $B_k$ , and  $C_k$  in the remainder of this section.

When  $B_k = 0, \forall k$ , problem (25) can be solved using the water-filling method [34]. In general cases, the following theorem elucidates the optimal properties of problem (25).

**Theorem 2.** *If each  $f_k(\rho_k)$  is strictly concave with  $\rho_k$ , the optimal solution to problem (25) is*

$$\rho_k^* = \max \{ \hat{\rho}_k(\mu^*), \gamma_{k,\text{th}} \}, \quad \forall k, \quad (26)$$

and satisfies  $\sum_{k=1}^K p_k(\rho_k^*) = P_{\text{sum}}$ , where  $\hat{\rho}_k(\mu)$  satisfies  $-f'_k(\hat{\rho}_k(\mu)) + \mu p'_k(\hat{\rho}_k(\mu)) = 0$ .

*Proof.* See Appendix F.  $\square$

Theorem 2 demonstrates that the optimal solution to problem (25) admits a water-filling structure and can be readily obtained using (26) once  $\mu^*$  is determined. The next result paves the way for computing  $\mu^*$ .

**Lemma 3.** *Let  $\mathcal{P}(\mu) \triangleq \sum_{k=1}^K p_k(\max \{ \hat{\rho}_k(\mu), \gamma_{k,\text{th}} \})$ . Then,  $\mathcal{P}(\mu)$  is a non-increasing function of  $\mu$  if  $f_k(\rho_k)$  is strictly concave with  $\rho_k$ .*

*Proof.* See Appendix G.  $\square$

Based on Lemma 3, the optimal  $\mu^*$  can be efficiently determined through the bisection method and the main step of solving problem (24) lies in determining  $\hat{\rho}_k(\mu)$ . Subsequently, we show how to obtain  $\hat{\rho}_k(\mu)$  for some utility functions.

1) *Maximize the WSR:* As detailed in Section II, for this scenario,  $f_k(\rho_k) = \omega_k \log(1 + \rho_k)$ . The closed-form expression for  $\hat{\rho}_k(\mu)$  can be derived from Theorem 2, as stated in the following corollary.

**Corollary 2.** *The function  $\hat{\rho}_k(\mu)$  for the WSR problem is*

$$\hat{\rho}_k^{\text{WSR}}(\mu) = (3\mu A_k)^{-2} x_k^2, \quad (27)$$

where  $x_k = \max_{j=1,2,3} \{ \mu B_k - \xi^j D_k - \Omega_k \xi^{-j} D_k^{-1} \in \mathbb{R} \}$ ,  $\xi = (-1 + \sqrt{3}j)/2$ ,  $\Omega_k = \mu^2 B_k^2 - 3\mu A_k(\mu A_k - \omega_k)$ ,  $\Delta_k = -2\mu^3 B_k^3 + 9\mu^2 A_k B_k(\mu A_k - \omega_k) - 27\mu^3 A_k^2 B_k$  and  $D_k = \sqrt[3]{\left( \Delta_k + \sqrt{\Delta_k^2 - 4\Omega_k^3} \right) / 2}, \forall k$ .

2) *Maximize the WSR with High SINR:* In this scenario,  $f_k(\rho_k) = \omega_k \log(\rho_k) \approx \omega_k \log(1 + \rho_k)$ , which is equivalent to minimizing the weighted geometric mean of user SINRs (WGMS) [30] and serves as a suitable approximation for the WSR problem when  $P_{\text{sum}}$  is large. Compared to the WSR problem,  $\hat{\rho}_k(\mu)$  for the WGMS problem admits a simpler closed-form solution.

**Corollary 3.** *The function  $\hat{\rho}_k(\mu)$  for the WGMS problem is*

$$\hat{\rho}_k^{\text{WGMS}}(\mu) = \left( \frac{-\mu B_k + \sqrt{\mu^2 B_k^2 + 4\mu A_k \omega_k}}{2\omega_k} \right)^{-2}. \quad (28)$$

$$\hat{P}(\boldsymbol{\rho}, \mathbf{Z}) = \sum_{k=1}^K A_k(\mathbf{Z}) \rho_k - 2B_k(\mathbf{Z}) \sqrt{\rho_k} + C_k(\mathbf{Z}), \quad (21)$$

$$A_k(\mathbf{Z}) = \alpha_k + \sum_{j \neq k} \alpha_j \beta_{jk} z_{jk}^2, \quad B_k(\mathbf{Z}) = \sum_{j \neq k} \alpha_k \beta_{kj} z_{kj}, \quad C_k(\mathbf{Z}) = \sum_{j \neq k} \alpha_j \beta_{jk} z_{jk}^2. \quad (22)$$

3) *Minimize the WMSE*: In this scenario,  $f_k(\rho_k) = -\omega_k(1 + \rho_k)^{-1}$ . Consequently,  $\hat{\rho}_k^{\text{WMSE}}(\mu)$  is the solution of the following equation with variable  $\rho_k$ :

$$-\frac{\omega_k}{(1 + \rho_k)^2} + \mu \left( A_k - \frac{B_k}{\sqrt{\rho_k}} \right) = 0. \quad (29)$$

Due to its monotonicity, the unique root of equation (29) can be determined using the bisection method. The initial range for finding  $\hat{\rho}_k^{\text{WMSE}}(\mu)$  is set as  $\rho_{k,\ell} = B_k^2/A_k^2$  and  $\rho_{k,u} = \max \left\{ 4B_k^2/A_k^2, \sqrt{\omega_k/(2\mu A_k)} - 1 \right\}$ .

4) *Minimize the WSER*: As discussed in Section II, the function  $f_k(\rho_k) = -\omega_k a_k \mathcal{Q}(\sqrt{b_k \rho_k})$  and  $\hat{\rho}_k^{\text{WSER}}(\mu)$  is the solution of the following equation with variable  $\rho_k$ :

$$-\frac{\omega_k a_k \sqrt{b_k}}{2\sqrt{2\pi} \rho_k} e^{-b_k \rho_k/2} + \mu \left( A_k - \frac{B_k}{\sqrt{\rho_k}} \right) = 0. \quad (30)$$

Similar to the WMSE problem,  $\hat{\rho}_k^{\text{WSER}}(\mu)$  can be computed using a bisection method with the initial range set as follows  $\rho_{k,\ell} = B_k^2/A_k^2$  and  $\rho_{k,u} = \max \left\{ 4B_k^2/A_k^2, 2\ln \left( (2\mu\sqrt{2\pi}B_k)^{-1} \omega_k a_k \sqrt{b_k} \right) / b_k \right\}$ .

Based on the consideration above, problem (25) can be optimally solved using the generalized water-filling (GWF) algorithm, as shown in Algorithm 1. Specifically, in some special situations, e.g.,  $f_k(\rho_k) = \log(1 + \rho_k)$  and  $f_k(\rho_k) = \log(\rho_k)$ ,  $\mathcal{P}(\mu)$  can be analytically obtained from Corollary 2 and Corollary 3, respectively.

---

**Algorithm 1:** Generalized Water-Filling Method (GWF) for Problem (25)

---

**Input :**  $\gamma_{\text{th}}, P_{\text{sum}}, \epsilon_\mu$

**Output:**  $\rho^*$

- 1 Initialize  $\mu_\ell, \mu_u$ ;
  - 2 **repeat**
  - 3      $\mu = (\mu_\ell + \mu_u) / 2$ ;
  - 4     Calculate  $\mathcal{P}(\mu)$  according to Lemma 3;
  - 5     **if**  $\mathcal{P}(\mu) < P_{\text{sum}}$  **then**  $\mu_\ell = \mu$  **else**  $\mu_u = \mu$ ;
  - 6 **until**  $\mu_u - \mu_\ell < \epsilon_\mu$ ;
  - 7 Calculate  $\rho^*$  using (26).
- 

C. Constructing  $\mathbf{W}^{\text{SBP}}$  from  $\rho^*$

An essential component of SBP involves constructing  $\mathbf{W}^{\text{SBP}}$  from the obtained  $\rho^*$ . Theoretically, if  $\rho^*$  is the optimal solution to problem (5), then the optimal precoding matrix to problem (4) is  $\mathbf{W}^*(\rho^*)$ , which can be computed using the FPI method described in [2]. However, the traditional FPI method is faced with two significant limitations: a slow convergence rate and substantial computational complexity. To tackle these disadvantages, we propose a modified fixed-point iteration (MFPI) method, as summarized in Algorithm 2.

In Algorithm 2, we modify the update rule of the traditional FPI, which typically sets  $\boldsymbol{\lambda}^{(n+1)} = \hat{\boldsymbol{\lambda}}^{(n+1)}$  [2], with an update rule  $\boldsymbol{\lambda}^{(n+1)} = (\boldsymbol{\rho} + \mathbf{1}_{K \times 1}) \odot \hat{\boldsymbol{\lambda}}^{(n+1)} - \boldsymbol{\rho} \odot \boldsymbol{\lambda}^{(n)}$  to accelerate convergence. The analysis for the convergence performance of Algorithm 2 is presented in the following result.

---

**Algorithm 2:** Modified Fixed-Point Iteration (MFPI) for Problem (7)

---

**Input :**  $\boldsymbol{\rho}, \mathbf{H}$

**Output:**  $\mathbf{W}^*(\boldsymbol{\rho})$

- 1 Initialize  $\boldsymbol{\lambda}^{(0)}$  and  $n = 0$ ;
  - 2 **repeat**
  - 3      $\mathbf{Q}^{(n)} = (\text{Diag}(\boldsymbol{\lambda}^{(n)}) + \bar{\mathbf{H}}^{-1})^{-1}$ ;
  - 4      $\hat{\boldsymbol{\lambda}}^{(n+1)} = (\mathbf{I} + \text{Diag}(\boldsymbol{\rho})^{-1})^{-1} \text{diag}(\mathbf{Q}^{(n)})$ ;
  - 5      $\boldsymbol{\lambda}^{(n+1)} = (\boldsymbol{\rho} + \mathbf{1}_{K \times 1}) \odot \hat{\boldsymbol{\lambda}}^{(n+1)} - \boldsymbol{\rho} \odot \boldsymbol{\lambda}^{(n)}$ ;
  - 6      $n = n + 1$ ;
  - 7 **until**  $\boldsymbol{\lambda}^{(n)}$  converges;
  - 8 Calculate  $\mathbf{W}^*(\boldsymbol{\rho})$  using Proposition 1.
- 

**Theorem 3.** Let  $\{\boldsymbol{\lambda}^{(n)}\}_{n=1}^\infty$  be the sequence generated by Algorithm 2. Then,  $\boldsymbol{\lambda}^{(n)}$  will converge to the optimal solution from any initial point  $\boldsymbol{\lambda}^{(0)} \in \mathbb{R}^{K \times 1}$ . Moreover, the following limits of convergence rate hold

$$\lim_{\boldsymbol{\rho} \rightarrow \infty} r_\lambda(\boldsymbol{\rho}) = 0, \quad \lim_{\boldsymbol{\rho} \rightarrow \infty} r_{\hat{\boldsymbol{\lambda}}}(\boldsymbol{\rho}) = 1, \quad (31)$$

where

$$r_x(\boldsymbol{\rho}) \triangleq \limsup_{n \rightarrow \infty} \frac{\|\mathbf{x}^{(n+1)} - \boldsymbol{\lambda}^*(\boldsymbol{\rho})\|_2^2}{\|\boldsymbol{\lambda}^{(n)} - \boldsymbol{\lambda}^*(\boldsymbol{\rho})\|_2^2}. \quad (32)$$

*Proof.* See Appendix H. □

*Remark 2:* Notice that Theorem 3 reveals that MFPI achieves a better convergence rate when  $\boldsymbol{\rho}$  is large. Furthermore, the computational complexity per iteration for MFPI is  $O(K^3)$ , no more than  $O(N^3)$  for the FPI method in [2, eq. (57)] and less than  $O(KN^3)$  for the method in [23, eq. (30)].

Since  $\rho^*$  is computed based on the approximate feasible SINR region  $\tilde{\mathcal{R}}$  rather than the true feasible SINR region  $\mathcal{R}$ , a further refinement of  $\mathbf{W}^*(\rho^*)$  is required. We propose a practical scaling method given by

$$\mathbf{W}^{\text{SBP}} = \mathbf{W}^*(\rho^*) \text{Diag}(\boldsymbol{\delta})^{\frac{1}{2}}, \quad (33)$$

where  $\boldsymbol{\delta} = [\delta_1, \dots, \delta_K]$  is the scaling factor chosen to satisfy the full power constraint and maintain the QoS equalities for weak users.

More specifically, let  $\mathcal{K}_u \triangleq \{k | \rho_k^* = \gamma_{k,\text{th}}\}$  be the set of weak users. We scale the  $\mathbf{w}_k^*(\rho^*)$  with a different scaling factor  $\sqrt{\delta_k}$  for  $k \in \mathcal{K}_u$  and by a common scaling factor  $\sqrt{\delta_0}$  for  $k \notin \mathcal{K}_u$ . Then,  $\bar{\boldsymbol{\delta}} = [\boldsymbol{\delta}^T, \delta_0]^T$  can be obtained by solving the following linear equations:

$$\begin{cases} \sum_{j=1}^K d_{kj} \delta_k = 1, & \forall k \in \mathcal{K}_u \\ \delta_k = \delta_0, & \forall k \notin \mathcal{K}_u, \\ \sum_{k=1}^K P_k \delta_k = P_{\text{sum}}, \end{cases} \quad (34)$$

where  $d_{kk} = \rho_k^{-1} |\mathbf{h}_k^H \mathbf{w}_k^*(\rho^*)|^2, \forall k \in \mathcal{K}_u, d_{kj} = -|\mathbf{h}_k^H \mathbf{w}_j^*(\rho^*)|^2, \forall k \neq j$  and  $P_k = \|\mathbf{w}_k^*(\rho^*)\|_2^2$ .

### D. Summary and Analysis of SBP

The complete procedure for the SBP framework is summarized in Algorithm 3.

---

#### Algorithm 3: SINR-Based Precoding (SBP) for Problem (4)

---

**Input :**  $\rho_{\text{th}}, P_{\text{sum}}, \mathbf{H}$   
**Output:**  $\mathbf{W}^{\text{SBP}}$

- 1 Initialize  $\rho^{(0)}$  and  $n = 0$ ;
- 2 **repeat**
- 3     Update  $\mathbf{Z}^{(n)}$  using (23);
- 4     Update  $\rho^{(n+1)}$  using Algorithm 1;
- 5      $n = n + 1$ ;
- 6 **until**  $\rho^{(n)}$  converges;
- 7 Calculate  $\mathbf{W}^*$  ( $\rho^{(n)}$ ) using Algorithm 2;
- 8 Calculate  $\delta$  by solving (34);
- 9 Calculate  $\mathbf{W}^{\text{SBP}}$  using (33).

---

The analysis of computational complexity and convergence properties is presented as follows:

1) *Computation Complexity:* Each iteration of the Algorithm 3 has a complexity of  $O(K^2 + K \log_2(\epsilon_\mu^{-1}))$ , stemming from the updates of  $\mathbf{Z}^{(n)}$   $O(K^2)$  and  $\rho^{(n)}$   $O(K \log_2(\epsilon_\mu^{-1}))$ . This complexity is substantially lower than that of WMMSE  $O(N^3)$  [10], RWMMSE  $O(K^3)$  [11], and QoS-constrained WMMSE  $O(T_i \log_2(\epsilon_\mu^{-1})(NK^2 + K^3))$  with  $T_i$  being the number of inner iterations [22].

Consequently, the total complexity of the SBP method is  $O(NK^2 + T_{\text{SBP}}(K^2 + K \log_2(\epsilon_\mu^{-1})) + T_{\text{MFPI}}K^3)$ , which arises from determining  $\{\alpha, \beta\}$ , iteratively calculating the optimized  $\rho^*$  and constructing the precoding matrix  $\mathbf{W}^{\text{SBP}}$  with  $T_{\text{SBP}}$  and  $T_{\text{MFPI}}$  representing the number of iterations for calculating  $\rho^*$  and  $\mathbf{W}^{\text{SBP}}$ , respectively. The total complexity of WMMSE method and RWMMSE method is  $O(T_{\text{W}}N^3)$  and  $O(NK^2 + T_{\text{R}}K^3)$ , where  $T_{\text{W}}$  and  $T_{\text{R}}$  represent the number of iterations for WMMSE and RWMMSE, respectively. The overall complexity of QoS-constrained WMMSE [25] is  $O(T_{\text{Q}}(NK^2 + T_i \log_2(\epsilon_\mu^{-1})(NK^2 + K^3)))$  with  $T_{\text{Q}}$  being the number of outer iterations. In contrast, as the low-complexity baseline, the total computational complexity of ZF precoding is  $O(NK^2 + K \log_2(\epsilon_\mu^{-1}))$ .

*Remark 3:* Notably, the average number of iterations  $T_{\text{SBP}}$  and  $T_{\text{MFPI}}$  are typically 10 and 5, respectively, while the WMMSE-based algorithms often require over 100 iterations to converge. To further demonstrate their convergence characteristics, we have plotted the performance against the average iteration number and running time in Section VI.

2) *Convergence Issue:* Note that Algorithm 3 addresses the original problem (4) by first determining the optimal  $\rho^*$  for problem (18) and subsequently constructing  $\mathbf{W}^{\text{SBP}}$  from  $\rho^*$ . Thus, the convergence lies in finding  $\rho^*$ .

**Theorem 4.** Let  $\{\rho^{(n)}\}_{n=1}^\infty$  be the sequence generated by Algorithm 3. Then,  $\rho^{(n)}$  will converge to a KKT-stationary point of problem (18) from any initial point  $\rho^{(0)} \in \tilde{\mathcal{R}}$  if each  $f_k(\rho_k)$  is strictly concave with  $\rho_k$  and  $\max_k [\sum_{j \neq k}^K \frac{\beta_{kj}}{1+\gamma_{j,\text{th}}}] < 1$ .

*Proof.* See Appendix I.  $\square$

*Remark 4:* Note that the utility functions discussed in Section II, such as maximizing WSR and minimizing WSER, are all strictly concave with  $\rho$ . Moreover, the condition  $\max_k [\sum_{j \neq k}^K \frac{\beta_{kj}}{1+\gamma_{j,\text{th}}}] < 1$  is generally satisfied. Otherwise, we can set  $\beta' = l\beta$  such that  $\max_k [\sum_{j \neq k}^K \frac{l\beta_{kj}}{1+\gamma_{j,\text{th}}}] < 1$ , where the resulting performance loss is negligible.

## V. FURTHER EXTENSIONS

The SBP framework developed above is based on the broadcast channel model (2). In this section, we extend the SBP framework to the interference channels as well as broadcast and interference channels with multi-antenna users under prefixed receivers.

### A. Interference Channels

In multi-cell or cell-free wireless networks, each user may not be served by all the antennas [10], [23]. For example, in a cell-free network, each user typically connects to only a few access points (APs), meaning that only the antennas on those relevant APs transmit signals to the user [23]. In this situation, the received signal at user  $k$  can be modeled as

$$y_k = \mathbf{h}_k^H \sum_{j=1}^K \mathbf{D}_j \mathbf{w}_j s_j + n_k, \quad \forall k, \quad (35)$$

where  $\mathbf{D}_k \in \mathbb{R}^N$  is a diagonal matrix with  $[\mathbf{D}_k]_{nn} = 1$  if user  $k$  is served by antenna  $n$  and  $[\mathbf{D}_k]_{nn} = 0$  otherwise. Notice that, broadcast channel model (2) cannot apply to (35) unless  $\mathbf{D}_1 = \mathbf{D}_2 = \dots = \mathbf{D}_K$ , i.e., all the users are served by the same antennas. To extend the SBP framework to the multi-cell transmission, we will introduce the interference channel model and provide a generalized result in this situation.

Consider a downlink multi-user system with  $K$  transmitter-receiver pairs, where each transmitter is equipped with  $N$  antennas and each receiver is equipped with a single antenna. Different from the broadcast channel model (2), the received signal of receiver  $k$  in the interference channel is

$$y_k = \mathbf{h}_{kk}^H \mathbf{w}_k s_k + \sum_{j \neq k}^K \mathbf{h}_{kj}^H \mathbf{w}_j s_j + n_k, \quad \forall k, \quad (36)$$

where  $\mathbf{h}_{kj} \in \mathbb{C}^{N \times 1}$  is the channel matrix between the  $k$ -th transmitter and the  $j$ -th receiver [35]. Specifically, when  $\mathbf{h}_{k1} = \mathbf{h}_{k2} = \dots = \mathbf{h}_{kK}, \forall k$ , the interference channel model (36) degrades into the previous broadcast channel model (2). Besides, by letting  $\mathbf{h}_{kj} = \mathbf{D}_j \mathbf{h}_k$ , equation (36) can cover the multi-cell or cell-free scenarios (35).<sup>4</sup>

Let  $P^{\text{IC}}(\rho)$  be the optimal objective value of the MP problem in interference channels:

$$\begin{aligned} & \text{minimize} && \sum_{k=1}^K \|\mathbf{w}_k\|^2 \\ & \{\mathbf{w}_k\}_{k=1}^K && \\ & \text{subject to} && 1 + \sum_{j \neq k}^K |\mathbf{h}_{kj}^H \mathbf{w}_j|^2 \leq \rho_k^{-1} |\mathbf{h}_{kk}^H \mathbf{w}_k|^2, \quad \forall k. \end{aligned} \quad (37)$$

<sup>4</sup>In multi-cell scenario, the total interference at user  $k$  consists of the intra-cell interference  $\sum_{j \neq k}^K \mathbf{h}_k^H \mathbf{D}_k \mathbf{D}_j \mathbf{w}_j s_j$  and inter-cell interference  $\sum_{j \neq k}^K \mathbf{h}_k^H (\mathbf{I} - \mathbf{D}_k) \mathbf{D}_j \mathbf{w}_j s_j$ .

Note that  $\mathbf{h}_{kk} \neq \mathbf{0}$  must hold for all  $k$ , otherwise, problem (37) becomes infeasible. An asymptotically tight approximation of  $P^{\text{IC}}(\boldsymbol{\rho})$  is provided in Theorem 5.

**Theorem 5.** An approximation of  $P^{\text{IC}}(\boldsymbol{\rho})$  is

$$\tilde{P}^{\text{IC}}(\boldsymbol{\rho}) = \sum_{k=1}^K \alpha_k^{\text{IC}} \rho_k \left( 1 - \sum_{j \in \mathcal{J}_k \setminus \{k\}} \frac{\beta_{kj}^{\text{IC}}}{1 + \rho_j} \right), \quad (38)$$

where  $\alpha_k^{\text{IC}} = [\bar{\mathbf{H}}_k^{-1}]_{\mathcal{I}_k^k \mathcal{I}_k^k}$ ,  $\beta_{kj}^{\text{IC}} = (\alpha_k^{\text{IC}} \alpha_j^{\text{IC}})^{-1} |[\bar{\mathbf{H}}_k^{-1}]_{\mathcal{I}_k^k \mathcal{I}_k^j}|^2$ ,  $\bar{\mathbf{H}}_k = (\mathbf{H}_k^{\text{IC}})^H \mathbf{H}_k^{\text{IC}}$ ,  $\mathbf{H}_k^{\text{IC}} = [\mathbf{h}_{\mathcal{I}_k^1, k}, \mathbf{h}_{\mathcal{I}_k^2, k}, \dots, \mathbf{h}_{\mathcal{I}_k^{|\mathcal{J}_k|}, k}]$ ,  $\mathcal{J}_k = \{j | \mathbf{h}_{jk} \neq \mathbf{0}\}$ , and  $\mathcal{I}_k^j$  is the index of  $j \in \mathcal{J}_k$ . Moreover, when  $\boldsymbol{\rho} \rightarrow \infty$ ,  $|\tilde{P}^{\text{IC}}(\boldsymbol{\rho}) - P^{\text{IC}}(\boldsymbol{\rho})| \rightarrow 0$ .

*Proof.* The proof is similar to that of Theorem 1.  $\square$

According to Theorem 5, the SBP framework can be directly applied to the interference channel by substituting  $\alpha_k$  and  $\beta_{kj}$  with  $\alpha_k^{\text{IC}}$  and  $\beta_{kj}^{\text{IC}}$ , respectively.

### B. Multi-Antenna Users with Pre-Fixed Receivers

Consider a downlink multiuser system consisting of a BS equipped with  $N$  antennas that simultaneously serves  $K$  users. In particular, user  $k$  is equipped with  $M_k$  antennas, and the BS transmits  $D_k$  data streams to user  $k$ , where  $1 \leq D_k \leq M_k$ .

Since the SBP framework can be directly extended from the broadcast channel to the interference channel, we focus on the broadcast channel model [27]. In this scenario, the received signal  $\mathbf{y}_k \in \mathbb{C}^{M_k \times 1}$  of user  $k$  is given by

$$\mathbf{y}_k = \mathbf{H}_k^H \left( \sum_{j=1}^K \mathbf{W}_j \mathbf{s}_j \right) + \mathbf{n}_k, \quad (39)$$

where  $\mathbf{s}_k = [s_{k,1}, s_{k,2}, \dots, s_{k,D_k}]^T \in \mathbb{C}^{D_k \times 1}$  is the transmitted symbol of user  $k$ ,  $\mathbf{H}_k \in \mathbb{C}^{N \times M_k}$  is the channel matrix to user  $k$ ,  $\mathbf{W}_k = [\mathbf{w}_{k,1}, \dots, \mathbf{w}_{k,D_k}]$  is the precoding matrix for user  $k$  and  $\mathbf{n}_k$  is the additive Gaussian noise with  $\mathbf{n}_k \sim \mathcal{CN}(\mathbf{0}_{M_k \times 1}, \mathbf{R}_k)$ , where  $\mathbf{R}_k = \mathbb{E}[\mathbf{n}_k \mathbf{n}_k^H] \in \mathbb{C}^{M_k \times M_k}$ .

Next, we consider the two common scenarios in downlink multiuser systems, i.e., each user is equipped with a fixed linear receiver or a fixed SIC receiver.

1) *Linear Receiver:* In this scenario, the recovered  $d$ -th data stream of user  $k$  is

$$\hat{s}_{k,d} = \mathbf{g}_{k,d}^H \left[ \mathbf{H}_k^H \left( \sum_{j=1}^K \mathbf{W}_j \mathbf{s}_j \right) + \mathbf{n}_k \right], \quad (40)$$

where  $\mathbf{g}_{k,d} \in \mathbb{C}^{M_k \times 1}$  is the receiver coefficient of the  $d$ -th data stream for user  $k$ . The SINR of the  $d$ -th stream of user  $k$  is

$$\rho_{k,d} \left( \{\mathbf{W}_k\}_{k=1}^K \right) = \frac{|\mathbf{g}_{k,d}^H \mathbf{H}_k^H \mathbf{w}_{k,d}|^2}{I_{k,d}^{\text{LIN}} + \mathbf{g}_{k,d}^H \mathbf{R}_k \mathbf{g}_{k,d}}, \quad (41)$$

where the interference power  $I_{k,d}^{\text{LIN}}$  is given by

$$I_{k,d}^{\text{LIN}} = \sum_{j \neq k} \sum_{m=1}^{D_j} |\mathbf{g}_{k,d}^H \mathbf{H}_k^H \mathbf{w}_{j,m}|^2 + \sum_{m \neq d}^{D_k} |\mathbf{g}_{k,d}^H \mathbf{H}_k^H \mathbf{w}_{k,m}|^2.$$

Denote the SINR of the  $d$ -th data stream of user  $k$  as  $\rho_{k,d}$ . Then, the minimum power required for the targeted SINRs  $\boldsymbol{\rho} \triangleq [\rho_{1,1}, \dots, \rho_{1,D_1}, \dots, \rho_{K,D_K}]^T$  can be approximated according to Corollary 1.

2) *SIC Receiver:* Without loss of generality, let the decoding order of the data streams of user  $k$  be  $s_{k,1}, s_{k,2}, \dots, s_{k,D_k}$ . Then, the  $d$ -th recovered data stream for user  $k$  is

$$\hat{s}_{k,d} = \mathbf{g}_{k,d}^H \left[ \mathbf{H}_k^H \left( \sum_{j \neq k} \mathbf{W}_j \mathbf{s}_j + \sum_{m \geq d}^{D_k} \mathbf{w}_{k,d} s_{k,m} \right) + \mathbf{n}_k \right], \quad (42)$$

where  $\mathbf{g}_{k,d} \in \mathbb{C}^{M_k \times 1}$  is the receiver coefficient of the  $d$ -th data stream of user  $k$ . The SINR of the  $d$ -th data stream of user  $k$  is

$$\rho_{k,d} \left( \{\mathbf{W}_k\}_{k=1}^K \right) = \frac{|\mathbf{g}_{k,d}^H \mathbf{H}_k^H \mathbf{w}_{k,d}|^2}{I_{k,d}^{\text{SIC}} + \mathbf{g}_{k,d}^H \mathbf{R}_k \mathbf{g}_{k,d}}, \quad (43)$$

where the interference power  $I_{k,d}^{\text{SIC}}$  is

$$I_{k,d}^{\text{SIC}} = \sum_{j \neq k} \sum_{m=1}^{D_j} |\mathbf{g}_{k,d}^H \mathbf{H}_k^H \mathbf{w}_{j,m}|^2 + \sum_{m > d}^{D_k} |\mathbf{g}_{k,d}^H \mathbf{H}_k^H \mathbf{w}_{k,m}|^2.$$

Similar to the linear receiver scenario, the minimum power required to achieve the target SINRs with fixed SIC receivers can be approximated according to Theorem 5.

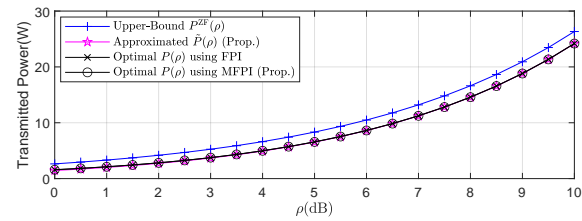
## VI. NUMERICAL ANALYSIS

In this section, we numerically evaluate the accuracy of the proposed approximation and verify the performance of the proposed SBP framework under different criteria. In the simulation, we employ QuaDRiGa [36]–[38] to generate the channel matrix  $\mathbf{H}$ , where the “3GPP 38.901 UMa NLOS” scenario is considered. The BS is located at 25 m high and consists of a uniform planar array (UPA) with  $N = N_v \times N_h$  antennas, where  $N_v$  and  $N_h$  represent the number of the antennas at each vertical column and horizontal row, respectively. Users are randomly located in the cell with a radius of 250 m. Unless otherwise specified, we set  $N = 4 \times 8$ ,  $K = 8$ ,  $\gamma_{k,\text{th}} = 1, \forall k$ , and each user is equipped with a single antenna.

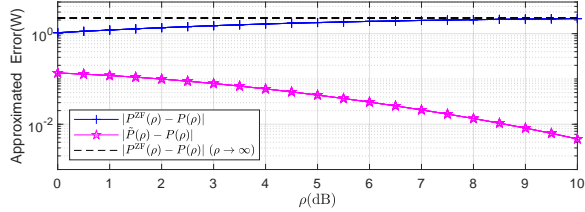
### A. Approximation of $\mathcal{R}$

In this subsection, we examine the accuracy of the proposed approximate feasible SINR region  $\tilde{\mathcal{R}}$ , which primarily depends on the approximation performance of the proposed  $\tilde{P}(\boldsymbol{\rho})$ . Following a similar approach to [11], we set the noise power to be equal for all users, given as  $\sigma_k^2 = 10^{\frac{1}{K} \sum_{k=1}^K \log_{10} \|\mathbf{h}_k\|_F^2} \times 10^{-\frac{\text{SNR}}{10}}$ , where SNR is the average received signal-to-noise ratio without precoding, with a default value of 15dB.

First, we validate the accuracy of our proposed approximation,  $\tilde{P}(\boldsymbol{\rho})$ , and verify our theoretical analysis of ZF precoding in Remark 1. We plot our approximation and the ZF solution against the true optimal value,  $P(\boldsymbol{\rho})$ , obtained from MFPI or FPI [2] with  $\boldsymbol{\rho} = \rho \mathbf{1}_{K \times 1}$ . From Fig. 1a, it can be observed that  $\tilde{P}(\boldsymbol{\rho})$  provides a precise approximation to  $P(\boldsymbol{\rho})$ . Furthermore, Fig. 1b shows that the gap between  $P(\boldsymbol{\rho})$  and  $P^{\text{ZF}}(\boldsymbol{\rho})$  tends to a constant  $C$  while the gap between  $P(\boldsymbol{\rho})$  and  $\tilde{P}(\boldsymbol{\rho})$  vanishes, where  $C = \sum_{k=1}^K \sum_{j \neq k} \alpha_k \beta_{kj}$  according to Corollary 1. This

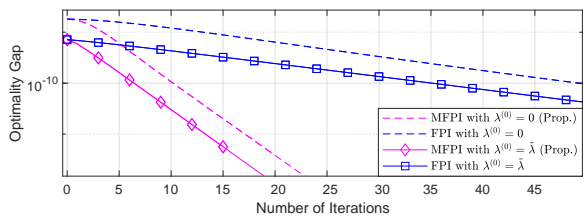


(a) Approximate value under different  $\rho$

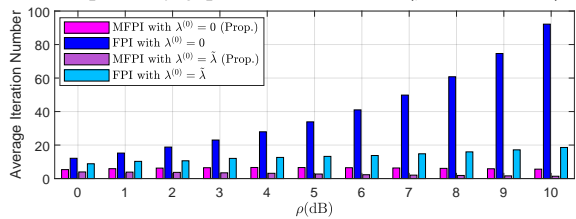


(b) Approximation absolute error under different  $\rho$

Fig. 1: Approximation performance of  $\tilde{P}(\rho)$ .



(a) Optimality gap of MFPI and FPI ( $\rho_k = 5\text{dB}, \forall k$ )



(b) Average iteration number under different  $\rho$

Fig. 2: Convergence performance of MFPI and FPI

indicates that ZF precoding is not asymptotically optimal for the power minimization problem.

Next, we compare the convergence performance of the proposed MFPI and FPI. In Fig. 2a, we compare the optimality gap of problem (7) using MFPI and FPI from different initial points, i.e.,  $\lambda^{(0)} = \mathbf{0}_{K \times 1}$  or  $\tilde{\lambda}(\rho)$ . It can be seen that MFPI achieves an optimality gap of  $10^{-10}$  in 6 iterations, whereas FPI requires 35 iterations when  $\lambda^{(0)} = \tilde{\lambda}(\rho)$ . In Fig. 2b, we plot the average number of iterations needed to reach an optimality gap of  $10^{-5}$ . As  $\rho$  increases, the average iteration number for FPI rises sharply. In contrast, the average iteration number for MFPI decreases, confirming the validity of Theorem 3. Furthermore, even in low-SINR scenarios, MFPI algorithm converges much faster than FPI.

Then, we examine the gap between the true SINR region  $\mathcal{R}$  defined in (8), the approximate SINR region  $\tilde{\mathcal{R}}$  defined in (16), and ZF precoding SINR region  $\mathcal{R}^{\text{ZF}}$  defined in Remark 1 under different values of  $P_{\text{sum}}$ . We consider a simple  $2 \times 2$  multiuser case, where  $\mathcal{R}$  can be analytically obtained from

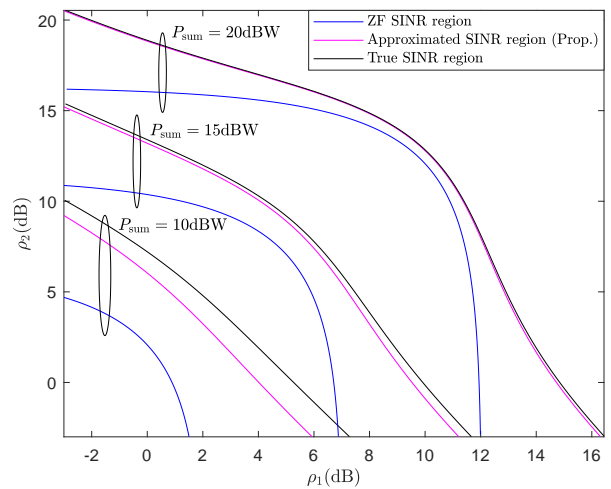
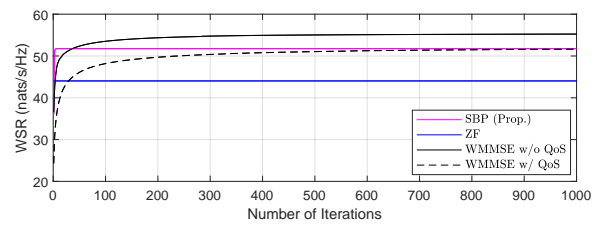
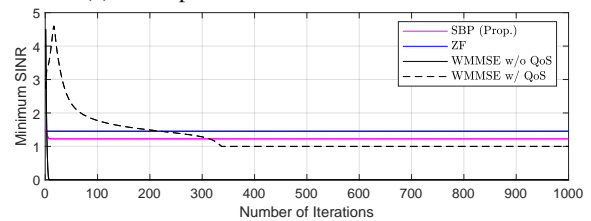


Fig. 3: Boundary of the feasible SINR region



(a) WSR performance vs. iteration number



(b) Minimum SINR vs. iteration number

Fig. 4: Convergence of WSR with different algorithms

Lemma 2. In the simulation, let  $\sigma_k^2 = 1$ ,  $\mathbf{H}$  be

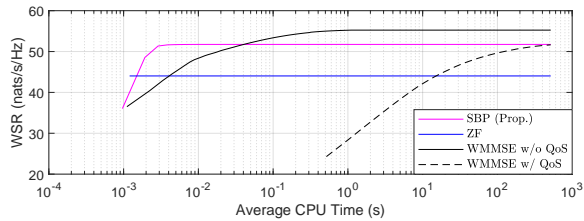
$$\mathbf{H} = \begin{bmatrix} -0.3430 - 0.0168j & -0.7667 + 0.6843j \\ 1.6663 + 0.3042j & 2.6355 + 0.0558j \end{bmatrix},$$

and  $\gamma_{k,\text{th}} = -3\text{dB}, \forall k$ . In Fig. 3, we illustrate the boundaries of  $\mathcal{R}$ ,  $\tilde{\mathcal{R}}$  and  $\mathcal{R}^{\text{ZF}}$ . It is evident that  $\tilde{\mathcal{R}}$  provides the closest approximation to  $\mathcal{R}$ , with the gap diminishing as  $P_{\text{sum}}$  increases.

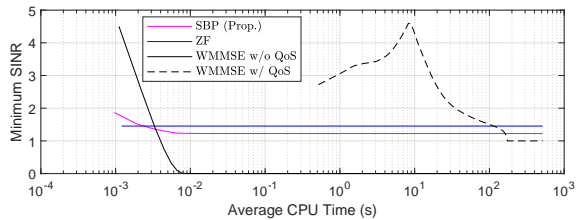
### B. WSR Performance

In this subsection, we provide a comprehensive performance evaluation of the proposed SBP method for WSR maximization with  $P_{\text{sum}} = 10\text{W}$ .

Fig. 4 compares the convergence performance of the proposed SBP and WMMSE-based methods with outer iterations. In particular, we refer to WMMSE w/o QoS as the method proposed in [9], and WMMSE w/ QoS as the method proposed in [22], respectively. We also compare with the ZF

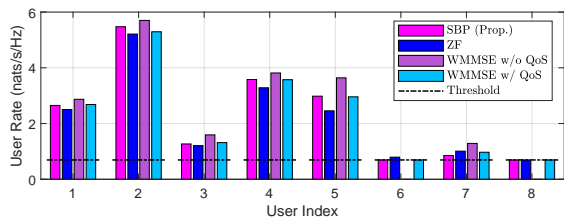


(a) WSR performance vs. CPU time

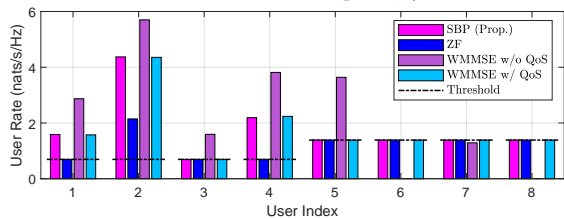


(b) Minimum SINR vs. CPU time

Fig. 5: Convergence of WSR with different algorithms



(a) User rate with equal  $\gamma_{k,\text{th}}$

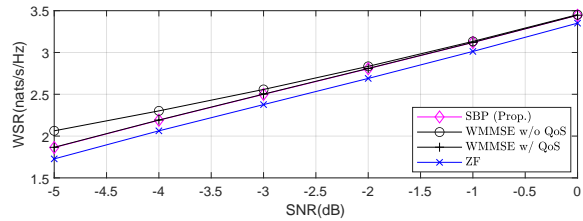


(b) User rate with unequal  $\gamma_{k,\text{th}}$

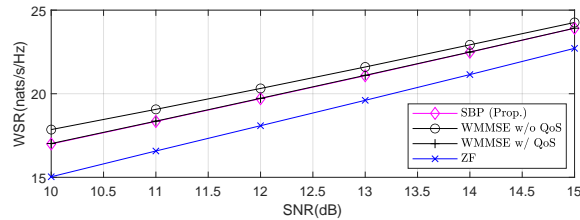
Fig. 6: User rate with different SINR threshold

method for WSR maximization [3]. For this simulation, we set  $\text{SNR} = 25\text{dB}$ ,  $\omega_k = 1, \forall k$ ,  $K = N/2 = 16$ , and use the ZF precoding matrix  $\mathbf{W}^{\text{ZF}}$  satisfying  $\gamma(\mathbf{W}^{\text{ZF}}) = \gamma_{\text{th}}$  for initialization. From Fig. 4a, it can be observed that SBP requires a few iterations to converge, whereas WMMSE-based methods require nearly hundreds of iterations. As shown in Fig. 4b, it is evident that SBP and WMMSE w/ QoS can guarantee the QoS requirements, while WMMSE w/o QoS fails to meet these constraints, with the minimum user SINR dropping close to zero. Although WMMSE w/o QoS achieves a slightly higher WSR, this is an expected outcome because it does not consider QoS constraints, thereby expanding its feasible solution space.

Notice that the average number of outer iterations may not reflect the actual convergence speed of different algorithms, since the computational complexity of each iteration is also critical. To provide a more practical comparison, we plot the WSR performance of the algorithms as a function of CPU time in Fig. 5. The results clearly show that the proposed method

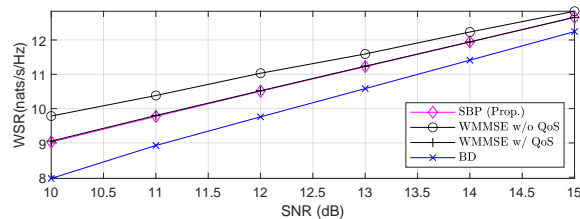


(a)  $N = 8, K = 2$

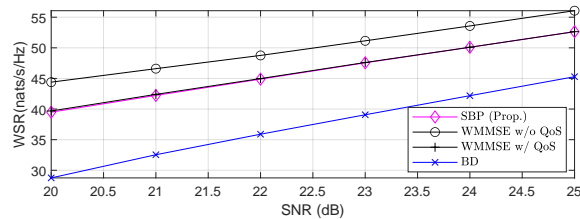


(b)  $N = 32, K = 8$

Fig. 7: WSR performance with single-antenna users



(a)  $N = 8, K = 2$



(b)  $N = 32, K = 8$

Fig. 8: WSR performance with multi-antenna users

exhibits the fastest convergence. In contrast, the WMMSE w/ QoS baseline requires substantially more computational time, since it requires solving a high-dimensional convex optimization problem in each iteration.

Fig. 6 assesses the SBP algorithm's performance across different QoS profiles by plotting the per-user rate  $r_k = \log(1 + \gamma_k(\mathbf{W}))$  for a fixed channel realization. We compare two distinct scenarios: a homogeneous profile with a uniform QoS target ( $\gamma_{k,\text{th}} = 1, \forall k$ ) in Fig. 6a, and a heterogeneous profile with mixed targets in Fig. 6b, with ( $\gamma_{k,\text{th}} = 1, \forall k \leq 4$ ) and ( $\gamma_{k,\text{th}} = 3, \forall 5 \leq k$ ). The proposed SBP algorithm can guarantee all QoS requirements in both settings. The WMMSE w/o QoS method, however, fails to meet these constraints, often sacrificing some users by suppressing their rates to zero.

Fig. 7 evaluates the WSR performance as a function of SNR for different methods with  $\omega_k = 1, \forall k$ . As the figure demonstrates, the proposed SBP method performs almost identically to the WMMSE w/ QoS method. Moreover, SBP and WMMSE-based algorithms exhibit similar WSR perfor-

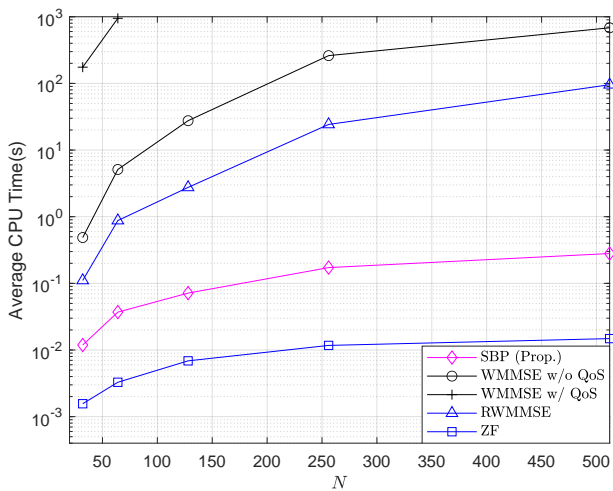
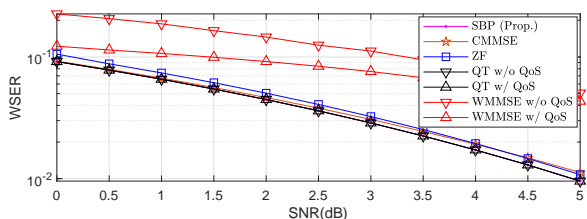
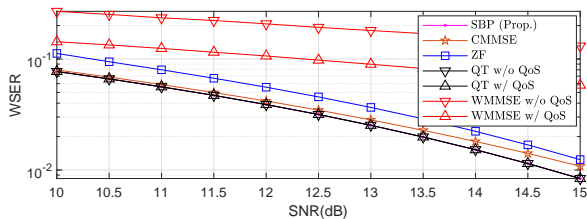


Fig. 9: Average CPU time vs. different  $N$



(a)  $N = 8, K = 2$



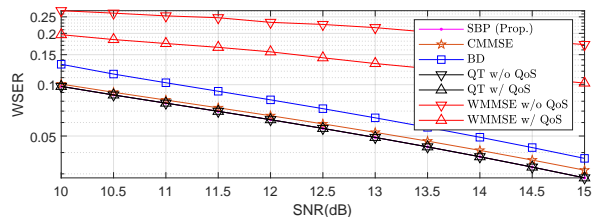
(b)  $N = 32, K = 8$

Fig. 10: WSER performance with single-antenna users

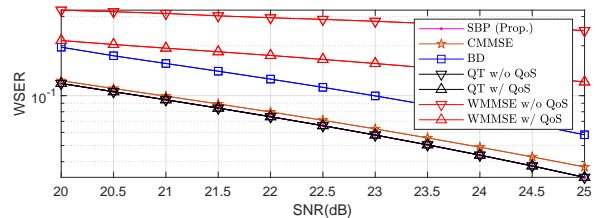
mance when SNR is high, as shown in Fig. 7a and Fig. 7b, respectively.

In Fig. 8, we also investigate the WSR performance with multi-antenna users, where we set  $\omega_k = 1$ ,  $M_k = D_k = 2$ , and the SINR threshold of each data stream of user  $k$  is set as  $\gamma_{k,\text{th}} = 2e^{-1/D_k} - 1$ . We also include the block-diagonalization (BD) method [27] as a low-complexity benchmark. In this scenario, the user antenna type is “3GPP-3D” and the pre-fixed linear receiver  $\mathbf{G}_k = [\mathbf{g}_{k,1}, \mathbf{g}_{k,2}] \in \mathbb{C}^{2 \times 2}$  is the right eigenvector matrix of  $\mathbf{H}_k$  corresponding to the  $D_k$  largest eigenvalues. It can be observed that the proposed SBP method can achieve nearly the same performance as WMMSE w/ QoS method.

Fig. 9 shows the average CPU time of SBP, WMMSE w/o QoS, WMMSE w/ QoS, RWMMSE [11], and ZF for the WSR maximization as a function of the number of BS antennas  $N$ . In the simulation, we set  $N \in \{32, 64, 128, 256, 512\}$  and  $K = N/4$ . The average CPU time is measured using MATLAB on an Intel Core i7-13700K. As illustrated, SBP

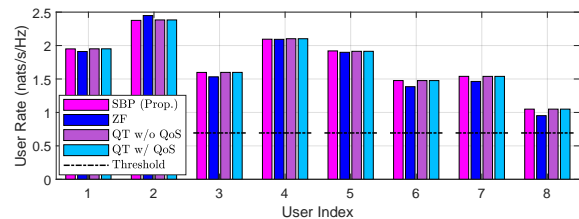


(a)  $N = 8, K = 2$

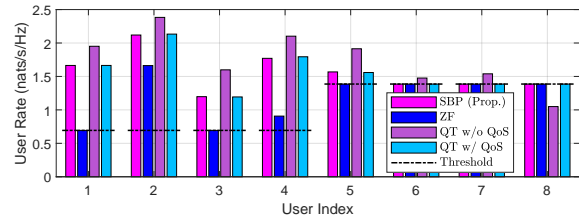


(b)  $N = 32, K = 8$

Fig. 11: WSR performance with multi-antenna users



(a) User rate with equal  $\gamma_{k,\text{th}}$



(b) User rate with unequal  $\gamma_{k,\text{th}}$

Fig. 12: User rate with different SINR threshold

holds a significant computational advantage over conventional WMMSE-based algorithms, since they possess a higher per-iteration complexity and usually require a large number of iterations to reach a solution. Furthermore, WMMSE w/ QoS exhibits the highest complexity, as they need to solve a high-dimensional optimization problem. When compared with the low-complexity ZF precoding, our SBP algorithm achieves a significant performance gain at the cost of an acceptable increase in computational complexity.

### C. WSR Performance

In this subsection, we analyze and present the performance of the SBP method in terms of the WSR metric.

Fig. 10 compares the average WSR performance of several algorithms under a sum power constraint, assuming quadrature phase shift keying (QPSK) modulation and uniform weights  $\omega_k = 1/K$ . The evaluated algorithms include our proposed SBP algorithm, designed to minimize WSR under QoS

constraints, and two methods based on the quadrature transform (QT) from [39]: one that minimizes WSER with QoS constraints (denoted as QT w/ QoS) and another without QoS constraints (denoted as QT w/o QoS). We also consider the constrained MMSE (CMMSE) precoding [17], which approximates the WSER minimization with a WMSE minimization problem, and ZF precoding for WSER minimization with QoS constraints. Besides, Fig. 11 shows the average WSER under multi-antenna users scenario using the same simulation parameters as Fig. 8. As we can see, the proposed SBP achieves nearly identical performance to the QT-based methods.

Fig. 12 plots the per-user rates for the WSER minimization problem under various QoS constraints. It is observed that, unlike the WSR maximization shown in Fig. 6, WSER minimization seldom sacrifices some users to enhance overall system performance, thereby making the WSER performance under QoS constraints similar to that without QoS constraints.

## VII. CONCLUSION

In this paper, we developed an efficient and unified framework for linear precoding design in downlink multiuser systems that accommodates various utility functions while ensuring QoS requirements under a sum power constraint. By analyzing the minimum power problem, we proposed an approximate SINR region that closely estimates the true SINR region, enabling us to transform the high-dimensional precoding design into a more manageable low-dimensional SINR allocation problem. Our SINR-based precoding (SBP) framework converges to a KKT-stationary solution and has been successfully extended to the interference channel and multi-antenna users scenario with pre-fixed receivers, achieving an order-of-magnitude reduction in the computational complexity in each iteration. Extensive simulation results demonstrate that our method achieves near-optimal performance comparable to existing methods while significantly reducing computational complexity.

## APPENDIX

### A. Proof of Lemma 1

Let  $\bar{\mathcal{R}} = \{\boldsymbol{\rho} | \boldsymbol{\rho} \geq \boldsymbol{\gamma}_{\text{th}}, P(\boldsymbol{\rho}) \leq P_{\text{sum}}\}$ . Then, we only need to prove that  $\mathcal{R} = \bar{\mathcal{R}}$ .

Select arbitrary  $\boldsymbol{\gamma}(\mathbf{W}_0) \in \mathcal{R}$ , then  $P(\boldsymbol{\gamma}(\mathbf{W}_0)) \leq \|\mathbf{W}_0\|_F^2 \leq P_{\text{sum}}$  according to the definition of  $P(\boldsymbol{\rho})$ . Since  $\boldsymbol{\gamma}(\mathbf{W}_0) \geq \boldsymbol{\gamma}_{\text{th}}$ , we obtain  $\boldsymbol{\gamma}(\mathbf{W}_0) \in \bar{\mathcal{R}}$ . Combined with the arbitrariness of  $\boldsymbol{\gamma}(\mathbf{W}_0)$ , we have  $\mathcal{R} \subset \bar{\mathcal{R}}$ .

On the other hand, for any  $\boldsymbol{\rho}_0 \in \bar{\mathcal{R}}$ , let  $\mathbf{W}_0^*$  be the optimal solution to problem (7) with QoS threshold  $\boldsymbol{\rho}_0$ . Then,  $\boldsymbol{\gamma}(\mathbf{W}_0^*) = \boldsymbol{\rho}_0$  holds [2]. Besides, we have  $P(\boldsymbol{\rho}_0) = \|\mathbf{W}_0^*\|_F^2 \leq P_{\text{sum}}$  and  $\boldsymbol{\gamma}(\mathbf{W}_0^*) \geq \boldsymbol{\gamma}_{\text{th}}$ . Consequently, we obtain  $\boldsymbol{\rho}_0 \in \mathcal{R}$  and  $\bar{\mathcal{R}} \subset \mathcal{R}$ .

Then, we have  $\mathcal{R} = \bar{\mathcal{R}}$  since  $\bar{\mathcal{R}} \subset \mathcal{R}$  and  $\mathcal{R} \subset \bar{\mathcal{R}}$ , which completes the proof of Lemma 1.

### B. Proof of Proposition 1

Denote  $\boldsymbol{\lambda} \triangleq [\lambda_1, \lambda_2, \dots, \lambda_K]^T$  and update function  $\mathbf{U}(\boldsymbol{\lambda}) \triangleq [U_1(\boldsymbol{\lambda}), U_2(\boldsymbol{\lambda}), \dots, U_K(\boldsymbol{\lambda})]^T$ , where  $U_k(\boldsymbol{\lambda})$  is

$$U_k(\boldsymbol{\lambda}) \triangleq \frac{\rho_k}{1 + \rho_k} \left[ \mathbf{h}_k^H (\mathbf{I} + \mathbf{H}\boldsymbol{\Lambda}\mathbf{H}^H)^{-1} \mathbf{h}_k \right]^{-1}. \quad (44)$$

According to [2],  $\boldsymbol{\lambda}^*(\boldsymbol{\rho})$  is the unique fixed-point of  $\mathbf{U}(\boldsymbol{\lambda})$ . Moreover,  $\mathbf{h}_k^H (\mathbf{I} + \mathbf{H}\boldsymbol{\Lambda}\mathbf{H}^H)^{-1} \mathbf{h}_k$  can be simplified into  $[(\bar{\mathbf{H}}^{-1} + \boldsymbol{\Lambda})^{-1}]_{kk}$ , according to the equation  $(\mathbf{I} + \mathbf{A}\mathbf{B})^{-1} \mathbf{A} = \mathbf{A} (\mathbf{I} + \mathbf{B}\mathbf{A})^{-1}$ . As a result, we have

$$\lambda_k^*(\boldsymbol{\rho}) \left[ (\bar{\mathbf{H}}^{-1} + \boldsymbol{\Lambda}^*)^{-1} \right]_{kk} = \frac{\rho_k}{1 + \rho_k}, \forall k, \quad (45)$$

which is equivalent to  $\boldsymbol{\Lambda}^*(\boldsymbol{\rho}) \odot (\boldsymbol{\Lambda}^*(\boldsymbol{\rho}) + \bar{\mathbf{H}}^{-1})^{-1} = (\mathbf{I} + \boldsymbol{\Gamma}^{-1})^{-1}$ . Moreover, the expression of  $\mathbf{W}^*(\boldsymbol{\rho})$  can be derived similarly to the method in [2], and thus the proof of Proposition 1 is completed.

### C. Proof of Lemma 2

When  $K = 2$ ,  $(\boldsymbol{\Lambda}^* + \bar{\mathbf{H}}^{-1})^{-1}$  possesses a closed-form expression, which is given by

$$\begin{aligned} (\boldsymbol{\Lambda}^* + \bar{\mathbf{H}}^{-1})^{-1} &= \frac{1}{(\lambda_1^* + h_{11})(\lambda_2^* + h_{22}) - |h_{21}|^2} \\ &\times \begin{bmatrix} \lambda_2^* + h_{22} & -h_{12} \\ -h_{21} & \lambda_1^* + h_{11} \end{bmatrix}. \end{aligned} \quad (46)$$

For simplicity, let  $\lambda_1^*(\boldsymbol{\rho})$  and  $\lambda_2^*(\boldsymbol{\rho})$  be  $\lambda_1^*$  and  $\lambda_2^*$ , respectively. According to Proposition 1,  $\lambda_1^*, \lambda_2^*$  is the positive solution of the following equations

$$\begin{cases} \lambda_1^* \lambda_2^* + h_{22} \lambda_1^* - \rho_1 h_{11} \lambda_2^* = \rho_1 (h_{11} h_{22} - |h_{21}|^2) \\ \lambda_1^* \lambda_2^* + h_{11} \lambda_2^* - \rho_2 h_{22} \lambda_1^* = \rho_2 (h_{11} h_{22} - |h_{21}|^2) \end{cases}. \quad (47)$$

Subtracting the second equation from the first equation in the (47), we obtain

$$\lambda_1^* = \frac{h_{11}(1 + \rho_1)}{h_{22}(1 + \rho_2)} \lambda_2^* + \frac{(\rho_1 - \rho_2)(h_{11} h_{22} - |h_{21}|^2)}{h_{22}(1 + \rho_2)}. \quad (48)$$

Substituting (48) back into (47), we find that  $\lambda_2^*$  is the positive root of a second-order polynomial, which can be derived analytically as in (12). Similarly,  $\lambda_1^*$  can be derived. The proof of Lemma 2 is completed.

### D. Proof of Theorem 1

**Lemma 4** ([40]). *Let  $\mathbf{P} \in \mathbb{C}^{N \times N}$  and suppose that  $\lim_{n \rightarrow \infty} \mathbf{P}^n = \mathbf{0}_{N \times N}$ . Then  $\mathbf{I} - \mathbf{P}$  is non-singular and*

$$(\mathbf{I} - \mathbf{P})^{-1} = \sum_{n=0}^{\infty} \mathbf{P}^n \quad (49)$$

We denote  $\boldsymbol{\lambda}^*(\boldsymbol{\rho})$  as  $\boldsymbol{\lambda}^*$  for simplicity. According to Proposition 1, the optimal solution to problem (9) satisfies

$$\begin{aligned} (1 + \rho_k^{-1})^{-1} &= \lambda_k^* \left[ (\bar{\boldsymbol{\Lambda}}^* - \bar{\mathbf{D}})^{-1} \right]_{kk} \\ &= \frac{\lambda_k^*}{\lambda_k^* + [\bar{\mathbf{H}}^{-1}]_{kk}} \left[ (\mathbf{I} - \bar{\mathbf{D}}\bar{\boldsymbol{\Lambda}}^{-1})^{-1} \right]_{kk}, \end{aligned} \quad (50)$$

where  $\bar{\mathbf{D}} = \bar{\mathbf{H}}^{-1} \odot \mathbf{I}_{K \times K}$ ,  $\bar{\boldsymbol{\Lambda}} = \boldsymbol{\Lambda}^*(\boldsymbol{\rho}) + \bar{\mathbf{D}}$ ,  $\bar{\mathbf{D}} = \bar{\mathbf{D}} - \bar{\mathbf{H}}^{-1}$ . Rearranging (50), we obtain

$$\frac{\lambda_k^*}{[\bar{\mathbf{H}}^{-1}]_{kk} \rho_k} = \frac{1}{\left[ (\mathbf{I} - \bar{\mathbf{D}}\bar{\boldsymbol{\Lambda}}^{-1})^{-1} \right]_{kk} (\rho_k + 1) - \rho_k}. \quad (51)$$

From the KKT condition of problem (7),  $\lambda_k^* \geq [\bar{\mathbf{H}}]_{kk}^{-1} \rho_k$ . Then  $\lim_{\rho \rightarrow \infty} \bar{\mathbf{\Lambda}}^{-1} = \mathbf{0}_{K \times K}$  according to  $0 < [\bar{\mathbf{\Lambda}}^{-1}]_{kk} \leq (\rho_k [\bar{\mathbf{H}}]_{kk}^{-1} + [\bar{\mathbf{H}}^{-1}]_{kk})^{-1}$ . As a result, when  $\rho$  is large enough, the maximum singular value  $\sigma_{\max}(\bar{\mathbf{D}}\bar{\mathbf{\Lambda}}^{-1}) < 1$ , which indicates that  $\lim_{n \rightarrow \infty} (\bar{\mathbf{D}}\bar{\mathbf{\Lambda}}^{-1})^n = \mathbf{0}_{K \times K}$ . So

$$\begin{aligned} & \lim_{\rho \rightarrow \infty} [(\mathbf{I} - \bar{\mathbf{D}}\bar{\mathbf{\Lambda}}^{-1})^{-1}]_{kk} (\rho_k + 1) - \rho_k \\ & \stackrel{(a)}{=} \lim_{\rho \rightarrow \infty} 1 + (\rho_k + 1) \sum_{n=2}^{\infty} [(\bar{\mathbf{D}}\bar{\mathbf{\Lambda}}^{-1})^n]_{kk} \stackrel{(b)}{=} 1, \end{aligned}$$

where (a) holds due to  $[\bar{\mathbf{D}}\bar{\mathbf{\Lambda}}^{-1}]_{kk} = 0, \forall k$  and (b) holds due to  $|\sum_{n=2}^{\infty} [(\bar{\mathbf{D}}\bar{\mathbf{\Lambda}}^{-1})^n]_{kk}| \leq \sum_{n=2}^{\infty} |[(\bar{\mathbf{D}}\bar{\mathbf{\Lambda}}^{-1})^n]_{kk}| \leq \sum_{n=2}^{\infty} \|\bar{\mathbf{D}}\|_F^n \max_{k=1}^K \|h_k\|^{2n} \rho_k^{-n} = o(\rho^{-1})^5$ . As a result, we have  $\lim_{\rho \rightarrow \infty} \frac{\lambda_k^*}{[\bar{\mathbf{H}}^{-1}]_{kk} \rho_k} = 1, \forall k$ .

Applying Lemma 4 again, we have

$$\begin{aligned} & 1 + (\rho_k + 1) \sum_{n=2}^{\infty} [(\bar{\mathbf{D}}\bar{\mathbf{\Lambda}}^{-1})^n]_{kk} \\ & \stackrel{(c)}{=} 1 + (\rho_k + 1) [\bar{\mathbf{D}}\bar{\mathbf{\Lambda}}^{-1} \bar{\mathbf{D}}\bar{\mathbf{\Lambda}}^{-1}]_{kk} + o(\rho^{-1}) \\ & = 1 + \frac{\rho_k + 1}{\lambda_k^* + [\bar{\mathbf{H}}^{-1}]_{kk}} \sum_{j=1}^K \frac{|\bar{\mathbf{D}}|_{kj}^2}{\lambda_j^* + [\bar{\mathbf{H}}^{-1}]_{jj}} + o(\rho^{-1}) \\ & \stackrel{(d)}{=} 1 + \frac{\rho_k + 1}{(\rho_k + 1) [\bar{\mathbf{H}}^{-1}]_{kk}} \sum_{j=1}^K \frac{|\bar{\mathbf{D}}|_{kj}^2}{(\rho_j + 1) [\bar{\mathbf{H}}^{-1}]_{jj}} + o(\rho^{-1}) \\ & \stackrel{(e)}{=} 1 + \sum_{j \neq k}^K \frac{\beta_{kj}}{\rho_j + 1} + o(\rho^{-1}), \end{aligned} \quad (52)$$

where (c) holds because  $|\sum_{n=3}^{\infty} [(\bar{\mathbf{D}}\bar{\mathbf{\Lambda}}^{-1})^n]_{kk}| \leq \sum_{n=3}^{\infty} |[(\bar{\mathbf{D}}\bar{\mathbf{\Lambda}}^{-1})^n]_{kk}| \leq \sum_{k=3}^{\infty} \|\bar{\mathbf{D}}\|_F^n \max_{k=1}^K \|h_k\|^{2n} \rho_k^{-n} = o(\rho^{-2})$ , (d) holds due to  $(\lambda_k^* + \alpha_k)^{-1} - (\alpha_k + \alpha_k \rho_k)^{-1} = (\alpha_k \rho_k - \lambda_k^*) / (\rho_k \lambda_k^* + \alpha_k \lambda_k^* + \alpha_k^2 \rho_k + \alpha_k^2) = o(\rho^{-1})$  and (e) holds due to  $[\bar{\mathbf{D}}]_{kk} = 0, \forall k$ . According to (51), we have

$$\lambda_k^* = \frac{\alpha_k \rho_k}{1 + \sum_{j \neq k}^K \frac{\beta_{kj}}{\rho_j + 1} + o(\rho^{-1})} = [\tilde{\lambda}]_k + o(1). \quad (53)$$

The last equation holds due to  $(1+x)^{-1} = 1-x+o(x)$ . The proof of Theorem 1 is completed.

### E. Proof of Proposition 2

First,  $\tilde{P}(\rho)$  can be rewritten as

$$\begin{aligned} & \tilde{P}(\rho) \\ & \stackrel{(a)}{=} \min_{\mathbf{Z}} \sum_{k=1}^K \alpha_k \rho_k + \sum_{k=1}^K \sum_{j \neq k}^K \alpha_k \beta_{kj} ((\rho_j + 1) z_{kj}^2 - 2\sqrt{\rho_k} z_{kj}) \\ & = \min_{\mathbf{Z}} \sum_{k=1}^K (A_k(\mathbf{Z}) \rho_k - 2B_k(\mathbf{Z}) \sqrt{\rho_k} + C_k(\mathbf{Z})) \\ & = \min_{\mathbf{Z}} \hat{P}(\rho, \mathbf{Z}) \end{aligned} \quad (54)$$

<sup>5</sup>For any function  $g(\rho)$ ,  $o(g(\rho))$  is defined as  $\lim_{\rho \rightarrow \infty} \frac{o(g(\rho))}{g(\rho)} = 0$ . For any vector-valued function  $\mathbf{g}(\rho)$ , we define  $o(\mathbf{g}(\rho)) \triangleq o(\max_k \{g_k(\rho)\})$ .

where  $\mathbf{Z} = (z_{kj})_{K \times K} \in \mathbb{R}^{K \times K}$  and (a) holds according to the following transformation [39]:

$$\frac{M(\mathbf{x})}{N(\mathbf{x})} = \max_{z \in \mathbb{R}} 2\sqrt{M(\mathbf{x})}z - N(\mathbf{x})z^2. \quad (55)$$

As a result,  $\tilde{P}(\rho) \leq \hat{P}(\rho, \mathbf{Z})$  for any  $\mathbf{Z} \in \mathbb{R}^{K \times K}$  and equation  $\tilde{P}(\rho) = \hat{P}(\rho, \mathbf{Z})$  holds when  $z_{ij} = \sqrt{\rho_i} / (\rho_j + 1), \forall j \neq i$ . Then, an equivalent expression of  $\tilde{\mathcal{R}}$  is

$$\begin{aligned} \tilde{\mathcal{R}} & = \left\{ \rho \mid \rho \geq \rho_{\text{th}}, \tilde{P}(\rho) \leq P_{\text{sum}} \right\} \\ & = \left\{ \rho \mid \rho \geq \rho_{\text{th}}, \min_{\mathbf{Z}} \hat{P}(\rho, \mathbf{Z}) \leq P_{\text{sum}} \right\} = \bigcup_{\mathbf{Z} \in \mathbb{R}^{K \times K}} \hat{\mathcal{R}}(\mathbf{Z}). \end{aligned} \quad (56)$$

This implies that the optimal  $\rho^*$  to problem (18) is also the optimal solution to problem (19). To further investigate the relationship between the KKT-stationary point of problems (18) and (19), we consider the KKT condition of problems (18) and (19). Denote the Lagrangian function of problem (18) and (19) be  $L_0(\rho, \mu, \mathbf{v}) = -f(\rho) + \mu(\tilde{P}(\rho) - P_{\text{sum}}) + \mathbf{v}^T(\rho - \gamma_{\text{th}})$  and  $L_1(\rho, \mathbf{Z}, \mu, \mathbf{v}) = -f(\rho) + \mu(\hat{P}(\rho, \mathbf{Z}) - P_{\text{sum}}) + \mathbf{v}^T(\rho - \gamma_{\text{th}})$ , respectively. Notice that when  $z_{kj} = \sqrt{\rho_k} / (\rho_j + 1)$ , we have

$$\frac{\partial \hat{P}(\rho, \mathbf{Z})}{\partial \rho_k} = \frac{\partial \tilde{P}(\rho)}{\partial \rho_k}, \quad \hat{P}(\rho, \mathbf{Z}) = \tilde{P}(\rho).$$

Let  $(\bar{\rho}, \bar{\mathbf{Z}}, \bar{\mu}, \bar{\mathbf{v}})$  be a KKT point of problem (19) satisfies  $\bar{z}_{ij} = \sqrt{\bar{\rho}_i} / (\bar{\rho}_j + 1)$ , then  $(\bar{\rho}, \bar{\mu}, \bar{\mathbf{v}})$  satisfies the KKT condition of problem (18) since  $\frac{\partial L_0(\bar{\rho}, \bar{\mu}, \bar{\mathbf{v}})}{\partial \rho_k} = \frac{\partial L_1(\bar{\rho}, \bar{\mathbf{Z}}, \bar{\mu}, \bar{\mathbf{v}})}{\partial \rho_k} = 0$  and  $\bar{\mu}(\tilde{P}(\bar{\rho}) - P_{\text{sum}}) = \bar{\mu}(\hat{P}(\bar{\rho}, \bar{\mathbf{Z}}) - P_{\text{sum}}) = 0$ . The proof of Proposition 2 is completed.

### F. Proof of Theorem 2

The KKT condition of problem (25) is given as follows:

$$-f'_k(\rho_k^*) + \mu^* p'_k(\rho_k^*) - v_k^* = 0, \quad \forall k, \quad (57)$$

$$\mu^* \left( \sum_{k=1}^K p_k(\rho_k^*) - P_{\text{sum}} \right) = 0, \quad (58)$$

$$v_k^* (\rho_k^* - \gamma_{k,\text{th}}) = 0, \quad \forall k, \quad (59)$$

$$\rho_k^* - \rho_{k,\text{th}} \geq 0, \quad \forall k, \quad (60)$$

$$\sum_{k=1}^K p_k(\rho_k^*) - P_{\text{sum}} \leq 0, \quad (61)$$

$$\mu^* \geq 0, \quad (62)$$

$$v_k^* \geq 0, \quad \forall k, \quad (63)$$

where  $\mu^*$  and  $v_k^*$  is the dual variable of  $\sum_{k=1}^K p_k(\rho_k) \leq P_{\text{sum}}$  and  $\rho_k \geq \rho_{k,\text{th}}$ , respectively.

We first prove  $\sum_{k=1}^K p_k(\rho_k^*) = P_{\text{sum}}$  by contradiction. If  $\sum_{k=1}^K p_k(\rho_k^*) < P_{\text{sum}}$ , according to (58), we have  $\mu^* = 0$ . Replacing  $\mu^*$  with 0 in (57), we have  $f'_k(\rho_k^*) = -v_k^* \leq 0$ , which contradicts the fact that  $f_k(\rho_k)$  is an increasing function. If  $\rho_k^* > \rho_{k,\text{th}}$ , then  $v_k^* = 0$ . Otherwise,  $v_k^* = -f'_k(\rho_{k,\text{th}}) + \mu^* p'_k(\rho_{k,\text{th}})$ . Define  $h_k(\rho_k) = -f'_k(\rho_k) + \mu^* p'_k(\rho_k)$ , then  $h_k(\rho_k)$  is monotonically increasing over  $\rho_k$ . In this case,  $\hat{\rho}_k(\mu^*) = h^{-1}(0) \leq h^{-1}(v_k^*) = \rho_{k,\text{th}}$ . Then the  $\rho_k^*$  can be summarized as  $\rho_k^* = \max\{\hat{\rho}_k(\mu^*), \gamma_{k,\text{th}}\}$ . Then the proof of Theorem 2 is completed.

### G. Proof of Lemma 3

Observing that  $\hat{\rho}_k(\mu)$  is a decreasing function over  $\mu$  is equivalent to for any  $\mu_2 > \mu_1 > 0$ ,  $\hat{\rho}_k(\mu_1) > \hat{\rho}_k(\mu_2)$  holds. Combined with  $p'_k(\hat{\rho}_k(\mu_1)) > 0$  and  $p_k(\hat{\rho}_k(\mu_2)) > 0$ , we have

$$\begin{aligned} 0 &= -f'_k(\hat{\rho}_k(\mu_2)) + \mu_2 p'_k(\hat{\rho}_k(\mu_2)) \\ &= -f'_k(\hat{\rho}_k(\mu_1)) + \mu_1 p'_k(\hat{\rho}_k(\mu_1)) \\ &< -f'_k(\hat{\rho}_k(\mu_1)) + \mu_2 p'_k(\hat{\rho}_k(\mu_1)). \end{aligned} \quad (64)$$

As function  $-f'_k(\rho_k) + \mu_2 p'_k(\rho_k)$  is monotonically increasing over  $\rho_k$ , we have  $\hat{\rho}_k(\mu_1) > \hat{\rho}_k(\mu_2)$ . Let  $\hat{\rho}_k(\mu) = \max\{\gamma_{k,\text{th}}, \hat{\rho}_k(\mu)\}$ , so  $\hat{\rho}_k(\mu)$  is non-increasing over  $\mu$ . Although we do not require that  $p_k(\rho_k)$  is monotonically increasing for all  $\rho_k \geq \gamma_{k,\text{th}}$ ,  $p_k(\hat{\rho}_k(\mu))$  is still non-decreasing over  $\mu > 0$  since  $p'_k(\hat{\rho}_k(\mu)) > 0$  for any  $\mu > 0$ . As a result,  $\mathcal{P}(\mu) = \sum_{k=1}^K p_k(\max\{\gamma_{k,\text{th}}, \hat{\rho}_k(\mu)\})$  is a non-increasing function of  $\mu$ . The proof of Lemma 3 is completed.

### H. Proof of Theorem 3

1) *Convergence of  $\{\lambda^{(n)}\}_{n=1}^\infty$* : The update rule of  $\lambda^{(n+1)}$  is given as  $\lambda^{(n+1)} = \mathbf{V}(\lambda^{(n)}) = (\boldsymbol{\rho} + \mathbf{1}) \odot \mathbf{U}(\lambda^{(n)}) - \boldsymbol{\rho} \odot \lambda^{(n)}$ , where  $\mathbf{U}(\lambda)$  is defined in Appendix B. Notice that the  $k$ -th element of  $\mathbf{V}(\lambda)$  can be equivalently transformed as

$$V_k(\lambda) = \rho_k [\mathbf{h}_k^H (\mathbf{I} + \mathbf{H}_{[k]}^H \boldsymbol{\Lambda}_{[k]} \mathbf{H}_{[k]})^{-1} \mathbf{h}_k]^{-1}, \quad (65)$$

where  $\mathbf{H}_{[k]} \triangleq [\mathbf{h}_1, \dots, \mathbf{h}_{k-1}, \mathbf{h}_{k+1}, \dots, \mathbf{h}_K]$  and  $\boldsymbol{\Lambda}_{[k]} = \text{Diag}(\lambda_1, \dots, \lambda_{k-1}, \lambda_{k+1}, \dots, \lambda_K)$ . It is easy to check that  $\mathbf{V}(\lambda)$  is a standard function in [2], thus  $\{\lambda^{(n)}\}_{n=1}^\infty$  can converge to the optimal solution from any initial point.

2) *The Analysis of Convergence Rate for  $\{\lambda^{(n)}\}_{n=1}^\infty$  and  $\{\hat{\lambda}^{(n)}\}_{n=1}^\infty$* : For simplicity, we denote  $\lambda^*(\boldsymbol{\rho})$  as  $\lambda^*$ . According to the Lagrange's mean value theorem, we have

$$\|\hat{\lambda}^{(n+1)} - \lambda^*\|_2^2 = \|\mathbf{F}(\lambda^{(n)}) (\lambda^{(n)} - \lambda^*)\|_2^2, \quad (66)$$

where  $[\mathbf{F}(\lambda^{(n)})]_{ij} \triangleq \partial U_i(\lambda) / \partial \lambda_j |_{\lambda = \eta_{ij}(\lambda^{(n)})}$ , where  $\eta_{ij}(\lambda^{(n)})$  is some point between  $\lambda^*$  and  $\lambda^{(n)}$ . As a result,

$$\sigma_{\min}(\mathbf{F}(\lambda^{(n)})) \leq \frac{\|\hat{\lambda}^{(n+1)} - \lambda^*\|_2}{\|\lambda^{(n)} - \lambda^*\|_2} \leq \sigma_{\max}(\mathbf{F}(\lambda^{(n)})). \quad (67)$$

Specifically, we define  $[\mathbf{F}^*]_{ij} \triangleq \partial U_i(\lambda^*) / \partial \lambda_j = \frac{\rho_i}{\rho_i + 1} \frac{[\mathbf{G}]_{ij}}{[\mathbf{G}]_{ii}}$ , where  $\mathbf{G}$  is defined in Proposition 1. To investigate the behavior of  $\frac{[\mathbf{G}]_{ij}}{[\mathbf{G}]_{ii}}$  when  $\boldsymbol{\rho} \rightarrow \infty$ , we define  $\mathbf{B} \triangleq (\mathbf{A}^* + \bar{\mathbf{H}}^{-1})^{-1}$  and then  $[\mathbf{G}]_{ij} = |[\mathbf{B}]_{ij}|^2$ . According to Lemma 4,

$$\mathbf{B} = (\mathbf{A}^*)^{-1} - (\mathbf{A}^*)^{-1} \bar{\mathbf{H}}^{-1} (\mathbf{A}^*)^{-1} + o(\boldsymbol{\rho}^{-2}), \quad (68)$$

then  $\frac{[\mathbf{B}]_{ij}}{[\mathbf{B}]_{ii}} \rightarrow \delta_{ij}$ , where  $\delta_{ij} = 1, \forall i = j$  and  $\delta_{ij} = 0, \forall i \neq j$ . As a result,  $\frac{[\mathbf{G}]_{ij}}{[\mathbf{G}]_{ii}} \rightarrow \delta_{ij}$ . Let  $\sigma_{\max}(\mathbf{A})$ ,  $\sigma_{\min}(\mathbf{A})$  denote the maximum singular value and the minimum singular value of  $\mathbf{A}$ , respectively. Then, an upper bound and lower bound of convergence rate  $r_{\hat{\lambda}}(\boldsymbol{\rho})$  is  $\sigma_{\max}(\mathbf{F}^*)$  and  $\sigma_{\min}(\mathbf{F}^*)$ , respectively. It is easy to check that

$$\lim_{\boldsymbol{\rho} \rightarrow \infty} \sigma_{\min}(\mathbf{F}^*) = \lim_{\boldsymbol{\rho} \rightarrow \infty} \sigma_{\max}(\mathbf{F}^*) = 1. \quad (69)$$

Similarly, in the case of MFPI, the upper and lower bound of  $r_{\lambda}(\boldsymbol{\rho})$  is  $\sigma_{\max}(\mathbf{M}^*)$  and  $\sigma_{\min}(\mathbf{M}^*)$ , where  $\mathbf{M}^* = (\text{Diag}(\boldsymbol{\rho}) + \mathbf{I}) \mathbf{F}^* - \text{Diag}(\boldsymbol{\rho})$ . Applying (68), we have

$$\lim_{\boldsymbol{\rho} \rightarrow \infty} \sigma_{\min}(\mathbf{M}^*) = \lim_{\boldsymbol{\rho} \rightarrow \infty} \sigma_{\max}(\mathbf{M}^*) = 0. \quad (70)$$

The proof of Theorem 3 is completed.

### I. Proof of Theorem 4

**Lemma 5** ([41]). *Let  $\phi: \mathcal{X} \rightarrow \mathcal{Y}$  be a continuous correspondence between topological spaces with nonempty compact values, and suppose  $F: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  is continuous. Define the value function  $m(\mathbf{x})$  by*

$$m(\mathbf{x}) = \max_{\mathbf{y} \in \phi(\mathbf{x})} F(\mathbf{x}, \mathbf{y}),$$

and the correspondence of maximisers by

$$\eta(\mathbf{x}) = \{\mathbf{Z} \in \phi(\mathbf{x}) | F(\mathbf{x}, \mathbf{Z}) = m(\mathbf{x})\}.$$

Then, the value function is continuous and the correspondence  $\eta(\mathbf{x})$  is upper hemicontinuous and has nonempty compact values.

First, we consider the trivial scenario that  $\tilde{\mathcal{R}}$  is a singleton, i.e.,  $\tilde{\mathcal{R}} = \{\boldsymbol{\rho}_{\text{th}}\}$ , which is equivalent to  $\tilde{P}(\boldsymbol{\rho}_{\text{th}}) = P_{\text{sum}}$  when  $\boldsymbol{\rho}_{\text{th}}$  satisfies  $\max_k \left[ \sum_{j \neq k} \frac{\beta_{kj}}{1 + \rho_j} \right] < 1$ . In this case,  $\{\boldsymbol{\rho}^{(n)}\}_{n=1}^\infty$  generated from Algorithm 3 is actually  $\boldsymbol{\rho}^{(n)} = \boldsymbol{\rho}_{\text{th}}$ , which is a trivial stationary point of problem (18) in this case.

Consider the non-trivial scenario that  $\tilde{\mathcal{R}}$  is not a singleton, i.e.  $\tilde{P}(\boldsymbol{\rho}_{\text{th}}) < P_{\text{sum}}$ . In this case,  $\tilde{\mathcal{R}}$  is a compact set since it is closed and bounded. According to the procedure of Algorithm 3,  $\{f(\boldsymbol{\rho}^{(n)})\}_{n=1}^\infty$  is a non-decreasing sequence with non-trivial upper bound  $f^{\text{ub}} = \max_{\boldsymbol{\rho} \in \tilde{\mathcal{R}}} f(\boldsymbol{\rho})$ . So,  $\{f(\boldsymbol{\rho}^{(n)})\}_{n=1}^\infty$  is convergent and let  $\lim_{n \rightarrow \infty} f(\boldsymbol{\rho}^{(n)}) = f_0$ . As a result, there exists a subsequence  $\{\boldsymbol{\rho}^{(n_j)}\}_{j=1}^\infty$  that converges to some  $\bar{\boldsymbol{\rho}} \in \tilde{\mathcal{R}}$ , satisfying  $f(\bar{\boldsymbol{\rho}}) = f_0$ .

Next, we prove that  $\{\boldsymbol{\rho}^{(n_j+1)}\}_{j=1}^\infty$  is convergent, which is mainly based on Lemma 5. Without the need to check the continuity of the correspondence  $\hat{\mathcal{R}}(\mathbf{Z})$  on  $\mathbb{R}^{K \times K}$ , we only need to consider the continuity on  $\mathcal{D}_{\mathbf{Z}}$ , where  $\mathcal{D}_{\mathbf{Z}} \triangleq \bigcup_{\boldsymbol{\rho} \in \tilde{\mathcal{R}}} \{\mathbf{Z} | [\mathbf{Z}]_{kj} = \sqrt{\bar{\rho}_k} / (1 + \rho_j), \forall k \neq j\}$  is a compact set since it is the continuous map from the compact set  $\tilde{\mathcal{R}}$  to  $\mathbb{R}^{K \times K}$ . The continuity of  $\hat{\mathcal{R}}(\mathbf{Z})$  consists of the upper-hemicontinuity (UHC) and lower-hemicontinuity (LHC) on  $\mathcal{D}_{\mathbf{Z}}$ . The UHC of  $\hat{\mathcal{R}}(\mathbf{Z})$  follows from the fact that the graph of  $\hat{\mathcal{R}}(\mathbf{Z})$  is closed and locally bounded on  $\mathcal{D}_{\mathbf{Z}}$ . In terms of LHC, for arbitrary point  $\mathbf{Z}_0 \in \mathcal{D}_{\mathbf{Z}}$ , assume  $\mathcal{G}$  is an open set satisfies  $\mathcal{G} \cap \hat{\mathcal{R}}(\mathbf{Z}_0) \neq \emptyset$ . Select any  $\boldsymbol{\rho}_0 \in \hat{\mathcal{R}}(\mathbf{Z}_0) \cap \mathcal{G}$ . If  $\tilde{P}(\boldsymbol{\rho}_0) < P_{\text{sum}}$ , then choose  $\mathcal{U} = \{\mathbf{Z} | \hat{P}(\boldsymbol{\rho}_0, \mathbf{Z}) < P_{\text{sum}}\}$ . Otherwise, select  $\boldsymbol{\rho}_1 \in \mathcal{G}$  satisfies  $\hat{P}(\boldsymbol{\rho}_1, \mathbf{Z}_0) < P_{\text{sum}}$  and  $\boldsymbol{\rho}_0 \neq \boldsymbol{\rho}_1 \geq \gamma_{\text{th}}$  and set  $\mathcal{U} = \{\mathbf{Z} | \hat{P}(\boldsymbol{\rho}_1, \mathbf{Z}) < P_{\text{sum}}\}$ . Note that there exists  $\boldsymbol{\rho}_1$  satisfies condition above unless  $\boldsymbol{\rho}_0 = \gamma_{\text{th}}$  and  $\tilde{P}(\boldsymbol{\rho}_0) = P_{\text{sum}}$ , which conflicts with the assumption that  $\tilde{P}(\boldsymbol{\rho}_0) < P_{\text{sum}}$ . Since  $\mathcal{U}$  includes an open neighborhood of  $\mathbf{Z}_0$  and for every  $\mathbf{Z} \in \mathcal{U}$ , we have  $\boldsymbol{\rho}_1 \in \hat{\mathcal{R}}(\mathbf{Z}) \cap \mathcal{G} \neq \emptyset$ , then  $\hat{\mathcal{R}}(\mathbf{Z}_0)$  is LHC at  $\mathbf{Z}_0$ . Combined with the arbitrariness of  $\mathbf{Z}_0$ ,  $\hat{\mathcal{R}}(\mathbf{Z}_0)$  is LHC over  $\mathcal{D}_{\mathbf{Z}}$ . As a result, by applying Lemma 5 the

correspondence of maximizers is hemicontinuous. Moreover, since  $f(\boldsymbol{\rho})$  is a strictly concave function, the correspondence is a singleton. Then we obtain the following limits:

$$\lim_{j \rightarrow \infty} \arg \max_{\boldsymbol{\rho} \in \hat{\mathcal{R}}(\mathbf{Z}^{(n_j)})} f(\boldsymbol{\rho}) = \arg \max_{\boldsymbol{\rho} \in \hat{\mathcal{R}}(\bar{\mathbf{Z}})} f(\boldsymbol{\rho}), \quad (71)$$

where  $[\bar{\mathbf{Z}}]_{ij} = \bar{\rho}_i^{\frac{1}{2}} (\bar{\rho}_j + 1)^{-1}$ . This limit reveals the convergence of the sequence  $\{\boldsymbol{\rho}^{(n_j+1)}\}_{j=1}^{\infty}$ . Denote the limit of  $\{\boldsymbol{\rho}^{(n_j+1)}\}_{j=1}^{\infty}$  is  $\bar{\boldsymbol{\rho}}$ . Then, we will show that sequences  $\{\boldsymbol{\rho}^{(n_j)}\}_{j=1}^{\infty}$  and  $\{\boldsymbol{\rho}^{(n_j+1)}\}_{n=1}^{\infty}$  converge to the same point.

Since  $P_{\text{sum}} \geq \hat{P}(\boldsymbol{\rho}^{(n_j+1)}, \mathbf{Z}^{(n_j)})$  holds for all  $j$  when  $\boldsymbol{\rho}^{(0)} \in \hat{\mathcal{R}}$  and when  $j$  tends to  $\infty$ , we have  $P_{\text{sum}} \geq \hat{P}(\bar{\boldsymbol{\rho}}, \bar{\mathbf{Z}})$ . As a result,  $\bar{\boldsymbol{\rho}} \in \hat{\mathcal{R}}(\bar{\mathbf{Z}})$ . As  $\{f(\boldsymbol{\rho}^{(n)})\}_{n=1}^{\infty}$  converges to  $f_0$ ,  $\lim_{j \rightarrow \infty} f(\boldsymbol{\rho}^{(n_j+1)}) = f(\bar{\boldsymbol{\rho}}) = f_0$ , which means  $\bar{\boldsymbol{\rho}}$  is the solution of the following optimization problem:

$$\begin{aligned} & \underset{\boldsymbol{\rho}}{\text{maximize}} && f(\boldsymbol{\rho}) \\ & \text{subject to} && \boldsymbol{\rho} \in \hat{\mathcal{R}}(\bar{\mathbf{Z}}). \end{aligned} \quad (72)$$

According to the condition that the objective function  $f(\boldsymbol{\rho})$  is a strictly concave function, then the solution of (72) is unique, i.e.,  $\bar{\boldsymbol{\rho}} = \bar{\boldsymbol{\rho}}$ . As a result, sequence  $\{\boldsymbol{\rho}^{(n_j+k)}\}_{j=1}^{\infty}$  will converge to  $\bar{\boldsymbol{\rho}}$  for any  $k \geq 0$ , i.e.,  $\{\boldsymbol{\rho}^{(n)}\}_{n=1}^{\infty}$  will converge.

According to the uniqueness of  $\bar{\boldsymbol{\rho}}$ ,  $\bar{\boldsymbol{\rho}}$  satisfies the KKT condition of problem (72). Replacing  $\bar{y}_{kj}$  with  $\frac{\sqrt{\bar{\rho}_k}}{\bar{\rho}_j+1}$ , we can prove that  $\bar{\boldsymbol{\rho}}$  satisfies the KKT condition of problem (18), thus the proof of Theorem 4 is completed.

## REFERENCES

- [1] J. Zhang, E. Björnson, M. Matthaiou, D. W. K. Ng, H. Yang, and D. J. Love, "Prospective multiple antenna technologies for beyond 5G," *IEEE J. Sel. Areas Commun.*, vol. 38, no. 8, pp. 1637–1660, Aug. 2020.
- [2] A. Wiesel, Y. Eldar, and S. Shamai, "Linear precoding via conic optimization for fixed MIMO receivers," *IEEE Trans. Signal Process.*, vol. 54, no. 1, pp. 161–176, Jan. 2006.
- [3] A. Wiesel, Y. C. Eldar, and S. Shamai, "Zero-forcing precoding and generalized inverses," *IEEE Trans. Signal Process.*, vol. 56, no. 9, pp. 4409–4418, Sep. 2008.
- [4] Y. Zhang, P. Mitran, and C. Rosenberg, "Joint resource allocation for linear precoding in downlink massive MIMO systems," *IEEE Trans. Commun.*, vol. 69, no. 5, pp. 3039–3053, May 2021.
- [5] J. Kaleva, A. Tölli, and M. Juntti, "Decentralized sum rate maximization with QoS constraints for interfering broadcast channel via successive convex approximation," *IEEE Trans. Signal Process.*, vol. 64, no. 11, pp. 2788–2802, Feb. 2016.
- [6] L.-N. Tran, M. F. Hanif, A. Tölli, and M. Juntti, "Fast converging algorithm for weighted sum rate maximization in multicell MISO downlink," *IEEE Signal Process. Lett.*, vol. 19, no. 12, pp. 872–875, Oct. 2012.
- [7] H. Kha, H. D. Tuan, and H. H. Nguyen, "Fast global optimal power allocation in wireless networks by local D.C. programming," *IEEE Trans. Wireless Commun.*, vol. 11, no. 2, pp. 510–515, Feb. 2012.
- [8] A. Alvarado, G. Scutari, and J.-S. Pang, "A new decomposition method for multiuser DC-programming and its applications," *IEEE Trans. Signal Process.*, vol. 62, no. 11, pp. 2984–2998, Apr. 2014.
- [9] S. S. Christensen, R. Agarwal, E. De Carvalho, and J. M. Cioffi, "Weighted sum-rate maximization using weighted MMSE for MIMO-BC beamforming design," *IEEE Trans. Wireless Commun.*, vol. 7, no. 12, pp. 4792–4799, Dec. 2008.
- [10] Q. Shi, M. Razaviyayn, Z.-Q. Luo, and C. He, "An iteratively weighted MMSE approach to distributed sum-utility maximization for a MIMO interfering broadcast channel," *IEEE Trans. Signal Process.*, vol. 59, no. 9, pp. 4331–4340, Sep. 2011.
- [11] X. Zhao, S. Lu, Q. Shi, and Z.-Q. Luo, "Rethinking WMMSE: Can its complexity scale linearly with the number of BS antennas?" *IEEE Trans. Signal Process.*, vol. 71, pp. 433–446, Feb. 2023.
- [12] T. Ma, Q. Shi, and E. Song, "QoS-constrained weighted sum-rate maximization in multi-cell multi-user MIMO systems: An ADMM approach," in *Proc. IEEE 2016 35th Chin. Control Conf. (CCC)*, Chengdu, China, Jul. 2016, pp. 6905–6910.
- [13] R. Hunger, W. Utschick, D. Schmidt, and M. Joham, "Alternating optimization for MMSE broadcast precoding," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, vol. 4, Toulouse, France, May 2006, pp. IV–IV.
- [14] C. W. Tan, M. Chiang, and R. Srikant, "Maximizing sum rate and minimizing MSE on multiuser downlink: Optimality, fast algorithms and equivalence via max-min SINR," *IEEE Trans. Signal Process.*, vol. 59, no. 12, pp. 6127–6143, Dec. 2011.
- [15] D. Palomar, "Convex primal decomposition for multicarrier linear MIMO transceivers," *IEEE Trans. Signal Process.*, vol. 53, no. 12, pp. 4661–4674, Dec. 2005.
- [16] D. Palomar, M. Bengtsson, and B. Ottersten, "Minimum BER linear transceivers for MIMO channels via primal decomposition," *IEEE Trans. Signal Process.*, vol. 53, no. 8, pp. 2866–2882, Aug. 2005.
- [17] J. Li, Y. Jiang, and L. Zhao, "Hybrid precoding design in multi-user mmWave massive MIMO systems for BER minimization," *IEEE Wireless Commun. Lett.*, vol. 13, no. 1, pp. 208–212, Jan. 2024.
- [18] M. Bengtsson and B. Ottersten, "Optimal and suboptimal transmit beamforming," in *Handbook of Antennas in Wireless Communications*, L. C. Godara, Ed. CRC Press, Aug. 2001.
- [19] E. Björnson, M. Bengtsson, and B. Ottersten, "Optimal multiuser transmit beamforming: A difficult problem with a simple solution structure [lecture notes]," *IEEE Signal Process. Mag.*, vol. 31, no. 4, pp. 142–148, Jul. 2014.
- [20] M. Schubert and H. Boche, "Solution of the multiuser downlink beamforming problem with individual SINR constraints," *IEEE Trans. Veh. Technol.*, vol. 53, no. 1, pp. 18–28, Jan. 2004.
- [21] Z.-Q. Luo and S. Zhang, "Dynamic spectrum management: Complexity and duality," *IEEE J. Sel. Top. Signal Process.*, vol. 2, no. 1, pp. 57–73, Feb. 2008.
- [22] Y. Huang, K. Chi, Q. Yang, Z. Yang, and Z. Zhang, "Soft actor-critic-based multi-user multi-tti MIMO precoding in multi-modal real-time broadband communications," *IEEE Trans. Wireless Commun.*, vol. 23, no. 12, pp. 18 286–18 301, Dec. 2024.
- [23] L. Han, J. Wang, R. Hou, S. He, D. W. K. Ng, L. Xia, and Q. Wang, "Resource efficient beamforming design for cell-free networks," *IEEE Trans. Commun.*, vol. 72, no. 12, pp. 7511–7525, Dec. 2024.
- [24] K. Chi, Y. Huang, Q. Yang, Z. Yang, and Z. Zhang, "MIMO precoding design with QoS and per-antenna power constraints," in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, Kuala Lumpur, Malaysia, Dec. 2023, pp. 3324–3329.
- [25] Z. Cao, W. Wang, and D. Qiao, "Statistical QoS-aware precoding design for downlink multiuser MIMO systems with per-antenna power constraints," *IEEE Wireless Commun. Letters*, vol. 13, no. 12, pp. 3628–3632, Dec. 2024.
- [26] Y. Qi and M. Vaezi, "Signaling design for MIMO-NOMA with different security requirements," *IEEE Trans. Signal Process.*, vol. 70, pp. 1389–1401, Mar. 2022.
- [27] Q. Spencer, A. Swindlehurst, and M. Haardt, "Zero-forcing methods for downlink spatial multiplexing in multiuser MIMO channels," *IEEE Trans. Signal Process.*, vol. 52, no. 2, pp. 461–471, Feb. 2004.
- [28] S. A. A. Fakoorian and A. L. Swindlehurst, "On the optimality of linear precoding for secrecy in the MIMO broadcast channel," *IEEE J. Sel. Areas Commun.*, vol. 31, no. 9, pp. 1701–1713, Sep. 2013.
- [29] J. Tong, P. J. Schreier, and S. R. Weller, "Linear precoding for MIMO systems with low-complexity receivers," *IEEE Trans. Wireless Commun.*, vol. 11, no. 8, pp. 2828–2837, Aug. 2012.
- [30] D. Palomar, J. Cioffi, and M. Lagunas, "Joint Tx-Rx beamforming design for multicarrier MIMO channels: a unified framework for convex optimization," *IEEE Trans. Signal Process.*, vol. 51, no. 9, pp. 2381–2401, Sep. 2003.
- [31] H. Boche and S. Stanczak, "Convexity of some feasible QoS regions and asymptotic behavior of the minimum total power in CDMA systems," *IEEE Trans. Commun.*, vol. 52, no. 12, pp. 2190–2197, Dec. 2004.
- [32] I. Stewart, *Galois theory*. Chapman and Hall/CRC, 2022.
- [33] L. Sanguinetti, A. L. Moustakas, E. Björnson, and M. Debbah, "Large system analysis of the energy consumption distribution in multi-user MIMO systems with mobility," *IEEE Trans. Wireless Commun.*, vol. 14, no. 3, pp. 1730–1745, Mar. 2015.
- [34] D. Palomar and J. Fonollosa, "Practical algorithms for a family of waterfilling solutions," *IEEE Trans. Signal Process.*, vol. 53, no. 2, pp. 686–695, Feb. 2005.

- [35] D. Park, "Weighted sum rate maximization of MIMO broadcast and interference channels with confidential messages," *IEEE Trans. Wireless Commun.*, vol. 15, no. 3, pp. 1742–1753, Mar. 2016.
- [36] S. Jaeckel, L. Raschkowski, K. Börner, and L. Thiele, "QuADriGa: A 3-D multi-cell channel model with time evolution for enabling virtual field trials," *IEEE Trans. Antennas Propag.*, vol. 62, no. 6, pp. 3242–3256, Jan. 2014.
- [37] S. Jaeckel, L. Raschkowski, F. Burkhardt, and L. Thiele, "Efficient sum-of-sinusoids-based spatial consistency for the 3GPP new-radio channel model," in *Proc. IEEE Globecom Workshops (GC Wkshps)*, Abu Dhabi, UAR, Dec. 2018, pp. 1–7.
- [38] S. Jaeckel, L. Raschkowski, S. Wu, L. Thiele, and W. Keusgen, "An explicit ground reflection model for mm-wave channels," in *Proc. IEEE Wireless Commun. Netw. Conf. Workshops (WCNCW)*, San Francisco, CA, USA, Mar 2017, pp. 1–5.
- [39] K. Shen and W. Yu, "Fractional programming for communication systems-part I: Power control and beamforming," *IEEE Trans. Signal Process.*, vol. 66, no. 10, pp. 2616–2630, May 2018.
- [40] M. Wu, B. Yin, G. Wang, C. Dick, J. R. Cavallaro, and C. Studer, "Large-scale MIMO detection for 3GPP LTE: Algorithms and FPGA implementations," *IEEE J. Sel. Topics in Signal Process.*, vol. 8, no. 5, pp. 916–929, Oct. 2014.
- [41] B. Beavis and I. M. Dobbs, *Optimization and stability theory for economic analysis*. Cambridge university press, 1990.



**Ruiling Hou** (Graduate Student Member, IEEE) received the B.E. degree in information engineering from the Chien-Shiung Wu College, Southeast University, Nanjing, China, in 2023. He is currently pursuing the Ph.D. degree in information engineering from the School of Information Science and Engineering, Southeast University, Nanjing. His research interests include optimization theory for wireless communications and information theory.



**Jiaheng Wang** (Senior Member, IEEE) received the Ph.D. degree in electronic and computer engineering from the Hong Kong University of Science and Technology, Kowloon, Hong Kong, in 2010, and the B.E. and M.S. degrees from the Southeast University, Nanjing, China, in 2001 and 2006, respectively.

He is currently a Full Professor at the National Mobile Communications Research Laboratory (NCRL), Southeast University, Nanjing, China. He is also with Purple Mountain Laboratories Nanjing, China. From 2010 to 2011, he was with the Signal Processing Laboratory, KTH Royal Institute of Technology, Stockholm, Sweden. He also held visiting positions at the Friedrich Alexander University Erlangen-Nürnberg, Nürnberg, Germany, and the University of Macau, Macau. His research interests are mainly on communication systems, wireless networks and network security.

Dr. Wang has published more than 200 articles on international journals and conferences. He serves as an Editor for the IEEE Transactions on Wireless Communications and an Editor for the IEEE Transactions on Communications. He was a Senior Area Editor for the IEEE Signal Processing Letters. He was a recipient of the Humboldt Fellowship for Experienced Researchers and the best paper awards of IEEE GLOBECOM 2019, ADHOCNETS 2019, and WCSP 2022 and 2014.



**Rui Zhou** (Member, IEEE) received the B.Eng. degree in information engineering from Southeast University, Nanjing, China, in 2017, and the Ph.D. degree from the Hong Kong University of Science and Technology (HKUST), Hong Kong, in 2021. He is currently a Research Scientist at the Shenzhen Research Institute of Big Data and an Adjunct Assistant Professor at the School of Science and Engineering, The Chinese University of Hong Kong, Shenzhen, China. His research interests include optimization algorithms, statistical signal processing,

machine learning, and financial engineering.



**Daniel P. Palomar** (Fellow, IEEE) received the Electrical Engineering and Ph.D. degrees from the Technical University of Catalonia (UPC), Barcelona, Spain, in 1998 and 2003, respectively, and was a Fulbright Scholar at Princeton University during 2004–2006. He is a Professor in the Department of Electronic & Computer Engineering at the Hong Kong University of Science and Technology (HKUST), Hong Kong, which he joined in 2006. He had previously held several research appointments, namely, at King's College London (KCL), London,

UK; Stanford University, Stanford, CA; Telecommunications Technological Center of Catalonia (CTTC), Barcelona, Spain; Royal Institute of Technology (KTH), Stockholm, Sweden; University of Rome "La Sapienza", Rome, Italy; and Princeton University, Princeton, NJ. His current research interests include data analytics, optimization methods, and deep learning in financial systems. Dr. Palomar has been recognized as a EURASIP Fellow (2024), an IEEE Fellow (2012), and, among others, has been awarded with the 2004/06 Fulbright Research Fellowship and the 2004, 2015, and 2020 Young Author Best Paper Awards by the IEEE Signal Processing Society. He is the author of many research articles and books, including *Portfolio Optimization: Theory and Application* and *Convex Optimization in Signal Processing and Communications*. He has been a Guest Editor of the IEEE Journal of Selected Topics in Signal Processing 2016 Special Issue on "Financial Signal Processing and Machine Learning for Electronic Trading", an Associate Editor of IEEE Transactions on Information Theory and of IEEE Transactions on Signal Processing, a Guest Editor of the IEEE Signal Processing Magazine 2010 Special Issue on "Convex Optimization for Signal Processing," the IEEE Journal on Selected Areas in Communications 2008 Special Issue on "Game Theory in Communication Systems," and the IEEE Journal on Selected Areas in Communications 2007 Special Issue on "Optimization of MIMO Transceivers for Realistic Communication Networks."



**Xiqi Gao** (Fellow, IEEE) received the Ph.D. degree in electrical engineering from Southeast University, Nanjing, China, in 1997.

He joined the Department of Radio Engineering, Southeast University, in April 1992. Since May 2001, he has been a professor of information systems and communications. From September 1999 to August 2000, he was a Visiting Scholar at Massachusetts Institute of Technology, Cambridge, MA, USA, and Boston University, Boston, MA. From August 2007 to July 2008, he visited the Darmstadt

University of Technology, Darmstadt, Germany, as a Humboldt scholar. His current research interests include broadband multicarrier communications, MIMO wireless communications, channel estimation, and turbo equalization, and multirate signal processing for wireless communications.

Dr. Gao received the Science and Technology Awards of the State Education Ministry of China in 1998, 2006, and 2009; the National Technological Invention Award of China in 2011, and the 2011 IEEE Communications Society Stephen O. Rice Prize Paper Award in the field of communication theory. From 2007 to 2012, he served as an Editor for the IEEE Transactions on Wireless Communications. From 2009 to 2013, he served as an Editor for the IEEE Transactions on Signal Processing. From 2015 to 2017, he served as an Editor for the IEEE Transactions on Communications.



**Björn Ottersten** (Fellow, IEEE) received the M.S. degree from Linköping University, Linköping, Sweden, in 1986, and the Ph.D. degree from Stanford University, Stanford, CA, USA, in 1990. He has held research positions with the Department of Electrical Engineering, Linköping University, the Information Systems Laboratory, Stanford University, the Katholieke Universiteit Leuven, Leuven, Belgium, and the University of Luxembourg, Luxembourg. From 1996 to 1997, he was the Director of Research with ArrayComm, Inc., a start-up in San Jose, CA,

USA, based on his patented technology. In 1991, he was appointed Professor of signal processing with the Royal Institute of Technology (KTH), Stockholm, Sweden where he has held positions as head of department and dean. He is the founding Director for the Interdisciplinary Centre for Security, Reliability and Trust, University of Luxembourg. He is a recipient of the IEEE Fourier Technical Field Award, the IEEE Signal Processing Society Technical Achievement Award, the EURASIP Group Technical Achievement Award, and the European Research Council (ERC) advanced research grant twice. He has co-authored journal papers that received the IEEE Signal Processing Society Best Paper Award in 1993, 2001, 2006, 2013, and 2019, and 9 IEEE conference papers best paper awards. He has been a board member of IEEE Signal Processing Society, EURASIP, the Swedish Research Council and currently serves on the ERC Scientific Council and the board of the Swedish Foundation for Strategic Research. Dr. Ottersten has served as Editor in Chief of EURASIP Signal Processing, and acted on the editorial boards of IEEE Transactions on Signal Processing, IEEE Signal Processing Magazine, IEEE Open Journal for Signal Processing, EURASIP Journal of Advances in Signal Processing and Foundations and Trends in Signal Processing. He is a fellow of IEEE, EURASIP, AAIA, and the Royal Swedish Academy of Engineering Sciences.