

Multicell Massive MIMO Multicast Transmission With Finite-Alphabet Inputs

Wenqian Wu[✉], *Student Member, IEEE*, Chengshan Xiao[✉], *Fellow, IEEE*, and Xiqi Gao[✉], *Fellow, IEEE*

Abstract—This paper investigates the precoding design for multicast transmission in multicell massive multiple-input multiple-output (MIMO) systems with finite-alphabet inputs. The users within each cell are interested in common information, and different cells provide distinct information. Focusing on the weighted max-min fairness (MMF) problem with only statistical channel state information at the base station, we provide the necessary conditions of the optimal precoding vectors to maximize the minimum weighted achievable ergodic rate, and an iterative algorithm is proposed to optimize the precoding vectors. To achieve lower computational complexity, we then derive a lower bound on the achievable ergodic rate for finite-alphabet inputs. Considering the problem of the minimum weighted rate lower bound maximization, we utilize the concave-convex procedure (CCCP) to develop a CCCP-based algorithm, which is proven to converge to a local optimum. Furthermore, exploiting the channel characteristic in massive MIMO systems, we prove that the optimal precoding vectors, maximizing the minimum weighted rate lower bound, are linear combinations of eigenvectors of transmit correlation matrices, and the original problem can be shifted into a lower dimensional space. Motivated by this insight, a relation-based algorithm is devised to obtain the optimal solution of the weighted MMF problem by using the duality between the MMF problem and the quality of service problem. Numerical results illustrate the tightness of the achievable ergodic rate lower bound and the significant performance of the devised algorithms.

Index Terms—Multicast transmission, massive MIMO, finite-alphabet signals, statistical CSI.

I. INTRODUCTION

WITH the advent of data-hungry applications and services, new wireless communication technologies are proposed to utilize energy and spectrum resources more efficiently [1]. In wireless networks, transmitting common information, such as

regular system updates, financial data or headline news, to several mobile terminals simultaneously is a typical scenario. To enable such service, multicast transmission has been considered in different releases of the Third Generation Partnership Project (3GPP) [2]. Generally, with channel state information at transmitter (CSIT), multicasting can be performed by precoding to increase the received signal-to-interference-plus-noise ratio (SINR). In the pioneering work of multicast beamforming for a single group of users [3], it was proven that both quality of service (QoS) and max-min fairness (MMF) problems were NP-hard. Then, the work in [3] was further extended to the case of multiple groups where a duality between the QoS and MMF problems was revealed [4]. In [5], joint multicast beamforming for multigroup multicell systems was investigated under the per-cell power constraints. Furthermore, the work in [6] proposed a cooperative multicast transmission scheme for the terrestrial-satellite network. In [3]–[6], with perfect channel state information (CSI), the semidefinite relaxation (SDR) method is adopted to find near-optimal solutions.

Recently, massive multiple-input multiple-output (MIMO) is regarded as a promising technology for next generation wireless systems to achieve significant performance in terms of spectrum efficiency and reliability [7]–[10]. In massive MIMO systems, the base station (BS) is equipped with a large number of antennas to simultaneously serve a number of mobile terminals in the same time-frequency resource. It was first proven in [7] that the intra-cell interferences and uncorrelated noises can be mitigated by employing unsophisticated beamforming in noncooperative massive MIMO systems with unlimited numbers of BS antennas. Since the publication of [7], massive MIMO systems has been investigated from various aspects, for example, in [11]–[14].

Motivated by the potential benefits of massive MIMO, several works were dedicated to the research of multicast transmission for massive MIMO communications [15]–[18]. Considering noncooperative multicell massive MIMO multicast transmission, the authors in [15] not only derived asymptotically optimal beamforming with perfect instantaneous CSIT, but also proposed a contamination-free pilot scheme to tackle the multicast beamforming with imperfect instantaneous CSIT. To find an efficient beamforming algorithm for multi-group multicasting in large-scale systems, [16] proposed a fast algorithm which adopts the concave-convex procedure (CCCP) [19] and the alternating direction method of multipliers (ADMM) [20] methods. Also, a two-layer precoding scheme was proposed in [17] for multigroup multicasting in large-scale antenna systems, aiming to reduce the computational burden of precoder design. In addition, with minimum mean-square error (MMSE) channel estimation, the performance of different multigroup multicast

Manuscript received January 16, 2019; revised April 16, 2019; accepted April 21, 2019. Date of publication May 15, 2019; date of current version July 16, 2019. The work of W. Wu and X. Gao was supported in part by National Natural Science Foundation of China under Grants 61320106003, 61761136016, 61801114, and 61631018, in part by the National Science and Technology Major Project of China under Grant 2017ZX03001002-004, and in part by the Huawei Cooperation Project. The work of C. Xiao was supported in part by US National Science Foundation under Grant ECCS-1827592. Part of this work was carried out while W. Wu was visiting Lehigh University, Bethlehem, PA. This paper was presented in part at the IEEE International Conference on Communications (ICC), Shanghai, China, 2019. The review of this paper was coordinated by Dr. Z. Ding. (*Corresponding author: Xiqi Gao.*)

W. Wu and X. Gao are with the National Mobile Communications Research Laboratory, Southeast University, Nanjing 210096, China (e-mail: wq_wu@seu.edu.cn; xqgao@seu.edu.cn).

C. Xiao is with the Department of Electrical and Computer Engineering, Lehigh University, Bethlehem, PA 18015 USA (e-mail: xiaoc@lehigh.edu).

Digital Object Identifier 10.1109/TVT.2019.2917069

precoding schemes was analyzed for six possible scenarios in massive MIMO systems [18]. Note that the aforementioned works on massive MIMO multicast transmission are founded on the assumption that the instantaneous CSI of users is available at the BS.

For massive MIMO systems, one serious challenge is the acquisition of accurate instantaneous CSIT. In general, the instantaneous CSIT can be acquired by exploiting the channel reciprocity in time-division duplex (TDD) systems or uplink feedback in frequency-division duplex (FDD) systems. Nevertheless, as mobility of terminals increases, the fluctuations of channel changes more rapidly and thus the accurate instantaneous CSIT acquisition can be challenging, where the delays obtaining the instantaneous CSIT are non-negligible with respect to the channel coherence time, resulting in outdated instantaneous CSIT. Therefore, in such case, exploiting statistical CSIT appears to be more reasonable for its robustness. Moreover, most existing precoding designs for massive MIMO multicast transmission focus on the received SINR, while achievable rate with finite-alphabet inputs can be a more specific metric to study the multicast transmission in massive MIMO systems. Although Gaussian inputs are information-theoretic optimal signals, it is important to note that practical transmit signals are often generated from finite constellation sets, e.g., phase-shift keying (PSK) and quadrature amplitude modulation (QAM). Additionally, the system throughput with finite-alphabet inputs is significantly different from using Gaussian inputs, which brings a substantial performance gap between the precoding schemes based on finite-alphabet inputs and those based on Gaussian inputs [21], [22]. In [22], the globally optimal linear precoder with finite-alphabet inputs was studied for point-to-point MIMO communication, and optimization problems of precoding matrix with finite-alphabet inputs for distinct scenarios were further investigated in [23]–[26]. With only statistical CSIT, precoding algorithms to improve the achievable ergodic rate have been investigated for finite-alphabet inputs [27]–[29]. However, to the best of our knowledge, the precoding design for multicell massive MIMO multicast transmission with finite-alphabet inputs and statistical CSIT is still an open and challenging problem.

This paper investigates the multicast transmission for multicell massive MIMO communications with finite-alphabet inputs, where each BS only has access to the statistical CSI of users. Our key contributions are summarized as follows:

- We provide the necessary conditions of the optimal precoding vectors to maximize the minimum weighted achievable ergodic rate, where finite-alphabet inputs and Gaussian interference approximation are considered. Based on these conditions, an iterative algorithm to search the optimal precoding vectors is proposed in this paper.
- To avoid the cumbersome computations of the achievable ergodic rate, we derive a lower bound on the achievable ergodic rate, which hinges on the transmit correlation matrices. Investigating the weighted MMF problem with the rate lower bound, we develop an iterative algorithm for precoding design by utilizing the CCCP method. It can be proven that the precoding vectors generated from the CCCP-based algorithm converge to a locally optimal solution of the weighted MMF problem.

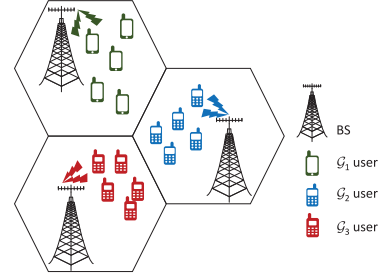


Fig. 1. Illustration of multicell multicast network with massive MIMO.

- By exploiting the channel characteristic in massive MIMO systems, we obtain the structure of the optimal precoding vectors to maximize the minimum weighted rate lower bound, which reveals that the optimal precoding vectors should be linear combinations of eigenvectors of transmit correlation matrices, and the optimization problem can be simplified into a lower dimensional space. Based on this result, an iterative algorithm is developed to search the optimal solution of the weighted MMF problem by using the duality between the QoS problem and the MMF problem.

Notation: Boldface lower case letters and boldface upper case letters denote vectors and matrices respectively. The operators $(\cdot)^H$, $(\cdot)^T$, and $(\cdot)^*$ denote the matrix conjugate-transpose, transpose and conjugate operations respectively. $\mathbf{1}$ denotes all-one vector and $\mathbf{0}$ denotes all-zero vector or matrix. The trace, ensemble expectation, real part, absolute value and the Euclidean norm operators are denoted by $\text{tr}(\cdot)$, $\mathbb{E}\{\cdot\}$, $\Re\{\cdot\}$, $|\cdot|$ and $\|\cdot\|$. we use $\mathbf{A} \succeq \mathbf{0}$ to denote a positive semidefinite Hermitian matrix \mathbf{A} and use $\mathbf{a} \succeq \mathbf{0}$ to denote vector $\mathbf{a} \in \mathbb{R}_+^{M \times 1}$. $\text{diag}(\mathbf{A})$ indicates a column vector, which is constituted by the main diagonal of \mathbf{A} . \odot denotes the Hadamard product.

II. SYSTEM MODEL AND PRELIMINARIES

We consider a multicell multicast transmission system depicted in Fig. 1, comprising L cells and K users per cell. Different from the widely investigated multiuser massive MIMO transmission for the user-specific information, we consider the scenario that there are a group of K users seeking for common information in each cell [15], [30]. For convenience, we define $\mathcal{L} \triangleq \{1, 2, \dots, L\}$ as the cell index set and \mathcal{G}_j as the index set of the users in the j th cell, where $|\mathcal{G}_j| = K$. Each BS is equipped with M transmit antennas, and K single-antenna users are uniformly distributed within each cell. Let x_j denote the common information signal for the K users in the j th cell with $\mathbb{E}\{|x_j|^2\} = 1$ and $\mathbf{p}_j \in \mathbb{C}^{M \times 1}$ denote the corresponding multicast precoding vector. The maximum transmit power at the j th BS is defined as P_j , i.e., $\|\mathbf{p}_j\|^2 \leq P_j$. Then, the received signal at the k th user in the j th cell can be expressed as

$$y_{j,k} = \mathbf{h}_{j,j,k}^H \mathbf{p}_j x_j + \sum_{i \neq j} \mathbf{h}_{i,j,k}^H \mathbf{p}_i x_i + z_{j,k} \quad (1)$$

where $z_{j,k} \sim \mathcal{CN}(0, 1)$. The vector $\mathbf{h}_{i,j,k} \in \mathbb{C}^{M \times 1}$ denotes the channel from the i th BS to the k th user in the j th cell. In this paper, the correlated fading channel is adopted, and $\mathbf{h}_{i,j,k}$ is

modeled as [31]

$$\mathbf{h}_{i,j,k} = \mathbf{R}_{i,j,k}^{\frac{1}{2}} \mathbf{g}_{i,j,k} \quad (2)$$

where $\mathbf{R}_{i,j,k} \in \mathbb{C}^{M \times M}$ is the transmit correlation matrix and $\mathbf{g}_{i,j,k} \in \mathbb{C}^{M \times 1}$ has independent identically distributed (i.i.d.) complex Gaussian entries with zero-mean and unit variance. In addition, the correlation matrix $\mathbf{R}_{i,j,k}$ can be expressed as

$$\mathbf{R}_{i,j,k} = \mathbf{V}_{i,j,k} \tilde{\mathbf{R}}_{i,j,k} \mathbf{V}_{i,j,k}^H \quad (3)$$

where $\tilde{\mathbf{R}}_{i,j,k} \in \mathbb{C}^{M \times M}$ is a diagonal matrix and $\mathbf{V}_{i,j,k} \in \mathbb{C}^{M \times M}$ is a deterministic unitary matrix. In massive MIMO systems, as $M \rightarrow \infty$, matrices $\mathbf{V}_{i,j,k}$ tend to be an identical deterministic unitary matrix \mathbf{V} , which only hinges on the topology of BS antenna array [32], i.e.,

$$\mathbf{R}_{i,j,k} \rightarrow \mathbf{V} \tilde{\mathbf{R}}_{i,j,k} \mathbf{V}^H, \quad \text{as } M \rightarrow \infty. \quad (4)$$

Especially, if the uniform linear array (ULA) is adopted at the BS, the discrete Fourier transform (DFT) matrix offers a good approximation to matrix \mathbf{V} [33].

In this paper, we assume that the transmit correlation matrices $\mathbf{R}_{i,j,k}$ are perfectly known at the BSs.¹ Instead of adopting the conventional assumption of Gaussian input data, we consider the transmit data symbols of the j th BS are generated from a discrete constellation of size Q_j with i.i.d. uniform distribution. Moreover, in a practical system, for the users in the j th cell, users may not have knowledge about the precoding vectors \mathbf{p}_i ($i \neq j$) of interfering users. Thus, it is practically impossible to exploit the finite-alphabet nature of the interferences. Consequently, we assume that each user views the interference as Gaussian while modeling its intended signal as finite-alphabet [24]. Based on the Gaussian interference approximation, for the k th user in the j th cell, the variance of the interference plus noise $\sum_{i \neq j} \mathbf{h}_{i,j,k}^H \mathbf{p}_i \mathbf{x}_i + z_{j,k}$ can be expressed as

$$\sigma_{j,k}^2 = \sum_{i \neq j} |\mathbf{h}_{i,j,k}^H \mathbf{p}_i|^2 + 1. \quad (5)$$

With the finite-alphabet signal Gaussian interference approximation, the received signal $y_{j,k}$ can be rewritten as

$$y_{j,k} = \mathbf{h}_{j,j,k}^H \mathbf{p}_j x_j + z'_{j,k} \quad (6)$$

where $z'_{j,k}$ is a zero-mean additive Gaussian noise with variance $\sigma_{j,k}^2$. Then, the achievable ergodic rate for the k th user in the j th cell can be expressed as [27]

$$R_{j,k} = \log Q_j - 1/Q_j \times \sum_{m=1}^{Q_j} \mathbb{E} \left\{ \log \sum_{n=1}^{Q_j} \exp \left(- \frac{|\mathbf{h}_{j,j,k}^H \mathbf{p}_j d_j^{m,n} + z'_{j,k}|^2 - |z'_{j,k}|^2}{\sum_{i \neq j} |\mathbf{h}_{i,j,k}^H \mathbf{p}_i|^2 + 1} \right) \right\} \quad (7)$$

where scalar $d_j^{m,n} \triangleq a_{j,m} - a_{j,n}$ and $a_{j,m}$ is the m th element in the constellation range for x_j .

¹Owing to the fact that the transmit correlation matrices are independent of sub-carriers [11], [34], the transmit correlation matrices acquisition process can be simplified. Moreover, with the statistical channel model [35], a relatively small number of parameters need to be counted, and the transmit correlation matrices $\mathbf{R}_{i,j,k}$ can be obtained efficiently by utilizing the channel correlation matrices acquisition method in [11].

III. PRECODING DESIGN WITH ERGODIC RATE

In this section, the multicast precoding design based on the achievable ergodic rate in (7) is investigated. By solving the corresponding weight MMF problem, we derive the necessary conditions of the optimal precoding vectors that maximizing the minimum weighted achievable ergodic rate. Then, we propose an iterative algorithm to optimize the precoding vectors with these necessary conditions.

For multicast transmission, the performance of the worst user is the bottleneck. Aiming to maximize the minimum weighted achievable ergodic rate among all KL served users in L cells, the corresponding MMF problem can be formulated as

$$\begin{aligned} \mathcal{F} : \quad & \max_{\{\mathbf{p}_j\}_{j=1}^L} \min_{j \in \mathcal{L}, k \in \mathcal{G}_j} \frac{1}{\theta_{j,k}} R_{j,k} \\ \text{s.t.} \quad & \|\mathbf{p}_j\|^2 \leq P_j, \forall j \in \mathcal{L} \end{aligned} \quad (8)$$

where $\theta_{j,k}$ is the predetermined target rate for the k th user in the j th cell. Note that specifying the target achievable ergodic rate for each user makes the problem \mathcal{F} more general, where each user's achievable ergodic rate is scaled by the weight factor $1/\theta_{j,k}$ to consider for different levels of service. Furthermore, for problem \mathcal{F} , since it is difficult to obtain the expression for the partial derivative of the objective function with respect to \mathbf{p}_j . To tackle problem \mathcal{F} easier, we introduce an auxiliary (positive real) variable t to lower bound the minimum weighted rate $\frac{1}{\theta_{j,k}} R_{j,k}$, where problem \mathcal{F} can be equivalently formulated as [3]–[6]

$$\begin{aligned} \mathcal{F}_t : \quad & \max_{\{\mathbf{p}_j\}_{j=1}^L, t} t \\ \text{s.t.} \quad & R_{j,k} \geq \theta_{j,k} t, \forall k \in \mathcal{G}_j, \forall j \in \mathcal{L} \\ & \|\mathbf{p}_j\|^2 \leq P_j, \forall j \in \mathcal{L}. \end{aligned} \quad (9)$$

Then, we derive the necessary conditions in the following proposition, which characterizes the optimal solution of problem \mathcal{F}_t .

Proposition 1: The optimal precoding design, which is the optimal solution of \mathcal{F}_t , satisfies the following conditions

$$\begin{aligned} & \sum_{k=1}^K \zeta_{j,k} \mathbb{E} \{ \varepsilon_{1,j,k} \mathbf{h}_{j,j,k}^* \} - \sum_{i \neq j} \sum_{k=1}^K \zeta_{i,k} \mathbb{E} \{ \varepsilon_{2,j,i,k} \mathbf{h}_{j,i,k}^* \} \\ & - \lambda_j \mathbf{p}_j^* = \mathbf{0}, \forall j \in \mathcal{L} \end{aligned} \quad (10a)$$

$$\zeta_{j,k} (R_{j,k} - \theta_{j,k} t) = 0, R_{j,k} \geq \theta_{j,k} t, \forall j \in \mathcal{L}, \forall k \in \mathcal{G}_j \quad (10b)$$

$$\lambda_j (\|\mathbf{p}_j\|^2 - P_j) = 0, \|\mathbf{p}_j\|^2 \leq P_j, \forall j \in \mathcal{L} \quad (10c)$$

$$1 - \sum_{j,k} \zeta_{j,k} \theta_{j,k} = 0 \quad (10d)$$

where $\zeta_{j,k} \geq 0$ and $\lambda_j \geq 0$. $\varepsilon_{1,j,k}$ and $\varepsilon_{2,j,i,k}$ are random variables given in (52) and (55), respectively.

Proof: See Appendix A. ■

Proposition 1 provides the necessary conditions of the optimal precoding vectors which maximize the minimum weighted achievable ergodic rate among the KL users. Generally, due to the non-convexity and the complex representation of problem \mathcal{F}_t , it is impractical to obtain a closed-form expression for the optimal precoding vectors \mathbf{p}_j . Nevertheless, Proposition 1 reveals the partial derivative of $R_{j,k}$ with respect to \mathbf{p}_j . Based

on these results, we develop an iterative algorithm to search the optimal precoding vectors numerically. Here, we first utilize the barrier method [36] to approximately reformulate the inequality constrained problem \mathcal{F}_t into the following unconstrained problem

$$\min_{\{\mathbf{p}_j\}_{j=1}^L, t} h(\{\mathbf{p}_j\}_{j=1}^L, t) = -\alpha t + \sum_j \phi_j(\mathbf{p}_j) + \sum_{j,k} \varphi_{j,k}(\{\mathbf{p}_i\}_{i=1}^L, t) \quad (11)$$

where $\alpha > 0$ sets the accuracy of the approximation, and the quality of the approximation improves as α grows. $\phi_j(\mathbf{p}_j)$ and $\varphi_{j,k}(\{\mathbf{p}_i\}_{i=1}^L, t)$ are the log barrier functions corresponding to the inequality constraints $\|\mathbf{p}_j\|^2 \leq P_j$ and $R_{j,k} \geq \theta_{j,k}t$, respectively. $\phi_j(\mathbf{p}_j)$ and $\varphi_{j,k}(\{\mathbf{p}_i\}_{i=1}^L, t)$ can be expressed as

$$\phi_j(\mathbf{p}_j) = \begin{cases} +\infty, & \|\mathbf{p}_j\|^2 \geq P_j \\ -\log(P_j - \mathbf{p}_j^H \mathbf{p}_j), & \|\mathbf{p}_j\|^2 < P_j \end{cases} \quad (12)$$

and

$$\varphi_{j,k}(\{\mathbf{p}_i\}_{i=1}^L, t) = \begin{cases} +\infty, & R_{j,k} \leq \theta_{j,k}t \\ -\log(R_{j,k} - \theta_{j,k}t), & R_{j,k} > \theta_{j,k}t. \end{cases} \quad (13)$$

Then, we solve problem in (11) with the gradient descent method [23], where the partial derivative of $h(\{\mathbf{p}_l\}_{l=1}^L, t)$ with respect to \mathbf{p}_j and t , respectively, are given by

$$\frac{\partial h(\{\mathbf{p}_l\}_{l=1}^L, t)}{\partial \mathbf{p}_j} = \frac{\mathbf{p}_j^*}{P_j - \mathbf{p}_j^H \mathbf{p}_j} - \sum_{k=1}^K \frac{\mathbb{E}\{\varepsilon_{1,j,k} \mathbf{h}_{j,j,k}^*\}}{R_{j,k} - \theta_{j,k}t} + \sum_{i \neq j} \sum_{k=1}^K \frac{\mathbb{E}\{\varepsilon_{2,j,i,k} \mathbf{h}_{j,i,k}^*\}}{R_{i,k} - \theta_{i,k}t} \quad (14)$$

$$\frac{\partial h(\{\mathbf{p}_l\}_{l=1}^L, t)}{\partial t} = -\alpha + \sum_{j,k} \frac{\theta_{j,k}}{R_{j,k} - \theta_{j,k}t}. \quad (15)$$

During each iteration, the backtracking line search method [36] is utilized to choose the step size. The detailed steps of the developed gradient-based algorithm are given in Algorithm 1.

Algorithm 1 provides a method to numerically solve problem \mathcal{F}_t . However, without closed-form expression, evaluating (14) and (15) with Monte-Carlo method is computationally cumbersome. Therefore, in the next section, we derive a lower bound on the achievable ergodic rate and devise two efficient algorithms, maximizing the minimum weighted rate lower bound among the KL users in L cells.

IV. PRECODING DESIGN WITH ACHIEVABLE ERGODIC RATE LOWER BOUND

In this section, we introduce a lower bound on the achievable ergodic rate in (7). Based on the rate lower bound, the MMF problem to maximize the minimum weighted rate lower bound is investigated, and an iterative algorithm with the CCCP is developed to find the local optimum of the weighted MMF problem. Then, since a large number of BS antennas are considered in this paper, we propose another iterative algorithm to search the optimal solution of the weighted MMF problem, where the

Algorithm 1: Gradient-based Precoding Design.

- 1: Initialization. Choose vectors \mathbf{p}_j as $\mathbf{p}_j = \sqrt{\frac{P_j}{M}} \mathbf{1}, \forall j \in \mathcal{L}$. Then, calculate the corresponding rate $R_{j,k}$ for each user and initialize t as $t = \min_{j,k} \frac{1}{\theta_{j,k}} R_{j,k}$. Set $\alpha := \alpha^{(0)} > 0$ and $c > 1$.
 - 2: **repeat**
 - 3: Set $u := 0, \mathbf{x}_j^{(0)} = \mathbf{p}_j$ and $t_{\text{in}}^{(0)} = t$.
 - 4: **repeat**
 - 5: Choose step size μ and ν by backtracking line search.
 - 6: Update $\mathbf{x}_j^{(u+1)} = \mathbf{x}_j^{(u)} + \mu \frac{\partial h(\{\mathbf{x}_i\}_{i=1}^L, t_{\text{in}})}{\partial \mathbf{x}_j}, t_{\text{in}}^{(u+1)} = t_{\text{in}}^{(u)} + \nu \frac{\partial h(\{\mathbf{x}_i\}_{i=1}^L, t_{\text{in}})}{\partial t_{\text{in}}}$, and calculate increment $\tau = |h(\{\mathbf{x}_i^{(u+1)}\}_{i=1}^L, t_{\text{in}}^{(u+1)}) - h(\{\mathbf{x}_i^{(u)}\}_{i=1}^L, t_{\text{in}}^{(u)})|$.
 - 7: Set $u = u + 1$.
 - 8: **until** $\tau < \epsilon_1$ (ϵ_1 is a predefined tolerance.)
 - 9: Update $\mathbf{p}_j = \mathbf{x}_j^{(u)}, t = t_{\text{in}}^{(u)}$, and $\alpha = c\alpha$.
 - 10: **until** $\frac{1}{\alpha} < \epsilon_2$ (ϵ_2 is a predefined tolerance.)
-

channel characteristic of massive MIMO systems and the duality between the QoS problem and the MMF problem are exploited.

A. Achievable Ergodic Rate Lower Bound

In general, it is impractical to derive the achievable ergodic rate $R_{j,k}$ in closed-form, which indicates that evaluating $R_{j,k}$ with Monte-Carlo method can be computationally inefficient. To avoid the cost of the Monte-Carlo averaging over the noise and channels, we introduce an achievable ergodic rate lower bound $R_{\text{lb},j,k}$.

Proposition 2: The achievable ergodic rate $R_{j,k}$ in (7) is lower bounded by

$$R_{\text{lb},j,k} = \log Q_j - (1/\ln 2 - 1) - \frac{1}{Q_j} \sum_{m=1}^{Q_j} \log \sum_{n=1}^{Q_j} \left(1 + \frac{1}{2} \frac{|d_j^{m,n}|^2 \mathbf{p}_j^H \mathbf{R}_{j,j,k} \mathbf{p}_j}{\sum_{i \neq j} \mathbf{p}_i^H \mathbf{R}_{i,j,k} \mathbf{p}_i + 1} \right)^{-1}. \quad (16)$$

Proof: See Appendix B. ■

In Proposition 2, we derive a closed-form lower bound on the achievable ergodic rate, where the closed-form lower bound $R_{\text{lb},j,k}$ only hinges on the modulation type, the precoding vectors \mathbf{p}_j and the transmit correlation matrices $\mathbf{R}_{i,j,k}$. Therefore, instead of tedious Monte-Carlo averaging over every channel realization, it is an efficient method to evaluate the system performance by using the rate lower bound. In order to reduce the complexity of precoding design, we turn to consider the MMF problem which maximizes the minimum weighted rate lower bound of the KL served users in L cells, and we can find that the MMF problem with the rate lower bound can be further simplified by utilizing the property of the lower bound. First, the MMF problem of the minimum weighted rate lower bound

maximization can be written as

$$\begin{aligned} \mathcal{F}_{\text{lb}} : \quad & \max_{\{\mathbf{p}_j\}_{j=1}^L} \min_{j \in \mathcal{L}, k \in \mathcal{G}_j} \frac{1}{\theta_{\text{lb},j,k}} R_{\text{lb},j,k} \\ \text{s.t.} \quad & \|\mathbf{p}_j\|^2 \leq P_j, \forall j \in \mathcal{L} \end{aligned} \quad (17)$$

where $\theta_{\text{lb},j,k}$ is the predetermined target rate lower bound for the k th user in the j th cell.

To solve \mathcal{F}_{lb} efficiently, we first analyze the expression of the rate lower bound $R_{\text{lb},j,k}$ in (16). Define the parameter $\gamma_{j,k}$ as

$$\gamma_{j,k} = \frac{\mathbf{p}_j^H \mathbf{R}_{j,j,k} \mathbf{p}_j}{\sum_i \mathbf{p}_i^H \mathbf{R}_{i,j,k} \mathbf{p}_i + 1}. \quad (18)$$

Then, the derivative of $R_{\text{lb},j,k}$ with respect to $\gamma_{j,k} \geq 0$ is given by

$$\begin{aligned} \frac{dR_{\text{lb},j,k}}{d\gamma_{j,k}} &= \frac{1}{\ln 2 \cdot Q_j} \sum_{m=1}^{Q_j} \frac{1}{\sum_{n=1}^{Q_j} 2 \left(1 + \frac{1}{2} |d_j^{m,n}|^2 \gamma_{j,k}\right)^{-1}} \\ &\quad \times \sum_{n=1}^{Q_j} \frac{|d_j^{m,n}|^2}{\left(1 + \frac{1}{2} |d_j^{m,n}|^2 \gamma_{j,k}\right)^2} \\ &\geq 0. \end{aligned} \quad (19)$$

From (19), it is easy to find that $R_{\text{lb},j,k}$ monotonously increases with the increase of $\gamma_{j,k}$, which means that the rate lower bound $R_{\text{lb},j,k}$ can be characterized by the parameter $\gamma_{j,k}$. Instead of solving the MMF optimization problem \mathcal{F}_{lb} with achievable ergodic rate lower bound, we consider the following equivalent optimization problem by introducing an auxiliary positive real variable t

$$\begin{aligned} \mathcal{F}_{\text{lb},t} : \quad & \max_{\{\mathbf{p}_j\}_{j=1}^L, t} t \\ \text{s.t.} \quad & \gamma_{j,k} \geq \rho_{j,k} t, \forall j \in \mathcal{L}, \forall k \in \mathcal{G}_j \\ & \|\mathbf{p}_j\|^2 \leq P_j, \forall j \in \mathcal{L} \end{aligned} \quad (20)$$

where $\rho_{j,k}$ is the predetermined target for $\gamma_{j,k}$. After reformulating the weighted MMF optimization problem, we propose two efficient algorithms, which are CCCP-based algorithm and relation-based algorithm, to solve problem $\mathcal{F}_{\text{lb},t}$.

B. CCCP-Based Algorithm

Generally, problem $\mathcal{F}_{\text{lb},t}$ is a non-convex problem with a non-convex feasible set. To tackle this problem, we define vector \mathbf{s}_j as $\mathbf{s}_j \triangleq [\mathbf{p}_j^T, t]^T$ and rewrite problem $\mathcal{F}_{\text{lb},t}$ into the following form

$$\begin{aligned} \hat{\mathcal{F}}_{\text{lb},t} : \quad & \max_{\{\mathbf{p}_j\}_{j=1}^L, t} t \\ \text{s.t.} \quad & f_{j,k}(\{\mathbf{p}_i\}_{i \neq j}) - g_{j,k}(\mathbf{s}_j) \leq 0, \forall j \in \mathcal{L}, \forall k \in \mathcal{G}_j \\ & \|\mathbf{p}_j\|^2 \leq P_j, \forall j \in \mathcal{L} \end{aligned} \quad (21)$$

where functions $f_{j,k}(\{\mathbf{p}_i\}_{i \neq j})$ and $g_{j,k}(\mathbf{s}_j)$ are defined as

$$f_{j,k}(\{\mathbf{p}_i\}_{i \neq j}) = \sum_{i \neq j} \mathbf{p}_i^H \mathbf{R}_{i,j,k} \mathbf{p}_i + 1 \quad (22)$$

$$g_{j,k}(\mathbf{s}_j) = \frac{\mathbf{p}_j^H \mathbf{R}_{j,j,k} \mathbf{p}_j}{\rho_{j,k} t}. \quad (23)$$

Meanwhile, the first-order Taylor expression [37] of the function $g_{j,k}(\mathbf{s}_j)$ at $\mathbf{s}_j^{(n)} \triangleq [(\mathbf{p}_j^{(n)})^T, t^{(n)}]^T$ is given by

$$\begin{aligned} \bar{g}_{j,k}(\mathbf{s}_j^{(n)}, \mathbf{s}_j) &= 2\Re \left\{ \frac{(\mathbf{p}_j^{(n)})^H \mathbf{R}_{j,j,k} \mathbf{p}_j}{\rho_{j,k} t^{(n)}} \right\} \\ &\quad - \frac{(\mathbf{p}_j^{(n)})^H \mathbf{R}_{j,j,k} \mathbf{p}_j^{(n)}}{\rho_{j,k} (t^{(n)})^2} t. \end{aligned} \quad (24)$$

Note that the first inequality constraint in $\hat{\mathcal{F}}_{\text{lb},t}$ refers to a difference of convex (DC) functions. Consequently, problem $\hat{\mathcal{F}}_{\text{lb},t}$ is a DC problem, which can be transformed into the following sequence of convex programs by utilizing the CCCP

$$\begin{aligned} \hat{\mathcal{F}}_{\text{lb},\text{iter},t} : \quad & [\mathbf{p}_1^{(n+1)}, \dots, \mathbf{p}_L^{(n+1)}, t^{(n+1)}] = \arg \max_{\{\mathbf{p}_j\}_{j=1}^L, t} t \\ \text{s.t.} \quad & f_{j,k}(\{\mathbf{p}_i\}_{i \neq j}) - \bar{g}_{j,k}(\mathbf{s}_j^{(n)}, \mathbf{s}_j) \leq 0, \forall j \in \mathcal{L}, \forall k \in \mathcal{G}_j \\ & \|\mathbf{p}_j\|^2 \leq P_j, \forall j \in \mathcal{L}. \end{aligned} \quad (25)$$

As can be seen from problem $\hat{\mathcal{F}}_{\text{lb},\text{iter},t}$, the main idea of the CCCP method is utilizing Taylor series expansion to linearize the convex part $g_{j,k}(\mathbf{s}_j)$ around a solution obtained from the previous iteration, which results in a convex feasible set. Then, the original problem $\hat{\mathcal{F}}_{\text{lb},t}$ is tackled as a sequence of convex programs. In order to further explain the feasibility of this approach, we present the following proposition.

Proposition 3: The sequence $\{\mathbf{p}_1^{(n)}, \dots, \mathbf{p}_L^{(n)}\}_{n=1}^\infty$ generated from problem $\hat{\mathcal{F}}_{\text{lb},\text{iter},t}$ converges to a stationary point of the original problem $\hat{\mathcal{F}}_{\text{lb},t}$.

Proof: See Appendix C. ■

Proposition 3 reveals that, utilizing the CCCP, we can obtain a locally optimal solution of the DC problem $\hat{\mathcal{F}}_{\text{lb},t}$, which indicates a candidate optimal precoding design. Then, for each iteration in $\hat{\mathcal{F}}_{\text{lb},\text{iter},t}$, it is a convex quadratically constrained quadratic program (QCQP), and we also utilize the barrier method to transfer $\hat{\mathcal{F}}_{\text{lb},\text{iter},t}$ into the following unconstrained problem

$$\begin{aligned} & [\mathbf{p}_1^{(n+1)}, \dots, \mathbf{p}_L^{(n+1)}, t^{(n+1)}] = \arg \min_{\{\mathbf{p}_j\}_{j=1}^L, t} -\alpha t \\ & \quad + \sum_j \phi_j(\mathbf{p}_j) + \sum_{j,k} \hat{\phi}_{j,k}^{(n)}(\{\mathbf{p}_i\}_{i=1}^L, t) \end{aligned} \quad (26)$$

where $\hat{\phi}_{j,k}^{(n)}(\{\mathbf{p}_i\}_{i=1}^L, t)$ is the log barrier function corresponding to the inequality constraint

$$q_{j,k}^{(n)}(\{\mathbf{p}_i\}_{i=1}^L, t) \triangleq f_{j,k}(\{\mathbf{p}_i\}_{i \neq j}) - \bar{g}_{j,k}(\mathbf{s}_j^{(n)}, \mathbf{s}_j) \leq 0 \quad (27)$$

and $\hat{\phi}_{j,k}^{(n)}(\{\mathbf{p}_i\}_{i=1}^L, t)$ is given by

$$\hat{\phi}_{j,k}^{(n)}(\{\mathbf{p}_i\}_{i=1}^L, t) = \begin{cases} +\infty, & q_{j,k}^{(n)}(\{\mathbf{p}_i\}_{i=1}^L, t) \geq 0 \\ -\log(-q_{j,k}^{(n)}(\{\mathbf{p}_i\}_{i=1}^L, t)), & \text{else.} \end{cases} \quad (28)$$

Then, the problem in (26) can be solved by the gradient descent method, where the step size in precoder update is chosen

Algorithm 2: CCCP-based Precoding Design.

-
- 1: Initialization. Set $n := 0$. Choose vectors $\mathbf{p}_j^{(0)}$ as $\mathbf{p}_j^{(0)} = \sqrt{\frac{P_j}{M}} \mathbf{1}$, and then set $t^{(0)}$ as
- $$t^{(0)} = \min_{j,k} \frac{(\mathbf{p}_j^{(0)})^H \mathbf{R}_{j,j,k} \mathbf{p}_j^{(0)}}{\rho_{j,k} \left(\sum_{i \neq j} (\mathbf{p}_i^{(0)})^H \mathbf{R}_{i,j,k} \mathbf{p}_i^{(0)} + 1 \right)}.$$
- 2: **repeat**
- 3: Set $\alpha := \alpha^{(0)} > 0$ and $c > 1$. Choose feasible vectors $\bar{\mathbf{p}}_1, \bar{\mathbf{p}}_2, \dots, \bar{\mathbf{p}}_L$ and \bar{t} .
- 4: **repeat**
- 5: Set $u := 0$, $\mathbf{x}_j^{(0)} = \bar{\mathbf{p}}_j$ and $t_{\text{in}}^{(0)} = \bar{t}$.
- 6: **repeat**
- 7: Choose step size μ, ν by backtracking line search.
- 8: Update $\mathbf{x}_j^{(u+1)} = \mathbf{x}_j^{(u)} + \mu \frac{\partial \hat{h}^{(n)}(\{\mathbf{x}_i\}_{i=1}^L, t_{\text{in}})}{\partial \mathbf{x}_j}$, $t_{\text{in}}^{(u+1)} = t_{\text{in}}^{(u)} + \nu \frac{\partial \hat{h}^{(n)}(\{\mathbf{x}_i\}_{i=1}^L, t_{\text{in}})}{\partial t_{\text{in}}}$, and calculate increment $\tau = |\hat{h}^{(n)}(\{\mathbf{x}_i^{(u+1)}\}_{i=1}^L, t_{\text{in}}^{(u+1)}) - \hat{h}^{(n)}(\{\mathbf{x}_i^{(u)}\}_{i=1}^L, t_{\text{in}}^{(u)})|$.
- 9: Set $u = u + 1$.
- 10: **until** $\tau < \epsilon_1$ (ϵ_1 is a predefined tolerance.)
- 11: Update $\bar{\mathbf{p}}_j = \mathbf{x}_j^{(u)}$, $\bar{t} = t_{\text{in}}^{(u)}$, and $\alpha = c\alpha$.
- 12: **until** $\frac{1}{\alpha} < \epsilon_2$ (ϵ_2 is a predefined tolerance.)
- 13: Update $\mathbf{p}_j^{(n+1)} = \bar{\mathbf{p}}_j$, $t^{(n+1)} = \bar{t}$ and $n = n + 1$.
- 14: **until** $|t^{(n)} - t^{(n-1)}| < \epsilon_3$ (ϵ_3 is a predefined tolerance.)
-

by backtracking line search method. We define the objective function in problem (26) as $\hat{h}^{(n)}(\{\mathbf{p}_l\}_{l=1}^L, t)$. The partial derivative of $\hat{h}^{(n)}(\{\mathbf{p}_l\}_{l=1}^L, t)$ with respect to \mathbf{p}_j and t , respectively, are given by

$$\begin{aligned} \frac{\partial \hat{h}^{(n)}(\{\mathbf{p}_l\}_{l=1}^L, t)}{\partial \mathbf{p}_j} &= \frac{\mathbf{p}_j^*}{P_j - \mathbf{p}_j^H \mathbf{p}_j} - \sum_{i \neq j} \sum_{k=1}^K \frac{\mathbf{R}_{j,i,k}^T \mathbf{p}_j^*}{q_{i,k}^{(n)}(\{\mathbf{p}_l\}_{l=1}^L, t)} \\ &\quad + \sum_{k=1}^K \frac{\mathbf{R}_{j,j,k}^T (\mathbf{p}_j^{(n)})^*}{\rho_{j,k} t^{(n)} q_{j,k}^{(n)}(\{\mathbf{p}_l\}_{l=1}^L, t)} \quad (29) \\ \frac{\partial \hat{h}^{(n)}(\{\mathbf{p}_l\}_{l=1}^L, t)}{\partial t} &= -\alpha - \sum_{j,k} \frac{(\mathbf{p}_j^{(n)})^H \mathbf{R}_{j,j,k} \mathbf{p}_j^{(n)}}{\rho_{j,k} (t^{(n)})^2 q_{j,k}^{(n)}(\{\mathbf{p}_l\}_{l=1}^L, t)}. \quad (30) \end{aligned}$$

Based on the above definitions, the steps of the proposed CCCP-based algorithm are given in Algorithm 2.

C. Relation-Based Algorithm

As mentioned in Section II, for massive MIMO systems, the eigenmatrices $\mathbf{V}_{i,j,k}$ of the transmit correlation matrices $\mathbf{R}_{i,j,k}$ become a unique deterministic unitary matrix \mathbf{V} as $M \rightarrow \infty$, where unitary matrix \mathbf{V} is independent of users. According to

this characteristic, we obtain the following result, which equivalently simplify the original problem $\mathcal{F}_{\text{lb},t}$ and reveals the structure of the optimal solution of $\mathcal{F}_{\text{lb},t}$.

Proposition 4: The optimal solution of problem $\mathcal{F}_{\text{lb},t}$ can be obtained by solving the following problem

$$\begin{aligned} \tilde{\mathcal{F}}_{\text{lb},t} : \quad & \max_{\{\mathbf{w}_j\}_{j=1}^L, t} t \\ \text{s.t.} \quad & \rho_{j,k} t \left(\sum_{i \neq j} \mathbf{w}_i^T \tilde{\mathbf{r}}_{i,j,k} + 1 \right) \\ & \quad - \mathbf{w}_j^T \tilde{\mathbf{r}}_{j,j,k} \leq 0, \forall k \in \mathcal{G}_j, \forall j \in \mathcal{L} \\ & \mathbf{1}^T \mathbf{w}_j \leq P_j, \forall j \in \mathcal{L} \\ & \mathbf{w}_j \succeq \mathbf{0}, \forall j \in \mathcal{L} \end{aligned} \quad (31)$$

where $\tilde{\mathbf{r}}_{i,j,k} \triangleq \text{diag}(\tilde{\mathbf{R}}_{i,j,k})$ and $\mathbf{w}_j \in \mathbb{R}^{M \times 1}$. For the optimal \mathbf{w}_j^* of problem $\tilde{\mathcal{F}}_{\text{lb},t}$, the optimal precoding vector at the j th BS maximizing the minimum weighted rate lower bound can be expressed as

$$\mathbf{p}_j^* = \sum_{m \in \mathcal{T}_j} \sqrt{w_{j,m}^*} \mathbf{v}_m \quad (32)$$

where $w_{j,m}^*$ is the m th entry of \mathbf{w}_j^* and \mathbf{v}_m is the m th column of matrix \mathbf{V} . The set \mathcal{T}_j is defined as

$$\mathcal{T}_j = \left\{ m \mid \sum_k \tilde{r}_{j,j,k,m} \neq 0, \nexists m' : (m, m') \text{ satisfies (34)} \right\} \quad (33)$$

where $\tilde{r}_{i,j,k,m}$ is the m th entry of $\tilde{\mathbf{r}}_{i,j,k}$ and (34) is given by

$$\begin{cases} \tilde{r}_{j,j,k,m} < \tilde{r}_{j,j,k,m'}, \forall k \in \mathcal{G}_j \\ \tilde{r}_{j,i,k,m} > \tilde{r}_{j,i,k,m'}, \forall k \in \mathcal{G}_i, \forall i \neq j. \end{cases} \quad (34)$$

Proof: See Appendix D. ■

From Proposition 4, the optimal precoding vectors, which maximize the minimum weighted rate lower bound among the KL served users, are linear combinations of the columns of a unique and deterministic matrix \mathbf{V} . Note that the columns of \mathbf{V} are the eigenvectors of different user's channel correlation matrix, which only hinge on the topology of the BS antenna array. Recalling the beam domain transmission proposed in [11], each column of \mathbf{V} is corresponding to a beam direction, which means that vectors $\tilde{\mathbf{r}}_{j,j,k}$ represent the beam gains and the entries of \mathbf{w}_j can be regarded as the power allocated to the corresponding beams. Furthermore, owing to the channel sparsity in massive MIMO systems, most elements in $\tilde{\mathbf{r}}_{j,j,k}$ are approximately zero [38]. Thus, when the number of users in each cell is relatively small, the value of $\sum_{k=1}^K \tilde{r}_{j,j,k,m}$ will approximate to zero for most m . For this case, the size of the set \mathcal{T}_j , i.e., $M_j = |\mathcal{T}_j|$, is much smaller than M , and we only need to optimize the vectors \mathbf{w}_j over a lower dimensional real space $M_j \times 1$, which implies that the beam domain transmission can maximize the minimum weighted rate lower bound with lower computational complexity.

Next, for problem $\tilde{\mathcal{F}}_{\text{lb},t}$, we utilize the duality between the MMF and QoS problem to obtain the optimal solution. Introducing an additional constraint $\sum_j \mathbf{1}^T \mathbf{w}_j \leq P_T$ and an auxiliary variable $s > 0$, we formulate the following

problem

$$\begin{aligned}
\tilde{\mathcal{F}}_{\text{lb},t,s} : \quad & \max_{\{\mathbf{w}_j\}_{j=1}^L, t, s \geq 0} t \\
\text{s.t.} \quad & \rho_{j,k} t \left(\sum_{i \neq j} \mathbf{w}_i^T \tilde{\mathbf{r}}_{i,j,k} + 1 \right) \\
& - \mathbf{w}_j^T \tilde{\mathbf{r}}_{j,j,k} \leq 0, \quad \forall k \in \mathcal{G}_j, \quad \forall j \in \mathcal{L} \\
& \mathbf{1}^T \mathbf{w}_j \leq P_j, \quad \forall j \in \mathcal{L} \\
& \mathbf{w}_j \succeq \mathbf{0}, \quad \forall j \in \mathcal{L} \\
& \sum_j \mathbf{1}^T \mathbf{w}_j + s = P_T. \tag{35}
\end{aligned}$$

Note that, if $P_T \geq \sum_j P_j$, the optimal solution of $\tilde{\mathcal{F}}_{\text{lb},t,s}$ is equal to that of $\tilde{\mathcal{F}}_{\text{lb},t}$, where the constraint $\sum_j \mathbf{1}^T \mathbf{w}_j + s = P_T$ can be ignored. Then, we consider the QoS precoding design problem which aims to minimize the total transmit power, while satisfying per-cell power constraints and the target $\gamma_{j,k}$ defined in (18). The QoS problem can be written as

$$\begin{aligned}
\tilde{\mathcal{Q}}_{\text{lb}} : \quad & \min_{\{\mathbf{w}_j\}_{j=1}^L} \sum_j \mathbf{1}^T \mathbf{w}_j \\
\text{s.t.} \quad & \rho_{j,k} \left(\sum_{i \neq j} \mathbf{w}_i^T \tilde{\mathbf{r}}_{i,j,k} + 1 \right) \\
& - \mathbf{w}_j^T \tilde{\mathbf{r}}_{j,j,k} \leq 0, \quad \forall k \in \mathcal{G}_j, \quad \forall j \in \mathcal{L} \\
& \mathbf{1}^T \mathbf{w}_j \leq P_j, \quad \forall j \in \mathcal{L} \\
& \mathbf{w}_j \succeq \mathbf{0}, \quad \forall j \in \mathcal{L}. \tag{36}
\end{aligned}$$

Find that the QoS problem $\tilde{\mathcal{Q}}_{\text{lb}}$ is a standard semidefinite programming (SDP) problem, and it can be tackled by classic convex optimization methods. Hence, we attempt to solve problem $\tilde{\mathcal{F}}_{\text{lb},t,s}$ by utilizing the duality between the QoS problem and the MMF problem. For concise representation, we define the maximum value of t in $\tilde{\mathcal{F}}_{\text{lb},t,s}$ as $t^* = \tilde{\mathcal{F}}_{\text{lb},t,s}(\boldsymbol{\kappa}, \boldsymbol{\rho}, P_T)$ and the minimum total transmitted power in $\tilde{\mathcal{Q}}_{\text{lb}}$ as $P_m^* = \tilde{\mathcal{Q}}_{\text{lb}}(\boldsymbol{\kappa}, \boldsymbol{\rho})$, where $\boldsymbol{\kappa} = [P_1, \dots, P_L]$ and $\boldsymbol{\rho} = [\rho_{1,1}, \dots, \rho_{L,K}]$. From [5], we have the following relations

$$t^* = \tilde{\mathcal{F}}_{\text{lb},t,s}(\boldsymbol{\kappa}, \boldsymbol{\rho}, \tilde{\mathcal{Q}}_{\text{lb}}(\boldsymbol{\kappa}, t^* \boldsymbol{\rho}) + s^*) \tag{37}$$

$$P_T - s^* = \tilde{\mathcal{Q}}_{\text{lb}}(\boldsymbol{\kappa}, \tilde{\mathcal{F}}_{\text{lb},t,s}(\boldsymbol{\kappa}, \boldsymbol{\rho}, P_T) \boldsymbol{\rho}) \tag{38}$$

where s^* is the optimal value of s in $\tilde{\mathcal{F}}_{\text{lb},t,s}(\boldsymbol{\kappa}, \boldsymbol{\rho}, P_T)$.

Based on the relations in (37) and (38), for given $\boldsymbol{\kappa}$, $\boldsymbol{\rho}$ and P_T , the optimal objective values of $\tilde{\mathcal{Q}}_{\text{lb}}(\boldsymbol{\kappa}, t \boldsymbol{\rho})$ and $\tilde{\mathcal{F}}_{\text{lb},t,s}(\boldsymbol{\kappa}, \boldsymbol{\rho}, P_T)$ are monotonically nondecreasing with t . Therefore, iteratively solving the QoS problem $\tilde{\mathcal{Q}}_{\text{lb}}(\boldsymbol{\kappa}, t \boldsymbol{\rho})$ with bisection search over t , we can find the optimal solution of problem $\tilde{\mathcal{F}}_{\text{lb},t,s}(\boldsymbol{\kappa}, \boldsymbol{\rho}, P_T)$, which is actually the optimal solution of problem $\tilde{\mathcal{F}}_{\text{lb},t}$ for $P_T \geq \sum_j P_j$. The detailed steps of the proposed relation-based algorithm are summarized in Algorithm 3.

For solving the SDP problem $\tilde{\mathcal{Q}}_{\text{lb}}(\boldsymbol{\kappa}, t \boldsymbol{\rho})$ in step 3 of Algorithm 3, we still utilize the barrier method in this paper. The

Algorithm 3: Relation-based Precoding Design.

1: Initialization. Decide the sets \mathcal{T}_j . Initialize $t_{\text{lb}} = 0$ and the upper bound of t as

$$t_{\text{ub}} = \max_{j,k,m} \frac{P_j \tilde{\mathbf{r}}_{j,j,k,m}}{\rho_{j,k}}.$$

2: **repeat**

3: Update $t = \frac{t_{\text{lb}} + t_{\text{ub}}}{2}$ and solve the QoS problem $\tilde{\mathcal{Q}}_{\text{lb}}(\boldsymbol{\kappa}, t \boldsymbol{\rho})$ via Algorithm 4 to obtain vectors \mathbf{w}_j .

4: If problem $\tilde{\mathcal{Q}}_{\text{lb}}(\boldsymbol{\kappa}, t \boldsymbol{\rho})$ is infeasible, set $t_{\text{ub}} = t$; otherwise set $t_{\text{lb}} = t$.

5: **until** $t_{\text{ub}} - t_{\text{lb}} \leq \epsilon$ (ϵ is a predefined tolerance.)

6: Obtain solution. Calculate $\mathbf{p}_j = \sum_{m \in \mathcal{T}_j} \sqrt{w_{j,m}} \mathbf{v}_m$.

corresponding unconstrained problem is given as follow

$$\begin{aligned}
\min_{\{\mathbf{w}_j\}_{j=1}^L} \tilde{h}(\{\mathbf{w}_j\}_{j=1}^L) = & \alpha \sum_j \mathbf{1}^T \mathbf{w}_j + \sum_j \psi_j(\mathbf{w}_j) \\
& + \sum_{j,k} \tilde{\varphi}_{j,k}(\{\mathbf{w}_i\}_{i=1}^L) \tag{39}
\end{aligned}$$

where $\psi_j(\mathbf{w}_j)$ and $\tilde{\varphi}_{j,k}(\{\mathbf{w}_i\}_{i=1}^L)$ are the log barrier functions corresponding to the inequality constraints in $\tilde{\mathcal{Q}}_{\text{lb}}(\boldsymbol{\kappa}, t \boldsymbol{\rho})$, which can be expressed as, respectively,

$$\psi_j(\mathbf{w}_j) = \begin{cases} +\infty, & \mathbf{1}^T \mathbf{w}_j \geq P_j \text{ or } \mathbf{w}_j \prec \mathbf{0} \\ -\log(P_j - \mathbf{1}^T \mathbf{w}_j) - \sum_m \log(w_{j,m}), & \text{else} \end{cases} \tag{40}$$

and

$$\tilde{\varphi}_{j,k}(\{\mathbf{w}_i\}_{i=1}^L) = \begin{cases} +\infty, & \tilde{q}_{j,k}(\{\mathbf{w}_i\}_{i=1}^L) \geq 0 \\ -\log(-\tilde{q}_{j,k}(\{\mathbf{w}_i\}_{i=1}^L)), & \text{else} \end{cases} \tag{41}$$

where function $\tilde{q}_{j,k}(\{\mathbf{w}_i\}_{i=1}^L)$ is defined as

$$\tilde{q}_{j,k}(\{\mathbf{w}_i\}_{i=1}^L) = t \rho_{j,k} \left(\sum_{i \neq j} \mathbf{w}_i^T \tilde{\mathbf{r}}_{i,j,k} + 1 \right) - \mathbf{w}_j^T \tilde{\mathbf{r}}_{j,j,k}. \tag{42}$$

Then, the partial derivative of $\tilde{h}(\{\mathbf{w}_i\}_{i=1}^L)$ with respect to \mathbf{w}_j is given by

$$\begin{aligned}
\frac{\partial \tilde{h}(\{\mathbf{w}_l\}_{l=1}^L)}{\partial \mathbf{w}_j} = & \alpha \mathbf{1} + \frac{\mathbf{1}}{P_j - \mathbf{1}^T \mathbf{w}_j} - \mathbf{w}_j^\dagger + \sum_{k=1}^K \frac{\tilde{\mathbf{r}}_{j,j,k}}{\tilde{q}_{j,k}(\{\mathbf{w}_l\}_{l=1}^L)} \\
& - \sum_{i \neq j} \sum_{k=1}^K \frac{t \rho_{i,k} \tilde{\mathbf{r}}_{j,i,k}}{\tilde{q}_{i,k}(\{\mathbf{w}_l\}_{l=1}^L)} \tag{43}
\end{aligned}$$

where $\mathbf{w}_j^\dagger \triangleq [w_{j,1}^{-1}, w_{j,2}^{-1}, \dots, w_{j,M}^{-1}]^T$.

We note that solving $\tilde{\mathcal{Q}}_{\text{lb}}(\boldsymbol{\kappa}, t \boldsymbol{\rho})$ with the barrier method requires a feasible starting point. However, with inappropriate value of t , problem $\tilde{\mathcal{Q}}_{\text{lb}}(\boldsymbol{\kappa}, t \boldsymbol{\rho})$ will be infeasible (the value of t is too large), and we should update $t_{\text{ub}} = t$ in Algorithm 3. Thus, before solving (39), we have to find a feasible solution of problem $\tilde{\mathcal{Q}}_{\text{lb}}(\boldsymbol{\kappa}, t \boldsymbol{\rho})$ or determine that none exists. To do this,

Algorithm 4: QoS Precoding Design.

-
- 1: Initialization. Solve problem \tilde{Q}_{ini} by barrier method.
During the procedure of solving \tilde{Q}_{ini} , once $(\mathbf{w}_1^{\text{ini}}, \dots, \mathbf{w}_L^{\text{ini}}, u)$ is feasible for \tilde{Q}_{ini} with $u \leq 0$, choose starting point $\mathbf{w}_j = \mathbf{w}_j^{\text{ini}}$. Otherwise, $\tilde{Q}_{\text{lb}}(\boldsymbol{\kappa}, t\rho)$ is infeasible, and terminate Algorithm 4.
 - 2: Set $\alpha := \alpha^{(0)} > 0$ and $c > 1$.
 - 3: **repeat**
 - 4: Set $u := 0$, $\mathbf{x}_j^{(0)} = \mathbf{w}_j$.
 - 5: **repeat**
 - 6: Determine step size μ by backtracking line search.
 - 7: Update $\mathbf{x}_j^{(u+1)} = \mathbf{x}_j^{(u)} + \mu \frac{\partial \tilde{h}(\{\mathbf{x}_i\}_{i=1}^L)}{\partial \mathbf{x}_j}$ and calculate increment $\tau = |\tilde{h}(\{\mathbf{x}_i^{(u+1)}\}_{i=1}^L) - \tilde{h}(\{\mathbf{x}_i^{(u)}\}_{i=1}^L)|$.
 - 8: Set $u = u + 1$.
 - 9: **until** $\tau < \epsilon_1$ (ϵ_1 is a predefined tolerance.)
 - 10: Update $\mathbf{w}_j = \mathbf{x}_j^{(u)}$ and $\alpha = c\alpha$.
 - 11: **until** $\frac{1}{\alpha} < \epsilon_2$ (ϵ_2 is a predefined tolerance.)
-

we consider the following problem

$$\begin{aligned}
 \tilde{Q}_{\text{ini}} : \quad & \min_{\{\mathbf{w}_j\}_{j=1}^L, u} u \\
 \text{s.t.} \quad & \tilde{q}_{j,k}(\{\mathbf{w}_i\}_{i=1}^L) \leq u, \forall j \in \mathcal{L}, \forall k \in \mathcal{G}_j \\
 & \mathbf{1}^T \mathbf{w}_j \leq P_j, \forall j \in \mathcal{L} \\
 & \mathbf{w}_j \geq \mathbf{0}, \forall j \in \mathcal{L}.
 \end{aligned} \tag{44}$$

The value of u can be interpreted as a bound on the infeasibility of inequalities $\tilde{q}_{j,k}(\{\mathbf{w}_i\}_{i=1}^L) \leq 0$. Specifically, if the optimal value $u^* \leq 0$, then $\tilde{Q}_{\text{lb}}(\boldsymbol{\kappa}, t\rho)$ has a feasible solution. Otherwise, if the optimal value $u^* > 0$, problem $\tilde{Q}_{\text{lb}}(\boldsymbol{\kappa}, t\rho)$ is infeasible. Since we can initialize $u = \max_{j,k} \tilde{q}_{j,k}(\{\mathbf{w}_i\}_{i=1}^L)$, problem

\tilde{Q}_{ini} is always feasible. Then, the barrier method can be applied to solve problem \tilde{Q}_{ini} . Notice that we do not need to solve \tilde{Q}_{ini} with high accuracy, and we can terminate when $u \leq 0$. Since if $(\mathbf{w}_1^{\text{ini}}, \dots, \mathbf{w}_L^{\text{ini}}, u)$ is feasible for \tilde{Q}_{ini} with $u < 0$, then $(\mathbf{w}_1^{\text{ini}}, \dots, \mathbf{w}_L^{\text{ini}})$ is a feasible starting point for $\tilde{Q}_{\text{lb}}(\boldsymbol{\kappa}, t\rho)$. The steps to solve problem $\tilde{Q}_{\text{lb}}(\boldsymbol{\kappa}, t\rho)$ are summarized in Algorithm 4.

Then, we discuss the computational complexities of the proposed algorithms. In each iteration, the main complexity of the CCCP-based algorithm comes from solving the QCQP problem, which is $O(LM + KL)^{3.5}$ [16]. For the relation-based algorithm, the complexity in each iteration is mainly contributed by the complexity of solving the SDP problem \tilde{Q}_{lb} . Because problem \tilde{Q}_{lb} is transformed into a lower dimensional space, where we need to optimize a $M_j \times 1$ real vector for the j th cell, $j = 1, \dots, L$. The complexity of one iteration in the relation-based algorithm will approximate to $O((\sum_j M_j)^{3.5} + KL(\sum_j M_j)^{1.5})$ [4]. Finally, we analyze the convergence of the proposed algorithms. For the outer iteration of the CCCP-based algorithm, Proposition 3 reveals that, utilizing the CCCP, the algorithm guarantee converge to a locally optimal solution. Considering its inner iteration, the corresponding convex QCQP problem is solved by the barrier method and the gradient descent method, where these methods have been proven

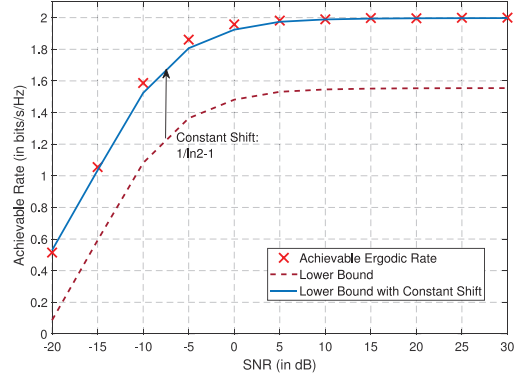


Fig. 2. Achievable ergodic rate and its lower bound.

to be convergent in [36]. As for the relation-base algorithm, since variable t is upper bounded and lower bounded by t_{ub} and t_{lb} , respectively, the bisection search over t is convergent. Meanwhile, the barrier method and the gradient descent method are also utilized to solve the convex inner iterative problem $\tilde{Q}_{\text{lb}}(\boldsymbol{\kappa}, t\rho)$, which makes the inner iteration of the relation-based algorithm convergent.

V. SIMULATION RESULTS

In this section, we present the simulation results to illustrate the performance of the developed algorithms. In our simulations, the WINNER II channel model is utilized to generate channel $\mathbf{h}_{i,j,k}$. For the WINNER II channel model, we consider the suburban scenario with the non-line-of-sight (NLOS) condition. Both the path loss and shadow fading model of WINNER II are utilized. For the topology of the cells, we consider a hexagonal grid of three-sectoral cells ($L = 3$), where the term cell refers to a 120° sector and each BS is located at the vertex of the hexagonal grid. K single-antenna users are uniformly distributed within each cell. Without loss of generality, we assume the same achievable ergodic rate lower bound target for all users, i.e., $\rho_{j,k} = 1, \forall j \in \mathcal{L}, \forall k \in \mathcal{G}_j$. In addition, users are uniformly distributed within each cell.

We note that, when the SNR approaches to 0 (i.e., $P_j \rightarrow 0$), the limits of the achievable ergodic rate $R_{j,k}$ in (7) and its lower bound $R_{\text{lb},j,k}$ in (16) are given by, respectively,

$$\lim_{\text{SNR} \rightarrow 0} R_{j,k} = 0 \tag{45}$$

$$\lim_{\text{SNR} \rightarrow 0} R_{\text{lb},j,k} = 1 - \frac{1}{\ln 2}. \tag{46}$$

Here, limits (45) and (46) indicate that the achievable ergodic rate and its lower bound has a constant gap in low SNR region. Because the objective function plus a constant value remains the optimized precoding vectors unchanged, we add a constant value $\frac{1}{\ln 2} - 1$ to the rate lower bound. Fig. 2 shows the achievable ergodic rate and the rate lower bound with a constant shift for one of the KL users, where the simulated curve of the achievable ergodic rate is obtained by the Monte-Carlo simulations. We consider the case of QPSK inputs with $M = 128$ and $K = 2$. Fig. 2 demonstrates that the rate lower bound with a constant shift is a good approximation to the achievable ergodic rate for QPSK inputs. By adding a constant, the rate lower bound are nearly identical to the achievable ergodic rate at low and high

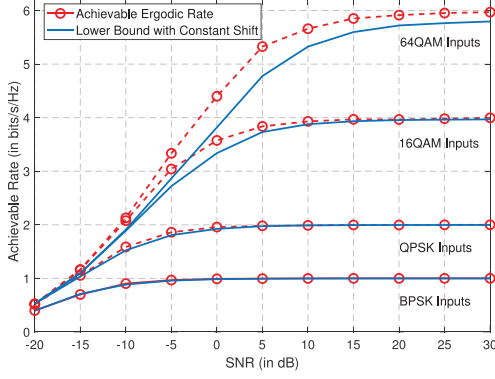


Fig. 3. Achievable ergodic rate and its lower bound with a constant shift for various input types.

SNR regions, and they approach to each other at medium SNR region. Besides, we find that the computational complexity of the rate lower bound is much lower than that of the achievable ergodic rate, where the cost of the Monte-Carlo averaging is avoided.

Fig. 3 further compares the achievable ergodic rate (7) and its lower bound with a constant shift for various input types. In low SNR region, the rate lower bound with a constant shift approximates to the achievable ergodic rate for all input types. As for medium and high SNR regions, the achievable ergodic rate is still well approximated by shifted rate lower bound for BPSK and QPSK input types. For the case of 16QAM inputs, the shifted rate lower bound and the achievable ergodic rate close to each other at high SNR region while there is a rate gap at medium SNR region, and for 64QAM inputs, the corresponding rate gap at medium and high SNR regions is larger than those of other input types in Fig. 3. This is because that, as the transmit power of each BS increases, inter-cell interferences dominate the tightness between the rate and its lower bound, especially for the input type with a large number of constellation points. Then, when SNR is high enough, we find that both the achievable ergodic rate and its lower bound will saturate.

Furthermore, in high SNR region, we can find that the achievable ergodic rate all saturate at the rate maximum $\log Q_j$ for various input types with the proposed precoding, which is a common result for the case of finite-alphabet inputs [27], [28]. However, the gap between the rate maximum $\log Q_j$ and the shifted rate lower bound increases with the increase of the order of input type. To interpret this result, we define $d_{j,m}^{\min} = \min_n |d_j^{m,n}|^2$, which indicates the minimum distance between the m th constellation point and the other constellation points. Particularly, for a equiprobable zero-mean constellation set, $d_j^{\min} \triangleq d_{j,1}^{\min} = \dots = d_{j,Q_j}^{\min}$. Then, we derive the following relations

$$\begin{aligned} \log Q_j - R_{lb,j,k} - \left(\frac{1}{\ln 2} - 1 \right) \\ &= \frac{1}{Q_j} \sum_{m=1}^{Q_j} \log_2 \sum_{n=1}^{Q_j} \left(1 + \frac{1}{2} |d_j^{m,n}|^2 \gamma_{j,k} \right)^{-1} \\ &\geq \frac{1}{Q_j} \sum_{m=1}^{Q_j} \log \left(1 + \frac{2}{2 + d_{j,m}^{\min} \gamma_{j,k}} \right) \end{aligned}$$

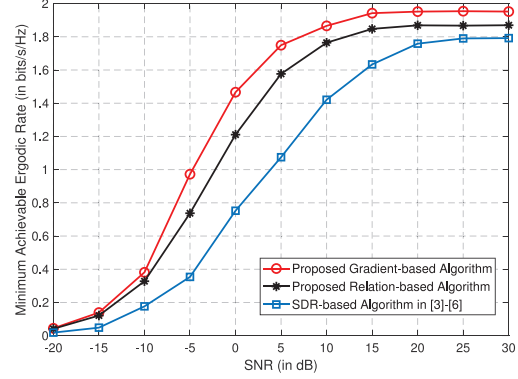


Fig. 4. Minimum achievable ergodic rate performance for different precoding algorithms.

$$\begin{aligned} &= \log \left(1 + \frac{2}{2 + d_j^{\min} \gamma_{j,k}} \right) \\ &\triangleq g(d_j^{\min}, \gamma_{j,k}). \end{aligned} \quad (47)$$

From (47), we note that the gap between the rate lower bound with a constant shift and the rate maximum $\log Q_j$ is lower bounded by $g(d_j^{\min}, \gamma_{j,k})$, where function $g(d_j^{\min}, \gamma_{j,k})$ monotonously increases with the decrease of d_j^{\min} for fixed $\gamma_{j,k}$. Then, for 64QAM inputs, whose corresponding d_j^{\min} is much smaller than those of other inputs types in Fig. 3, the gap between the rate lower bound with a constant shift and rate maximum $\log Q_j$ is larger than those of other inputs types in Fig. 3 at high SNR region.

Fig. 4 shows the minimum achievable ergodic rate for different precoding algorithms with $M = 64$, $K = 3$ and QPSK inputs. In Fig. 4, the SDR-based algorithm studied in [3]–[6] refers to an algorithm solving the following problem

$$\begin{aligned} &\max_{\{\mathbf{X}_j\}_{j=1}^L, t} t \\ &\text{s.t. } \rho_{j,k} t \left(\sum_{i \neq j} \text{tr}(\mathbf{X}_i \mathbf{R}_{i,j,k}) + 1 \right) \\ &\quad - \text{tr}(\mathbf{X}_j \mathbf{R}_{j,j,k}) \leq 0, \forall k \in \mathcal{G}_j, \forall j \in \mathcal{L} \\ &\quad \text{tr}(\mathbf{X}_j) \leq P_j, \forall j \in \mathcal{L} \\ &\quad \mathbf{X}_j \succeq \mathbf{0}, \forall j \in \mathcal{L} \end{aligned} \quad (48)$$

where $\mathbf{X}_j \triangleq \mathbf{p}_j \mathbf{p}_j^H$. From [3]–[6], the problem in (48) is transformed from the original problem $\mathcal{F}_{lb,t}$, while dropping the rank-1 constraint, i.e., $\text{rank}(\mathbf{X}_j) = 1, \forall j \in \mathcal{L}$. After obtaining the solution \mathbf{X}_j , the corresponding precoding vectors are generated by Gaussian randomization method with appropriate power scaling factors [5], [6]. Simulation results show that the developed gradient-based algorithm outperforms the other two algorithms, which is because that the gradient-based algorithm directly optimizes the precoding vectors with the achievable ergodic rate in (7). However, without closed-form expression for the gradients in (14) and (15), the gradient-based algorithm is computationally inefficient for high-order modulation. To reduce the computational complexity, we can adopt our proposed relation-based algorithm which avoids the Monte-Carlo averaging. As shown

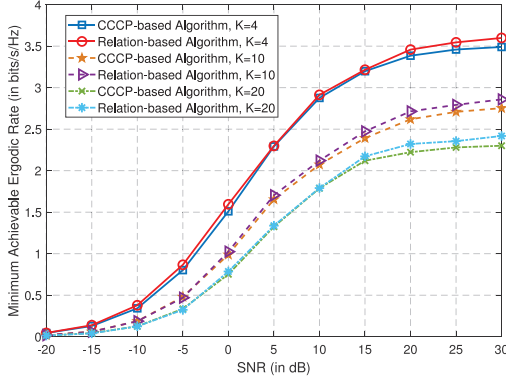


Fig. 5. Minimum achievable ergodic rate performance for the proposed CCCP-based algorithm and the proposed relation-based algorithm with different number of users in each cell.

in Fig. 4, the relation-based algorithm achieves slightly worse performance than that of the gradient-based algorithm, while outperforming the SDR-based algorithm. In addition, we note that the computational complexity of the SDR-based algorithm is much higher than that of the proposed relation-based algorithm, since the SDR-based algorithm lifts the original problem into higher dimensional space.

Fig. 5 shows the performance of the max-min achievable ergodic rate for the proposed CCCP-based algorithm and relation-based algorithm with different number of users. Here, we consider the 16QAM inputs with $M = 128$. For the number of users in each cell, we set $K = 4$, $K = 10$ and $K = 20$, respectively. We can find that, for both CCCP-based algorithm and relation-based algorithm, the minimum achievable ergodic rate of all KL users degrades as K increases. This is intuitive because the more users there are, the more likely users' channels are in deep fading condition, which results in the reduction of the minimum achievable ergodic rate. Moreover, as shown in Fig. 5, the CCCP-based algorithm and the relation-based algorithm can achieve nearly the same minimum achievable ergodic rate at low and medium SNR regions, and the relation-based perform slightly better when the SNR is high. It is because that the precoding vectors generated from the CCCP-based algorithm guarantee to converge to a locally optimal solution of the problem $\mathcal{F}_{lb,t}$, while the relation-based algorithm can generate the optimal solution of the problem $\mathcal{F}_{lb,t}$ with perfect bisection search of parameter t . Nevertheless, it should be pointed out that the relation-based algorithm is founded on the assumption that the number of BS antennas goes to infinity, without which the CCCP-based algorithm can still work.

To further illustrate the performance of the CCCP-based algorithm and the relation-based algorithm. Fig. 6 compares the max-min rate lower bound performance of the CCCP-based algorithm and the relation-based algorithm for different number of BS antennas. Considering 16QAM inputs, we set SNR = 20 dB and $K = 8$. From Fig. 6, we can find that the CCCP-based algorithm achieves a higher minimum rate lower bound than that of the relation-based algorithm when M is relatively small. However, as M increases, the minimum rate lower bound performance of the relation-based algorithm become better than

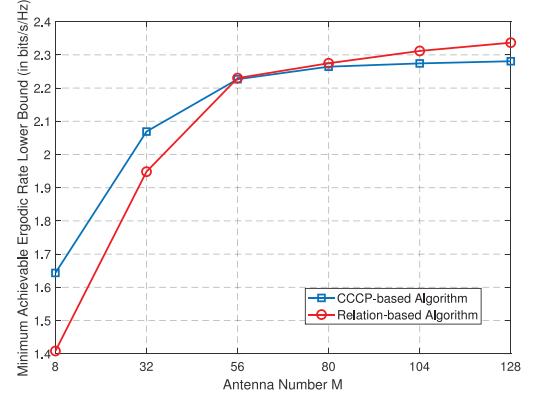


Fig. 6. Comparison of the rate lower bound between the proposed CCCP-based algorithm and the proposed relation-based algorithm with different number of BS antennas.

that of the CCCP-based algorithm. Notice that the optimality of the relation-based algorithm is founded on the channel characteristic, where the eigenmatrices of different user's transmit correlation matrices tend to be a unique deterministic matrix \mathbf{V} as $M \rightarrow \infty$. When the number of BS antennas is relatively small, the users' channel may not satisfy this characteristic. Therefore, the CCCP-based algorithm can achieve a better performance, where it can obtain a locally optimal solution without this characteristic. Then, as each BS is equipped with a large number of antennas, where the channel characteristic is nearly satisfied, the relation-based algorithm outperforms the CCCP-based algorithm.

Meanwhile, as illustrated in Fig. 6, increasing the number of BS antennas, the performance of both the CCCP-based algorithm and the relation-based algorithm can be improved. It is because that increasing the number of BS antennas allows the resolution of more paths, resulting in less inter-cell interference. Then, when the number of BS antennas M is large enough, the gain obtained from the increase of the number of antennas will be smaller, since the inter-cell interference has been substantially canceled.

VI. CONCLUSION

We investigated multicell massive MIMO multicast transmission with finite-alphabet signals, where only statistical CSI of users is available at the BSs. For finite-alphabet inputs, the necessary conditions of the optimal precoding vectors to maximize the minimum weighted achievable ergodic rate among all users were obtained, assuming that the interference is Gaussian noise. Based on these conditions, we proposed the gradient-based algorithm to iteratively search the optimal precoding vectors. Then, owing to the high computational burden of the achievable ergodic rate with finite-alphabet inputs, we derived a lower bound on the achievable ergodic rate, which depends on the transmit correlation matrices. With this lower bound, we investigated the corresponding weighted MMF problem and devised an efficient iterative algorithm for precoding design by exploiting the CCCP method, where the precoding vectors generated from the CCCP-based algorithm can be proven to converge to a local optimum. Moreover, in massive MIMO systems, the eigenmatrices of the

$$\varepsilon_{2,j,i,k} = \frac{1}{\ln 2 \cdot Q_i} \sum_{m=1}^{Q_i} \frac{\sum_{n=1}^{Q_i} \exp(-f_{m,n,i,k}) \left(|\mathbf{h}_{i,i,k}^H \mathbf{p}_i d_i^{m,n} + z'_{i,k}|^2 - |z'_{i,k}|^2 - \Re\{(z'_{i,k})^* \mathbf{h}_{i,i,k}^H \mathbf{p}_i d_i^{m,n}\} \right) \mathbf{p}_j^H \mathbf{h}_{j,i,k}}{\left(\sum_{i' \neq i} |\mathbf{h}_{i',i,k}^H \mathbf{p}_{i'}|^2 + 1 \right)^2 \cdot \left(\sum_{n=1}^{Q_i} \exp(-f_{m,n,i,k}) \right)} \quad (55)$$

transmit correlation matrices become an identical deterministic unitary matrix as the number of BS antennas goes to infinity. Utilizing this characteristic, we proved that the optimal precoding vectors should be linear combinations of eigenvectors of transmit correlation matrices, where the weighted MMF problem can be further simplified. For obtaining the optimal solution of the simplified weighted MMF problem, we proposed an iterative algorithm by using the duality between the QoS problem and the MMF problem. The tightness of the achievable ergodic rate lower bound and the significant performance of the devised algorithms were illustrated through simulation results.

APPENDIX A PROOF OF PROPOSITION 1

To solve problem \mathcal{F}_t , we develop the cost function of the optimal precoders as

$$\mathcal{C} = t + \sum_{j,k} \zeta_{j,k} (R_{j,k} - \theta_{j,k} t) - \sum_{j=1}^L \lambda_j (\|\mathbf{p}_j\|^2 - P_j) \quad (49)$$

where $\zeta_{j,k} \geq 0$ and $\lambda_j \geq 0$ are Lagrange multipliers associated with the inequality constraints $R_{j,k} \geq \theta_{j,k} t$ and $\|\mathbf{p}_j\|^2 \leq P_j$, respectively. The Karush-Kuhn-Tucker (KKT) conditions enable us to establish the equations for the optimal precoders \mathbf{p}_j , $\forall j \in \mathcal{L}$, which can be expressed as

$$\frac{\partial \mathcal{C}}{\partial t} = 1 - \sum_{j,k} \zeta_{j,k} \theta_{j,k} = 0 \quad (50a)$$

$$\frac{\partial \mathcal{C}}{\partial \mathbf{p}_j} = \sum_{i,k} \zeta_{i,k} \frac{\partial R_{i,k}}{\partial \mathbf{p}_j} - \lambda_j \mathbf{p}_j^* = \mathbf{0}, \quad \forall j \in \mathcal{L} \quad (50b)$$

$$\zeta_{j,k} (R_{j,k} - \theta_{j,k} t) = 0, \quad R_{j,k} \geq \theta_{j,k} t, \quad \forall j \in \mathcal{L}, \forall k \in \mathcal{G}_j \quad (50c)$$

$$\lambda_j (\|\mathbf{p}_j\|^2 - P_j) = 0, \quad \|\mathbf{p}_j\|^2 \leq P_j, \quad \forall j \in \mathcal{L}. \quad (50d)$$

In (50b), the calculation of the derivative $\sum_{i,k} \frac{\partial R_{i,k}}{\partial \mathbf{p}_j}$ includes two parts. For $i = j$

$$\frac{\partial R_{j,k}}{\partial \mathbf{p}_j} = \mathbb{E}\{\varepsilon_{1,j,k} \mathbf{h}_{j,j,k}^* \mathbf{p}_j\} \quad (51)$$

where the random variable $\varepsilon_{1,j,k}$ given as (52),

$$\varepsilon_{1,j,k} = \frac{1}{\ln 2 \cdot Q_j} \sum_{m=1}^{Q_j} \frac{\sum_{n=1}^{Q_j} \exp(-f_{m,n,j,k}) \left(|d_j^{m,n}|^2 \mathbf{p}_j^H \mathbf{h}_{j,j,k} + d_j^{m,n} (z'_{j,k})^* \right)}{\left(\sum_{i \neq j} |\mathbf{h}_{i,j,k}^H \mathbf{p}_i|^2 + 1 \right) \cdot \left(\sum_{n=1}^{Q_j} \exp(-f_{m,n,j,k}) \right)} \quad (52)$$

In (52), $f_{m,n,j,k}$ is defined as

$$f_{m,n,j,k} = \frac{|\mathbf{h}_{j,j,k}^H \mathbf{p}_j d_j^{m,n} + z'_{j,k}|^2}{\sum_{i \neq j} |\mathbf{h}_{i,j,k}^H \mathbf{p}_i|^2 + 1}. \quad (53)$$

For the case of $i \neq j$, we note that $z'_{i,k}$ is a zero-mean Gaussian noise with variance $\sum_{i' \neq i} |\mathbf{h}_{i',i,k}^H \mathbf{p}_{i'}|^2 + 1$, and the corresponding derivatives of $R_{i,k}$ is given as

$$\frac{\partial R_{i,k}}{\partial \mathbf{p}_j} = -\mathbb{E}\{\varepsilon_{2,j,i,k} \mathbf{h}_{j,i,k}^* \mathbf{p}_j\} \quad (54)$$

where $\varepsilon_{2,j,i,k}$ given by (55), shown at the top of this page. Thus, condition (50b) can be expressed as

$$\sum_{k=1}^K \zeta_{j,k} \mathbb{E}\{\varepsilon_{1,j,k} \mathbf{h}_{j,j,k}^* \mathbf{p}_j\} - \sum_{i \neq j} \sum_{k=1}^K \zeta_{i,k} \mathbb{E}\{\varepsilon_{2,j,i,k} \mathbf{h}_{j,i,k}^* \mathbf{p}_j\} - \lambda_j \mathbf{p}_j^* = \mathbf{0}. \quad (56)$$

APPENDIX B PROOF OF PROPOSITION 2

For the k th user in the j th cell, the corresponding achievable ergodic rate $R_{j,k}$ can be rewritten as

$$R_{j,k} = \mathbb{E}\{I_{j,k}(\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_L)\} \quad (57)$$

where $I_{j,k}(\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_L)$ defined by (58).

$$I_{j,k}(\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_L) = \log Q_j - \frac{1}{Q_j} \sum_{m=1}^{Q_j} \mathbb{E}_{\mathbf{h}_{j,j,k}, z'_{j,k}} \left\{ \log \sum_{n=1}^{Q_j} \exp \left(- \frac{|\mathbf{h}_{j,j,k}^H \mathbf{p}_j d_j^{m,n} + z'_{j,k}|^2 - |z'_{j,k}|^2}{\sum_{i \neq j} |\mathbf{h}_{i,j,k}^H \mathbf{p}_i|^2 + 1} \right) \right\} \quad (58)$$

Using the similarly method in [27], it is easy to find that $I_{j,k}(\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_L)$ is lower bounded by

$$I_{j,k}(\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_L) \geq \log Q_j - (1/\ln 2 - 1) - 1/Q_j \times \sum_{m=1}^{Q_j} \log \sum_{n=1}^{Q_j} \left(1 + \frac{1}{2} \frac{|d_j^{m,n}|^2 \mathbf{p}_j^H \mathbf{R}_{j,j,k} \mathbf{p}_j}{\sum_{i \neq j} |\mathbf{h}_{i,j,k}^H \mathbf{p}_i|^2 + 1} \right)^{-1}. \quad (59)$$

Then, we consider function $\log \sum_i (1 + \frac{b_i}{x})^{-1}$ with respect to $x > 0$, where $b_i \geq 0$. We have

$$\frac{d^2}{dx^2} \log \sum_i \left(1 + \frac{b_i}{x} \right)^{-1} \leq 0. \quad (60)$$

Therefore, function $\log \sum_i (1 + \frac{b_i}{x})^{-1}$ is a concave function with respect to $x > 0$ and the Jensen's inequality can be adopted to obtain the lower bound on $R_{j,k}$. Combining (57) and (59), we

can obtain the following relation by invoking Jensen's inequality

$$\begin{aligned}
R_{j,k} &\geq \log Q_j - (1/\ln 2 - 1) - 1/Q_j \\
&\times \sum_{m=1}^{Q_j} \log \sum_{n=1}^{Q_j} \left(1 + \frac{1}{2} \frac{|d_j^{m,n}|^2 \mathbf{p}_j^H \mathbf{R}_{j,j,k} \mathbf{p}_j}{\sum_{i \neq j} |\mathbf{h}_{i,j,k}^H \mathbf{p}_i|^2} + 1 \right)^{-1} \\
&= \log Q_j - (1/\ln 2 - 1) - 1/Q_j \\
&\times \sum_{m=1}^{Q_j} \log \sum_{n=1}^{Q_j} \left(1 + \frac{1}{2} \frac{|d_j^{m,n}|^2 \mathbf{p}_j^H \mathbf{R}_{j,j,k} \mathbf{p}_j}{\sum_{i \neq j} \mathbf{p}_i^H \mathbf{R}_{i,j,k} \mathbf{p}_i + 1} \right)^{-1}. \quad (61)
\end{aligned}$$

APPENDIX C PROOF OF PROPOSITION 3

Due to the convexity of function $g_{j,k}(\mathbf{s}_j)$, for the n th iteration result $(\{\mathbf{p}_i^{(n)}\}_{i \neq j}, \mathbf{s}_j^{(n)})$ in problem $\hat{\mathcal{F}}_{\text{lb}, \text{iter}, t}$, we have

$$g_{j,k}(\mathbf{s}_j^{(n)}) \geq \bar{g}_{j,k}(\mathbf{s}_j^{(n-1)}, \mathbf{s}_j^{(n)}). \quad (62)$$

Then, the following relations hold

$$\begin{aligned}
&f_{j,k}(\{\mathbf{p}_i^{(n)}\}_{i \neq j}) - \bar{g}_{j,k}(\mathbf{s}_j^{(n)}, \mathbf{s}_j^{(n)}) \\
&= f_{j,k}(\{\mathbf{p}_i^{(n)}\}_{i \neq j}) - g_{j,k}(\mathbf{s}_j^{(n)}) \\
&\leq f_{j,k}(\{\mathbf{p}_i^{(n)}\}_{i \neq j}) - \bar{g}_{j,k}(\mathbf{s}_j^{(n-1)}, \mathbf{s}_j^{(n)}) \\
&\leq 0. \quad (63)
\end{aligned}$$

From (63), the n th iteration results $(\{\mathbf{p}_i^{(n)}\}_{i \neq j}, \mathbf{s}_j^{(n)})$ is also a feasible solution for the $(n+1)$ th iteration, which means that $\{t^{(n)}\}_{n=1}^{\infty}$ is monotonic. Moreover, the set of $\{\mathbf{p}_j\}_{j=1}^L$ is closed and bounded, while the value of $\{t^{(n)}\}_{n=1}^{\infty}$ is upper bounded to satisfy the constraints. Consequently, the sequence $\{\{\mathbf{p}_i^{(n)}\}_{i \neq j}, \mathbf{s}_j^{(n)}\}_{n=1}^{\infty}$ generated from problem $\hat{\mathcal{F}}_{\text{lb}, \text{iter}, t}$ will converge [39], and we have

$$\begin{aligned}
\lim_{n \rightarrow \infty} f_{j,k}(\{\mathbf{p}_i^{(n)}\}_{i \neq j}) - g_{j,k}(\mathbf{s}_j^{(n)}) \\
= f_{j,k}(\{\mathbf{p}_i^*\}_{i \neq j}) - g_{j,k}(\mathbf{s}_j^*) \quad (64)
\end{aligned}$$

where $(\{\mathbf{p}_i^*\}_{i \neq j}, \mathbf{s}_j^*)$ is the generalized fixed point of CCCP. For problem $\hat{\mathcal{F}}_{\text{lb}, t}$, it is easy to find that $(\{\mathbf{p}_i^*\}_{i \neq j}, \mathbf{s}_j^*)$ satisfies the KKT conditions of $\hat{\mathcal{F}}_{\text{lb}, t}$ [40], resulting in a stationary point of problem $\hat{\mathcal{F}}_{\text{lb}, t}$.

APPENDIX D PROOF OF PROPOSITION 4

Recalling (4), in massive MIMO systems, we have

$$\mathbf{R}_{i,j,k} = \mathbf{V} \tilde{\mathbf{R}}_{i,j,k} \mathbf{V}^H, \quad \forall i, j \in \mathcal{L}, \quad \forall k \in \mathcal{G}_j. \quad (65)$$

Thus, by defining vectors $\tilde{\mathbf{p}}_j \triangleq \mathbf{V}^H \mathbf{p}_j$, problem $\mathcal{F}_{\text{lb}, t}$ can be rewritten as

$$\begin{aligned}
&\max_{\{\tilde{\mathbf{p}}_j\}_{j=1}^L, t} t \\
&\text{s.t. } \tilde{\mathbf{p}}_j^H \tilde{\mathbf{R}}_{j,j,k} \tilde{\mathbf{p}}_j - \rho_{j,k} t \left(\sum_{i \neq j} \tilde{\mathbf{p}}_i^H \tilde{\mathbf{R}}_{i,j,k} \tilde{\mathbf{p}}_i + 1 \right) \geq 0, \\
&\quad \forall j \in \mathcal{L}, \quad \forall k \in \mathcal{G}_j \\
&\|\tilde{\mathbf{p}}_j\|^2 \leq P_j, \quad \forall j \in \mathcal{L}. \quad (66)
\end{aligned}$$

Observing problem (66), where $\tilde{\mathbf{R}}_{i,j,k}$ are diagonal matrices, problem $\mathcal{F}_{\text{lb}, t}$ can be further equivalently reformulated as

$$\begin{aligned}
&\tilde{\mathcal{F}}_{\text{lb}, t} : \max_{\{\mathbf{w}_j\}_{j=1}^L, t} t \\
&\text{s.t. } \rho_{j,k} t \left(\sum_{i \neq j} \mathbf{w}_i^T \tilde{\mathbf{r}}_{i,j,k} + 1 \right) \\
&\quad - \mathbf{w}_j^T \tilde{\mathbf{r}}_{j,j,k} \leq 0, \quad \forall k \in \mathcal{G}_j, \quad \forall j \in \mathcal{L} \\
&\mathbf{1}^T \mathbf{w}_j \leq P_j, \quad \forall j \in \mathcal{L} \\
&\mathbf{w}_j \succeq \mathbf{0}, \quad \forall j \in \mathcal{L} \quad (67)
\end{aligned}$$

where $\mathbf{w}_j \triangleq \tilde{\mathbf{p}}_j \odot \tilde{\mathbf{p}}_j^*$ and $\tilde{\mathbf{r}}_{i,j,k} \triangleq \text{diag}(\tilde{\mathbf{R}}_{i,j,k})$. Then, for the optimal \mathbf{w}_j^* of problem $\tilde{\mathcal{F}}_{\text{lb}, t}^2$, the optimal precoding vectors \mathbf{p}_j^* can be expressed as

$$\mathbf{p}_j^* = \sum_{m=1}^M \sqrt{w_{j,m}^*} \mathbf{v}_m \quad (68)$$

where $w_{j,m}^*$ is the m th entry of \mathbf{w}_j^* and \mathbf{v}_m is the m th column of matrix \mathbf{V} .

Then, letting the m th entry of $\tilde{\mathbf{r}}_{i,j,k}$ be $\tilde{r}_{i,j,k,m}$, we investigate the conditions such that $w_{j,m}^* = 0$. First, for any m which satisfies $\sum_k \tilde{r}_{j,j,k,m} = 0$, we rewrite $\tilde{\mathcal{F}}_{\text{lb}, t}$ as

$$\begin{aligned}
&\max_{\{\mathbf{w}_l\}_{l=1}^L, t} t \\
&\text{s.t. } w_{j,m} \tilde{r}_{j,j,k,m} + \sum_{m' \neq m} w_{j,m'} \tilde{r}_{j,j,k,m'} \\
&\quad - \rho_{j,k} t \left(\sum_{l \neq j} \mathbf{w}_l^T \tilde{\mathbf{r}}_{l,j,k} + 1 \right) \geq 0, \quad \forall k \in \mathcal{G}_j \\
&\mathbf{w}_i^T \tilde{\mathbf{r}}_{i,i,k} - \rho_{i,k} t \left(w_{j,m} \tilde{r}_{j,i,k,m} + \sum_{m' \neq m} w_{j,m'} \tilde{r}_{j,i,k,m'} \right. \\
&\quad \left. + \sum_{l \neq i, l \neq j} \mathbf{w}_l^T \tilde{\mathbf{r}}_{l,i,k} + 1 \right) \geq 0, \quad \forall i \in \mathcal{G}_i, \quad \forall i \neq j \\
&\mathbf{1}^T \mathbf{w}_l \leq P_l, \quad \forall l \in \mathcal{L} \\
&\mathbf{w}_l \succeq \mathbf{0}, \quad \forall l \in \mathcal{L}. \quad (69)
\end{aligned}$$

²Since we consider the power constraints $\mathbf{1}^T \mathbf{w}_j \leq P_j$, there might exist different \mathbf{w}_j and \mathbf{w}_j' achieving the maximal value of the objective function in problem $\tilde{\mathcal{F}}_{\text{lb}, t}$. Here, the optimal solution refers to the vectors which achieve the maximum objective function value while the total transmit power $\sum_j \mathbf{1}^T \mathbf{w}_j$ is lowest.

Given the condition $\sum_k \tilde{r}_{j,j,k,m} = 0$, we can prove that $\tilde{r}_{j,j,k,m} = 0$ for $\forall k \in \mathcal{G}_j$, which results in $w_{j,m} \tilde{r}_{j,j,k,m} = 0$ for $\forall k \in \mathcal{G}_j$. Thus, observing problem (69), it is easy to find that $w_{j,m}^* = 0$.

Next, we consider another case where (m, m') satisfies the following condition

$$\begin{cases} \tilde{r}_{j,j,k,m} < \tilde{r}_{j,j,k,m'}, \forall k \in \mathcal{G}_j \\ \tilde{r}_{j,i,k,m} > \tilde{r}_{j,i,k,m'}, \forall k \in \mathcal{G}_i, \forall i \neq j. \end{cases} \quad (70)$$

Note that, for $\forall j \in \mathcal{L}, \forall k \in \mathcal{G}_j$, the KKT conditions for problem $\tilde{\mathcal{F}}_{lb,t}$ can be expressed as

$$1 - \sum_{j,k} \xi_{j,k} \rho_{j,k} \left(\sum_{i \neq j} \mathbf{w}_i^T \tilde{\mathbf{r}}_{i,j,k} + 1 \right) = 0 \quad (71a)$$

$$\sum_{k=1}^K \xi_{j,k} \tilde{\mathbf{r}}_{j,j,k} - \sum_{i \neq j} \sum_{k=1}^K \xi_{i,k} \rho_{i,k} t \tilde{\mathbf{r}}_{j,i,k} - \lambda_{lb,j} \mathbf{1} + \boldsymbol{\eta}_j = \mathbf{0} \quad (71b)$$

$$\xi_{j,k} \left(\rho_{j,k} t \left(\sum_{i \neq j} \mathbf{w}_i^T \tilde{\mathbf{r}}_{i,j,k} + 1 \right) - \mathbf{w}_j^T \tilde{\mathbf{r}}_{j,j,k} \right) = 0 \quad (71c)$$

$$\rho_{j,k} t \left(\sum_{i \neq j} \mathbf{w}_i^T \tilde{\mathbf{r}}_{i,j,k} + 1 \right) - \mathbf{w}_j^T \tilde{\mathbf{r}}_{j,j,k} \leq 0 \quad (71d)$$

$$\lambda_{lb,j} (\mathbf{1}^T \mathbf{w}_j - P_j) = 0, \mathbf{1}^T \mathbf{w}_j \leq P_j, \mathbf{w}_j \succeq \mathbf{0} \quad (71e)$$

$$\boldsymbol{\eta}_j \odot \mathbf{w}_j = \mathbf{0} \quad (71f)$$

where $\xi_{j,k} \geq 0$, $\lambda_{lb,j} \geq 0$ and $\boldsymbol{\eta}_j \succeq \mathbf{0}$ are Lagrange multipliers associated with the inequality constraints in problem $\tilde{\mathcal{F}}_{lb,t}$. Then, from the equation (71b), we have

$$\sum_{k=1}^K \xi_{j,k} \tilde{\mathbf{r}}_{j,j,k,m} - \sum_{i \neq j} \sum_{k=1}^K \xi_{i,k} \rho_{i,k} t \tilde{\mathbf{r}}_{j,i,k,m} + \eta_{j,m} = \lambda_{lb,j} \quad (72)$$

$$\sum_{k=1}^K \xi_{j,k} \tilde{\mathbf{r}}_{j,j,k,m'} - \sum_{i \neq j} \sum_{k=1}^K \xi_{i,k} \rho_{i,k} t \tilde{\mathbf{r}}_{j,i,k,m'} + \eta_{j,m'} = \lambda_{lb,j}. \quad (73)$$

Here, $\eta_{j,m}$ and $\eta_{j,m'}$ are the m th and m' th entries of $\boldsymbol{\eta}_j$, respectively. Given the pair (m, m') satisfying (70), subtracting (72) from (73) and using the condition $\eta_{j,m'} \geq 0$, we can obtain the following relations

$$\begin{aligned} \eta_{j,m} &\geq \eta_{j,m} - \eta_{j,m'} \\ &= \sum_{k=1}^K \xi_{j,k} (\tilde{r}_{j,j,k,m'} - \tilde{r}_{j,j,k,m}) \\ &\quad - \sum_{i \neq j} \sum_{k=1}^K \xi_{i,k} \rho_{i,k} t (\tilde{r}_{j,i,k,m'} - \tilde{r}_{j,i,k,m}) \\ &> 0. \end{aligned} \quad (74)$$

Note that the condition in (71a) indicates that there is at least one $\xi_{j,k} > 0$. Thus, the last inequality in (74) strictly holds. Then, owing to the KKT condition in (71f), we have $\eta_{j,m} w_{j,m}^* = 0$, and consequently, $w_{j,m}^* = 0$.

Define the set \mathcal{T}_j as

$$\mathcal{T}_j = \left\{ m \mid \sum_k \tilde{r}_{j,j,k,m} \neq 0, \nexists m' : (m, m') \text{ satisfies (70)} \right\}. \quad (75)$$

Then, based on the above discussions, the optimal precoding vectors \mathbf{p}_j^* can be expressed as

$$\mathbf{p}_j^* = \sum_{m \in \mathcal{T}_j} \sqrt{w_{j,m}^*} \mathbf{v}_m. \quad (76)$$

REFERENCES

- [1] "Cisco visual networking index: Global mobile data traffic forecast update, 2016-2021 white paper," Cisco, White Paper, Mar. 2017. [Online]. Available: <https://www.cisco.com/c/en/us/solutions/collateral/service-provider/visual-networking-index-vni/mobile-white-paper-c11-520862.html>
- [2] D. Lecomte and F. Gabin, "Evolved multimedia broadcast/multicast service (eMBMS) in LTE-advanced: Overview and Rel-11 enhancements," *IEEE Commun. Mag.*, vol. 50, no. 11, pp. 68–74, Nov. 2012.
- [3] N. D. Sidiropoulos, T. N. Davidson, and Z.-Q. Luo, "Transmit beamforming for physical-layer multicasting," *IEEE Trans. Signal Process.*, vol. 54, no. 6, pp. 2239–2251, Jun. 2006.
- [4] E. Karipidis, N. D. Sidiropoulos, and Z.-Q. Luo, "Quality of service and max-min fair transmit beamforming to multiple cochannel multicast groups," *IEEE Trans. Signal Process.*, vol. 56, no. 3, pp. 1268–1279, Mar. 2008.
- [5] G. Hsu, B. Liu, H. Wang, and H. Su, "Joint beamforming for multicell multigroup multicast with per-cell power constraints," *IEEE Trans. Veh. Technol.*, vol. 66, no. 5, pp. 4044–4058, May 2017.
- [6] X. Zhu, C. Jiang, L. Yin, L. Kuang, N. Ge, and J. Lu, "Cooperative multigroup multicast transmission in integrated terrestrial-satellite networks," *IEEE J. Sel. Areas Commun.*, vol. 36, no. 5, pp. 981–992, May 2018.
- [7] T. L. Marzetta, "Noncooperative cellular wireless with unlimited numbers of base station antennas," *IEEE Trans. Wireless Commun.*, vol. 9, no. 5, pp. 3590–3600, Nov. 2010.
- [8] E. G. Larsson, O. Edfors, F. Tufvesson, and T. L. Marzetta, "Massive MIMO for next generation wireless systems," *IEEE Commun. Mag.*, vol. 52, no. 2, pp. 186–195, Feb. 2014.
- [9] L. Lu, G. Y. Li, A. L. Swindlehurst, A. Ashikhmin, and R. Zhang, "An overview of massive MIMO: Benefits and challenges," *IEEE J. Sel. Topics Signal Process.*, vol. 8, no. 5, pp. 742–758, Oct. 2014.
- [10] E. Björnson, E. G. Larsson, and T. L. Marzetta, "Massive MIMO: Ten myths and one critical question," *IEEE Commun. Mag.*, vol. 54, no. 2, pp. 114–123, Feb. 2016.
- [11] C. Sun, X. Q. Gao, S. Jin, M. Matthaiou, Z. Ding, and C. Xiao, "Beam division multiple access transmission for massive MIMO communications," *IEEE Trans. Commun.*, vol. 63, no. 6, pp. 2170–2184, Jun. 2015.
- [12] Y. Wu, R. Schober, D. W. K. Ng, C. Xiao, and G. Caire, "Secure massive MIMO transmission with an active eavesdropper," *IEEE Trans. Inf. Theory*, vol. 62, no. 7, pp. 3880–3900, Jul. 2016.
- [13] L. You, X. Q. Gao, G. Y. Li, X. G. Xia, and N. Ma, "BDMA for millimeter-wave/terahertz massive MIMO transmission with per-beam synchronization," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 7, pp. 1550–1563, Jul. 2017.
- [14] B. Wang, F. Gao, S. Jin, H. Lin, and G. Y. Li, "Spatial- and frequency-wideband effects in millimeter-wave massive MIMO systems," *IEEE Trans. Signal Process.*, vol. 66, no. 13, pp. 3393–3406, Jul. 2018.
- [15] Z. Xiang, M. Tao, and X. Wang, "Massive MIMO multicasting in noncooperative cellular networks," *IEEE J. Sel. Areas Commun.*, vol. 32, no. 6, pp. 1180–1193, Jun. 2014.
- [16] E. Chen and M. Tao, "ADMM-based fast algorithm for multi-group multicast beamforming in large-scale wireless systems," *IEEE Trans. Commun.*, vol. 65, no. 6, pp. 2685–2698, Jun. 2017.
- [17] M. Sadeghi, L. Sanguinetti, R. Couillet, and C. Yuen, "Reducing the computational complexity of multicasting in large-scale antenna systems," *IEEE Trans. Wireless Commun.*, vol. 16, no. 5, pp. 2963–2975, May 2017.
- [18] M. Sadeghi, E. Björnson, E. G. Larsson, C. Yuen, and T. L. Marzetta, "Max-min fair transmit precoding for multi-group multicasting in massive MIMO," *IEEE Trans. Wireless Commun.*, vol. 17, no. 2, pp. 1358–1373, Feb. 2018.
- [19] A. L. Yuille and A. Rangarajan, "The concave-convex procedure," *Neural Computation*, vol. 15, no. 4, pp. 915–936, Apr. 2003.

- [20] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Found. Trends Mach. Learn.*, vol. 3, no. 1, pp. 1–122, Jan. 2011.
- [21] A. Lozano, A. M. Tulino, and S. Verdú, "Optimum power allocation for parallel Gaussian channels with arbitrary input distributions," *IEEE Trans. Inf. Theory*, vol. 52, no. 7, pp. 3033–3051, Jul. 2006.
- [22] C. Xiao, Y. R. Zheng, and Z. Ding, "Globally optimal linear precoders for finite alphabet signals over complex vector Gaussian channels," *IEEE Trans. Signal Process.*, vol. 59, no. 7, pp. 3301–3314, Jul. 2011.
- [23] Y. Wu, M. Wang, C. Xiao, Z. Ding, and X. Q. Gao, "Linear precoding for MIMO broadcast channels with finite-alphabet constraints," *IEEE Trans. Wireless Commun.*, vol. 11, no. 8, pp. 2906–2920, Aug. 2012.
- [24] W. Wu, K. Wang, W. Zeng, Z. Ding, and C. Xiao, "Cooperative multicell MIMO downlink precoding with finite-alphabet inputs," *IEEE Trans. Commun.*, vol. 63, no. 3, pp. 766–779, Mar. 2015.
- [25] S. R. Aghdam and T. M. Duman, "Joint precoder and artificial noise design for MIMO wiretap channels with finite-alphabet inputs based on the cut-off rate," *IEEE Trans. Wireless Commun.*, vol. 16, no. 6, pp. 3913–3923, Jun. 2017.
- [26] J. Jin, Y. R. Zheng, W. Chen, and C. Xiao, "Hybrid precoding for millimeter wave MIMO systems: A matrix factorization approach," *IEEE Trans. Wireless Commun.*, vol. 17, no. 5, pp. 3327–3339, May 2018.
- [27] W. Zeng, C. Xiao, M. Wang, and J. Lu, "Linear precoding for finite-alphabet inputs over MIMO fading channels with statistical CSI," *IEEE Trans. Signal Process.*, vol. 60, no. 6, pp. 3134–3148, Jun. 2012.
- [28] Y. Wu, C. Wen, C. Xiao, X. Q. Gao, and R. Schober, "Linear precoding for the MIMO multiple access channel with finite alphabet inputs and statistical CSI," *IEEE Trans. Wireless Commun.*, vol. 14, no. 2, pp. 983–997, Feb. 2015.
- [29] Y. Wu, D. W. K. Ng, C. Wen, R. Schober, and A. Lozano, "Low-complexity MIMO precoding for finite-alphabet signals," *IEEE Trans. Wireless Commun.*, vol. 16, no. 7, pp. 4571–4584, Jul. 2017.
- [30] Z. Xiang, M. Tao, and X. Wang, "Coordinated multicast beamforming in multicell networks," *IEEE Trans. Wireless Commun.*, vol. 12, no. 1, pp. 12–21, Jan. 2013.
- [31] S. Wagner, R. Couillet, M. Debbah, and D. T. M. Slock, "Large system analysis of linear precoding in correlated MISO broadcast channels under limited feedback," *IEEE Trans. Inf. Theory*, vol. 58, no. 7, pp. 4509–4537, Jul. 2012.
- [32] L. You, X. Q. Gao, X.-G. Xia, N. Ma, and Y. Peng, "Pilot reuse for massive MIMO transmission over spatially correlated Rayleigh fading channels," *IEEE Trans. Wireless Commun.*, vol. 14, no. 6, pp. 3352–3366, Jun. 2015.
- [33] A. Adhikary, J. Nam, J.-Y. Ahn, and G. Caire, "Joint spatial division and multiplexing: The large-scale array regime," *IEEE Trans. Inf. Theory*, vol. 59, no. 10, pp. 6441–6463, Oct. 2013.
- [34] K. Liu, V. Raghavan, and A. M. Sayeed, "Capacity scaling and spectral efficiency in wide-band correlated MIMO channels," *IEEE Trans. Inf. Theory*, vol. 49, no. 510, pp. 2504–2526, Oct. 2003.
- [35] W. Weichselberger, M. Herdin, H. Özcelik, and E. Bonek, "A stochastic MIMO channel model with joint correlation of both link ends," *IEEE Trans. Wireless Commun.*, vol. 5, no. 1, pp. 90–100, Jan. 2006.
- [36] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, U.K.: Cambridge Univ. Press, 2004.
- [37] D. H. Brandwood, "A complex gradient operator and its application in adaptive array theory," *IEE Proc., Parts F and H*, vol. 130, no. 1, pp. 11–16, Feb. 1983.
- [38] L. You, X. Q. Gao, A. L. Swindlehurst, and W. Zhong, "Channel acquisition for massive MIMO-OFDM with adjustable phase shift pilots," *IEEE Trans. Signal Process.*, vol. 64, no. 6, pp. 1461–1476, Mar. 2016.
- [39] G. R. Lanckriet and B. K. Sriperumbudur, "On the convergence of the concave-convex procedure," in *Proc. Adv. Neural Inf. Process. Syst.*, 2009, pp. 1759–1767.
- [40] C. Sun, X. Q. Gao, and Z. Ding, "BDMA in multicell massive MIMO communications: Power allocation algorithms," *IEEE Trans. Signal Process.*, vol. 65, no. 11, pp. 2962–2974, Jun. 2017.



Wenqian Wu (S'15) received the B.S. degree with electrical engineering from Southeast University, Nanjing, China, in 2015. He is currently working toward the Ph.D. degree with the National Mobile Communications Research Laboratory, Southeast University, Nanjing, China. From February 2018, he is a Visiting Student with the Department of Electrical and Computer Engineering, Lehigh University, Bethlehem, PA, USA. His research interests include massive MIMO communications and physical layer security.



Chengshan Xiao (M'99–SM'02–F'10) received the B.Sc. degree in electronic engineering from the University of Electronic Science and Technology of China, Chengdu, China, in 1987, the M.Sc. degree in electronic engineering from Tsinghua University, Beijing, China, in 1989, and the Ph.D. degree in electrical engineering from the University of Sydney, Australia, in 1997.

He is the Chandler Weaver Professor and Chair of the Department of Electrical and Computer Engineering, Lehigh University, Bethlehem, PA, USA. He is a fellow of the Canadian Academy of Engineering. Previously, he was a Program Director with the Division of Electrical, Communications and Cyber Systems, USA National Science Foundation. He was a Senior Member of scientific staff with Nortel Networks, Ottawa, Canada; a Faculty Member with Tsinghua University, Beijing, China; the University of Alberta, Edmonton, Canada; the University of Missouri, Columbia, MO; and Missouri University of Science and Technology, Rolla, MO. He also held Visiting Professor positions in Germany and Hong Kong. He is the holder of several patents granted in USA, Canada, China, and Europe. His invented algorithms have been implemented into Nortel's base station radio products after successful technical field trials and network integration. His research interests include wireless communications, signal processing, and underwater acoustic communications.

Dr. Xiao is the Awards Committee Chair and Elected Member-at-Large of Board of Governors of IEEE Communications Society. Previously, he served on the IEEE Technical Activity Board (TAB) Periodical Committee, he was an Elected Member-at-Large of Board of Governors, a member of Fellow Evaluation Committee, Director of Conference Publications, Distinguished Lecturer of the IEEE Communications Society, and Distinguished Lecturer of the IEEE Vehicular Technology, Society. He also served as an Editor, Area Editor, and the Editor-in-Chief of the IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS, an Associate Editor of the IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY, and IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS-I. He was the Technical Program Chair of the 2010 IEEE International Conference on Communications, Cape Town, South Africa, a Technical Program Co-Chair of the 2017 IEEE Global Communications Conference, Singapore. He served as the Founding Chair of the IEEE Wireless Communications Technical Committee. He received several distinguished awards including 2014 Humboldt Research Award, 2014 IEEE Communications Society Joseph LoCicero Award, 2015 IEEE Wireless Communications Technical Committee Recognition Award, and 2017 IEEE Communications Society Harold Sobol Award.



Xiqi Gao (S'92–AM'96–M'02–SM'07–F'15) 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 with the Massachusetts Institute of Technology, Cambridge, 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. From 2007 to 2012, he served as an Editor for the IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS. From 2009 to 2013, he served as an Associate Editor for the IEEE TRANSACTIONS ON SIGNAL PROCESSING. From 2015 to 2017, he served as an Editor for the IEEE TRANSACTIONS ON COMMUNICATIONS.

Dr. Gao was the recipient of 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 communications theory.