ADMM-based Fast Algorithm for Multi-group Multicast Beamforming in Large-Scale Wireless Systems

Erkai Chen and Meixia Tao

Abstract

Multi-group multicast beamforming in wireless systems with large antenna arrays and massive audience is investigated in this paper. Multicast beamforming design is a well-known non-convex quadratically constrained quadratic programming (QCQP) problem. A conventional method to tackle this problem is to approximate it as a semi-definite programming problem via semi-definite relaxation, whose performance, however, deteriorates considerably as the number of per-group users goes large. A recent attempt is to apply convex-concave procedure (CCP) to find a stationary solution by treating it as a difference of convex programming problem, whose complexity, however, increases dramatically as the problem size increases. In this paper, we propose a low-complexity high-performance algorithm for multi-group multicast beamforming design in large-scale wireless systems by leveraging the alternating direction method of multipliers (ADMM) together with CCP. In specific, the original non-convex QCQP problem is first approximated as a sequence of convex subproblems via CCP. Each convex subproblem is then reformulated as a novel ADMM form. Our ADMM reformulation enables that each updating step is performed by solving multiple small-size subproblems with closed-form solutions in parallel. Numerical results show that our fast algorithm maintains the same favorable performance as state-of-the-art algorithms but reduces the complexity by orders of magnitude.

Index Terms

Physical layer multicasting, large-scale optimization, non-convex quadratically constrained quadratic programming (QCQP), convex-concave procedure (CCP), alternating direction method of multipliers (ADMM).

I. Introduction

Multicasting is a promising approach to deliver a common message to multiple receivers by exploiting the broadcast nature of wireless medium. It has great potential in many applications such as live video streaming, venue casting, mobile application updates, advertisements, and public group communications [2]–[4]. It can also be used in heterogeneous networks (HetNets) with wireless backhaul to push common information from a macro base station (BS) to multiple small BSs [5]. Recently, multicasting is shown to be useful for content delivery even when user demands are different in wireless cache networks [6], [7].

Physical layer multicasting via beamforming further boosts the energy and spectrum efficiencies by exploiting channel state information at the transmitter. Multicast beamforming is first considered in [8] for a single group of users. The similar problem for multiple co-channel groups is later studied in [9]. In [10], coordinated multicast beamforming among BSs in multi-cell networks is investigated. The design of multicast beamforming in cellular networks with massive multiple-input multiple-output (MIMO) is studied in [11]. In general, the design of multicast beamforming is a non-convex quadratically constrained quadratic programming (QCQP) problem and its global optimal solution is difficult to obtain. A prevailing method to tackle this problem is to approximate it as a convex semi-definite programming (SDP) problem via semi-definite relaxation (SDR) [8]–[10], [12]. By solving the relaxed problem, one can obtain not only a bound on the optimal performance, but also possibly a global optimal solution

This work is supported by the National Natural Science Foundation of China under grants 61571299 and 61521062. This paper will be presented in part at the IEEE ICC 2017 [1].

E. Chen and M. Tao are with the Department of Electronic Engineering at Shanghai Jiao Tong University, Shanghai, P. R. China (emails: {cek1006, mxtao}@sjtu.edu.cn).

of the original problem if the SDP solution is rank-one. If the SDP solution is not rank-one, then it is followed by a post-processing randomization procedure to generate high-quality approximate solutions. However, the SDR-based algorithm is not computationally efficient as it requires lifting the problem into higher dimensional space. Meanwhile, it has been demonstrated in [9] that, the performance of SDR with Gaussian randomization deteriorates considerably as the number of per-group users increases. A more recent method to tackle the multicast beamforming problem is to apply successive convex approximation (SCA) [13] or convex-concave procedure (CCP) [14] to find a stationary solution [15]–[18]. In particular, the CCP-based algorithms treat the problem as a difference of convex (DC) programming problem and then iteratively solve a sequence of convex subproblems constructed by replacing the concave parts of the DC functions with their first-order Taylor expansions. Both the SCA and CCP algorithms have demonstrated better performance than the SDR methods when the rank-one probability of SDP solutions is low [17], [18]. Note that the SCA and CCP methods have also been applied recently to tackle non-convex resource allocation and beamforming design problems in various wireless networks, e.g., [19], [20].

Note that none of the aforementioned SDR or CCP algorithms for multicast beamforming design have found closed-form solutions for the approximated problems and they all rely on general-purpose solvers to find the solutions numerically. In the simplest scenario with single-group multicast beamforming, if a convex solver based on interior-point methods (e.g., CVX [21]) is used, then the worst-case complexity is $\mathcal{O}(N^2+K)^{3.5}$ for SDR [8] and is $\mathcal{O}(N+K)^{3.5}$ for each CCP iteration [15], where N is the number of antennas at the transmitter and K is the number of users. Clearly, the computational costs of these second-order algorithms are not scalable when the problem size increases. Recently, ultra-dense small cell deployments [22] and massive MIMO [23] have become important candidate technologies for the future generations of wireless systems (5G). They can provide high-volume and diversified data services for a large set of devices, including not only smart mobile devices operated by humans but also special-purpose machine-to-machine devices [24]. The scale of these wireless systems is significantly large, dealing possibly with hundreds of antennas and users. It is thus of particular importance to investigate low-complexity and high-performance algorithms for multicast beamforming in large-scale systems.

A. Related Works

The authors in [25] propose a fast algorithm to solve the max-min fairness (MMF) problem for single-group multicast beamforming. The original problem is first approximated by replacing the MMF objective with its proportional fairness. The additive update and multiplicative update algorithms are then introduced by iteratively maximizing two different concave approximations of the new objective, both of which are updated in closed form. The algorithms are demonstrated to achieve comparable performance as CCP-based algorithms at a much lower complexity. This method, however, cannot be extended to the more general multi-group multicasting scenario.

The alternating direction method of multipliers (ADMM) [26] is a powerful first-order method well suited to large-scale convex optimization. It has recently been used in wireless networks for distributed or parallel computing of various resource allocation problems and beamforming design problems. The authors in [27] present a two-stage optimization framework to efficiently solve large-scale power minimization and network utility maximization problems for dense wireless cooperative networks which can be exactly reformulated as second-order cone programming (SOCP) problems. In the two-stage approach, the original SOCP problem is first transformed into a standard cone programming form with matrix stuffing and then solved using the ADMM algorithm (i.e., the operator splitting method) proposed in [28]. In [29], the authors propose a consensus form of ADMM to solve a general QCQP problem. It first reformulates the QCQP problem in consensus optimization form by introducing a local copy of the optimization variables (i.e., global variable) for each single quadratic constraint and then applies ADMM to update the local and global variables alternatively. In particular, the ADMM updating step for each local copy of variables is a QCQP with only one constraint (QCQP-1) and hence can be done efficiently with possibly closed form. The authors in [29] then apply the consensus ADMM for single-group multicast beamforming problem directly.

Its convergence, however, cannot be guaranteed due to the non-convexity of the multicast beamforming problem.

Both the two-stage approach for SOCPs [27] and consensus ADMM for QCQPs [29] deal with general optimization frameworks, but require a large number of auxiliary variables and hence lift the original problem into a much higher dimension space. Customized algorithms for specific problems to enable efficient or parallel computing are developed in [30], [31] by exploiting the structures of the problems. Specifically, in [30], the authors consider the robust coordinated beamforming in multi-cell networks. By using SDR approximation and S-Procedure methods, the original non-convex beamforming problem is first reformulated as a tractable convex SDP. Then an ADMM-based distributed algorithm is proposed, which is provably able to converge to the global optimum of the centralized SDP problem. The ADMM update in each BS, however, still relies on general-purpose solvers. In [31], the joint BS activation and beamforming design for power minimization problem in HetNets is first reformulated as an SOCP using a sparsity regularizer. An efficient algorithm based on ADMM is developed to solve the SOCP, in which each updating step is in closed form and can be carried out distributively among multiple BSs.

B. Contributions

The main contribution of this work is the development of a low-complexity high-performance algorithm for multi-group multicast beamforming design by adopting the ADMM approach in conjunction with the CCP method. Both the quality-of-service (QoS) problem and MMF problems are considered, where the QoS problem is to minimize the total transmit power subject to an individual signal-to-interference-and-noise ratio (SINR) constraint for each user and a peak power constraint for each antenna, and the MMF problem is to maximize the minimum weighted SINR subject to the per-antenna peak power constraints. The proposed algorithm is first designed for the QoS problem and then extended to the MMF problem. The main technical novelty and research findings of this work are summarized as follows:

- The proposed algorithm exploits the advantages of both the CCP principle and the ADMM approach. Specifically, the QoS problem in the non-convex QCQP form is first approximated with a sequence of convex subproblems by adopting CCP, which enables superior performance over the conventional SDR methods. Each convex CCP subproblem is then reformulated as a novel ADMM form that facilitates close-form solutions and parallel computing in large-scale systems.
- In the ADMM reformulation of each convex CCP subproblem, each updating step is decomposed into multiple small-size subproblems which are solved in parallel with closed-form solutions. Compared with the two-stage approach [27] and the consensus ADMM method [29] proposed for general frameworks, our new ADMM utilizes the specific structure of the considered multi-group multicast beamforming problem and hence requires much less auxiliary variables and reduces the complexity significantly.
- We further propose an efficient ADMM-based method to obtain a starting point for CCP. In the special case when the number of transmit antennas exceeds the number of users, we also find a closed-form starting point.
- Numerical simulations are conducted in large-scale systems with transmit antenna number $N \in [40, 250]$ and user number $K \in [50, 140]$. The results show that the proposed fast algorithm maintains the same favorable performance as the existing CCP algorithm (solved by interior-point solvers), which is within 1dB close to the SDR lower bound for the QoS problem or within 0.5dB close to the SDR upper bound for the MMF problem, but reduces the running time by orders of magnitude.

C. Organization and Notations

The rest of the paper is organized as follows. Section II introduces the system model and problem formulations, including the QoS problem and MMF problem. Section III provides the details of the proposed fast algorithm for the QoS problem. The extension of the algorithm to the MMF problem is

introduced in Section IV. Simulation results are provided in Section V. Finally, we conclude the paper in Section VI.

Notations: In the remainder of this paper, boldface lower-case and upper-case letters denote vectors and matrices respectively. Calligraphy letters denote sets or problems, depending on the context. \mathbb{R} and \mathbb{C} denote the real and complex domains, respectively. $\mathbb{E}(\cdot)$ denotes the expectation of a random variable. $\mathcal{CN}(\delta, \sigma^2)$ represents a complex Gaussian distribution with mean δ and variance σ^2 . The operators $(\cdot)^T$, $(\cdot)^H$, $(\cdot)^{-1}$, $(\cdot)^{\dagger}$, and $\mathrm{Tr}(\cdot)$ correspond to the transpose, the Hermitian transpose, inverse, Moore-Penrose inverse, and trace respectively. $|\cdot|$ and $||\cdot||_2$ denote the absolute value and Euclidean norm, respectively. The real part of a complex number x is denoted by $\Re\{x\}$. Finally, \mathbf{I}_N denotes the $N \times N$ identity matrix.

II. PROBLEM SETTING

We consider a multi-group multicasting system as in [9] where one transmitter, equipped with N antennas, serves M groups of single-antenna users. The users within each group desire a common multicast message, which is independent for different groups. Let \mathcal{G}_m denote the set of users in multicast group m, for all $m \in \mathcal{M} = \{1, \ldots, M\}$. Let $\mathcal{K} = \{1, \ldots, K\}$ denote the set of total users in all groups. Each user participates in only one multicast group, thus $\mathcal{G}_i \cap \mathcal{G}_j = \emptyset, i \neq j, \forall i, j \in \mathcal{M}$ and $\bigcup_{m=1}^M \mathcal{G}_m = \mathcal{K}$. Let $\mathbf{w}_m \in \mathbb{C}^N$ denote the beamforming vector of group m, for all $m \in \mathcal{M}$ and $\mathbf{h}_k \in \mathbb{C}^N$ denote the

Let $\mathbf{w}_m \in \mathbb{C}^N$ denote the beamforming vector of group m, for all $m \in \mathcal{M}$ and $\mathbf{h}_k \in \mathbb{C}^N$ denote the channel vector from the transmitter to the k-th user, for all $k \in \mathcal{K}$. Each channel is modeled as a complex, random vector which is flat in frequency and quasi-static in time. And it is assumed to be perfectly available at the transmitter. The corresponding received signal at user k in group m can be written as

$$y_k = \underbrace{\mathbf{h}_k^H \mathbf{w}_m x_m}_{\text{desired signal}} + \underbrace{\sum_{j \neq m} \mathbf{h}_k^H \mathbf{w}_j x_j}_{\text{inter-group interference}} + \underbrace{n_k}_{\text{noise}}, \ \forall k \in \mathcal{G}_m,$$

$$(1)$$

where $x_m \in \mathbb{C}$ is the data symbol transmitted to multicast group m with $\mathbb{E}[|x_m|^2] = 1$, and $n_k \sim \mathcal{CN}(0, \sigma_k^2)$ is the additive white Gaussian noise at user k. The maximum power radiated by each antenna n is denoted as P_n . We thus have

$$\sum_{m=1}^{M} \mathbf{w}_{m}^{H} \mathbf{R}_{n} \mathbf{w}_{m} \leq P_{n}, \ \forall n \in \mathcal{N},$$
(2)

where $\mathcal{N} = \{1, \dots, N\}$ is the index set of all the antennas at the transmitter and $\mathbf{R}_n \in \{0, 1\}^{N \times N}$ is the all-zero matrix except the *n*-th diagonal entry being 1, for all $n \in \mathcal{N}$. The received SINR of user $k \in \mathcal{G}_m$ is expressed as

$$SINR_k = \frac{|\mathbf{h}_k^H \mathbf{w}_m|^2}{\sum_{j \neq m} |\mathbf{h}_k^H \mathbf{w}_j|^2 + \sigma_k^2}, \ \forall k \in \mathcal{G}_m.$$
(3)

Similar to [9], we consider two problem formulations for the multi-group multicast beamforming design, the QoS problem and MMF problem. The QoS problem is to minimize the total radiated power of the transmitter subject to an individual target SINR constraint for each user and a peak power constraint for each transmit antenna. It is expressed as

$$Q: \underset{\mathbf{w}}{\text{minimize}} \sum_{m=1}^{M} \|\mathbf{w}_m\|_2^2$$
 (4a)

subject to
$$\frac{|\mathbf{h}_k^H \mathbf{w}_m|^2}{\sum_{j \neq m} |\mathbf{h}_k^H \mathbf{w}_j|^2 + \sigma_k^2} \ge \gamma_k, \ \forall k \in \mathcal{G}_m, \forall m \in \mathcal{M},$$
 (4b)

$$\sum_{m=1}^{M} \mathbf{w}_{m}^{H} \mathbf{R}_{n} \mathbf{w}_{m} \leq P_{n}, \ \forall n \in \mathcal{N},$$

$$(4c)$$

where $\mathbf{w} \triangleq \{\mathbf{w}_m | m \in \mathcal{M}\}$ is the set of all the beamforming vectors and γ_k is the minimum received SINR required by user k.

The MMF problem is to maximize the minimum weighted SINR over all users subject to a peak power constraint for each transmit antenna. It is expressed as

$$\mathcal{F}$$
: maximize t (5a)

subject to
$$\frac{1}{g_k} \frac{|\mathbf{h}_k^H \mathbf{w}_m|^2}{\sum_{j \neq m} |\mathbf{h}_k^H \mathbf{w}_j|^2 + \sigma_k^2} \ge t, \ \forall k \in \mathcal{G}_m, \forall m \in \mathcal{M},$$
 (5b)

$$\sum_{m=1}^{M} \mathbf{w}_{m}^{H} \mathbf{R}_{n} \mathbf{w}_{m} \leq P_{n}, \ \forall n \in \mathcal{N},$$
(5c)

where $\{g_k\}_{k=1}^K$ is the set of weights to account for possibly different service grades among the users.

Slightly different from the original QoS problem and MMF problem formulated in [9], our problems have included the per-antenna peak power constraints, which are more realistic in practical systems. As mentioned in the introduction, both QoS and MMF problems can be approximately solved using the SDR method and the CCP method. In the next section, we focus on the QoS problem and propose a fast algorithm to find a stationary solution (possibly local optimum) based on the CCP method and the ADMM approach. The extension to the MMF problem shall be briefly discussed in Section IV.

Note that while the MMF problem \mathcal{F} is always feasible, the QoS problem \mathcal{Q} can be infeasible if the SINR targets $\{\gamma_k\}$ and the peak power constraints $\{P_n\}$ are too stringent, or the channels of users in different multicast groups are highly correlated. In [10], a necessary condition for the multi-cell multicast QoS beamforming problem to be feasible is given. In this paper, we derive a sufficient condition for checking the feasibility of \mathcal{Q} , which shall be introduced in Section III-C. In the following, we only discuss \mathcal{Q} when it is feasible.

III. FAST ALGORITHM FOR QOS PROBLEM

Recently, the authors in [18] formulate a novel sparse multicast beamforming problem in cache-enabled cloud radio access networks which is solved using CCP-based algorithms. Since the multi-group multicast beamforming problem considered in this work is a special case of the problem in [18] ¹, the same CCP-based algorithm in [18] applies here. The new contribution in this work is to develop its low-complexity implementation using ADMM. We name the proposed fast algorithm as CCP-ADMM algorithm. In this section, we first briefly outline the CCP-based algorithm, then present a novel ADMM approach to solve each CCP subproblem with closed-form solutions for the updating steps. After that, we introduce a new method to find a starting point for the CCP algorithm also based on ADMM.

A. CCP Algorithm

As in [18], the SINR constraints (4b) in problem Q can be written as a DC form

$$\gamma_k \left(\sum_{j \neq m} |\mathbf{h}_k^H \mathbf{w}_j|^2 + \sigma_k^2 \right) - |\mathbf{h}_k^H \mathbf{w}_m|^2 \le 0, \ \forall k \in \mathcal{G}_m, \forall m \in \mathcal{M}.$$
 (6)

The CCP-based algorithm is to convexify the SINR constraints in the above DC form by replacing the concave parts with their first-order Taylor expansions, then solve a sequence of convex subproblems successively. As such, CCP designed for DC programming problems is also known as a special case of SCA, which is for general non-convex problems. This iterative procedure is guaranteed to converge to a

¹The optimization problem in [18] is to minimize the weighted sum of total transmission power and backhaul cost subject to a minimum SINR constraint for each multicast group. It reduces to the QoS problem considered in this work when the backhaul cost is ignored.

stationary point of the original problem Q, according to [14]. Specifically, in the t-th iteration, we need to solve

$$Q^{(t)}: \underset{\mathbf{w}}{\text{minimize}} \sum_{m=1}^{M} \|\mathbf{w}_m\|_2^2$$
 (7a)

subject to
$$\gamma_k \left(\sum_{j \neq m} |\mathbf{h}_k^H \mathbf{w}_j|^2 + \sigma_k^2 \right) - 2 \mathfrak{Re} \left\{ (\mathbf{w}_m^{(t)})^H \mathbf{h}_k \mathbf{h}_k^H \mathbf{w}_m \right\}$$

$$+ |\mathbf{h}_k^H \mathbf{w}_m^{(t)}|^2 \le 0, \ \forall k \in \mathcal{G}_m, \forall m \in \mathcal{M},$$
 (7b)

$$\sum_{m=1}^{M} \mathbf{w}_{m}^{H} \mathbf{R}_{n} \mathbf{w}_{m} \leq P_{n}. \ \forall n \in \mathcal{N},$$
 (7c)

where $\mathbf{w}_m^{(t)}$ is the optimal solution obtained from the previous iteration. Each $\mathcal{Q}^{(t)}$ is a convex QCQP problem and can be solved using a general-purpose solver through interior-point methods. In what follows, we exploit the specific structure of problem $\mathcal{Q}^{(t)}$ and find its optimal solution using an ADMM-based fast algorithm.

Note that a parallel SCA algorithm is proposed in [20] to solve a resource allocation problem in a multi-cell cellular system with coordinated multi-point transmission (CoMP). In each iteration of the parallel SCA, all BSs are enabled to update their variables through parallelly solving multiple convexified and decoupled subproblems, one for each BS. Different from [20], where the optimization variables from different BSs are separated in the constraints, the optimization variables $\{\mathbf{w}_m\}$ of our problem are coupled in the constraints (i.e., (4b) and (4c)). We could use the similar idea to approximate the original problem \mathcal{Q} by a sequence of parallel subproblems with each \mathbf{w}_m decoupled in both the objective and the constraints. However, the solutions of these parallel subproblems in each SCA iteration may not be feasible to the original problem \mathcal{Q} as the originally coupled constraints may no longer be satisfied.

B. ADMM algorithm for each CCP subproblem

We first introduce two sets of auxiliary variables (their significance will be discussed in the end of this subsection):

$$\Gamma_{k,m} = \mathbf{h}_k^H \mathbf{w}_m, \ \forall k \in \mathcal{K}, m \in \mathcal{M},$$
 (8)

$$\mathbf{v}_m = \mathbf{w}_m, \ \forall m \in \mathcal{M},\tag{9}$$

where $\Gamma_{k,m}$ represents the interference level at user k caused by the beamforming vector of group m, and $\mathbf{v}_m \in \mathbb{C}^N$ is a copy of the original beamforming vector \mathbf{w}_m . To ease the notation, we define $\Gamma \triangleq \{\Gamma_{k,m} | k \in \mathcal{K}, m \in \mathcal{M}\}$, and $\mathbf{v} \triangleq \{\mathbf{v}_m | m \in \mathcal{M}\}$.

Then problem $\mathcal{Q}^{(t)}$ can be equivalently expressed as

$$\underset{\{\mathbf{\Gamma}, \mathbf{v}, \mathbf{w}\}}{\text{minimize}} \quad \sum_{m=1}^{M} \|\mathbf{w}_m\|_2^2 \tag{10a}$$

subject to
$$\Gamma_{k,m} - \mathbf{h}_k^H \mathbf{w}_m = 0, \ \forall k \in \mathcal{K}, m \in \mathcal{M},$$
 (10b)

$$\mathbf{v}_m - \mathbf{w}_m = 0, \ \forall m \in \mathcal{M}, \tag{10c}$$

$$\gamma_k \left(\sum_{i \neq m} |\Gamma_{k,j}|^2 + \sigma_k^2 \right) - 2 \mathfrak{Re} \left\{ (\mathbf{w}_m^{(t)})^H \mathbf{h}_k \Gamma_{k,m} \right\}$$

$$+ |\mathbf{h}_k^H \mathbf{w}_m^{(t)}|^2 \le 0, \ \forall k \in \mathcal{G}_m, \forall m \in \mathcal{M},$$
(10d)

$$\sum_{m=1}^{M} \mathbf{v}_{m}^{H} \mathbf{R}_{n} \mathbf{v}_{m} \leq P_{n}, \ \forall n \in \mathcal{N}.$$
 (10e)

We define the feasible region of constraint (10d) as C, and its indicator function as

$$\mathbb{I}_{\mathcal{C}}(\Gamma) = \begin{cases} 0, & \text{if } \Gamma \in \mathcal{C}, \\ +\infty, & \text{otherwise.} \end{cases}$$
(11)

Similarly, we define the feasible region of constraint (10e) as \mathcal{D} , and its indicator function as

$$\mathbb{I}_{\mathcal{D}}(\mathbf{v}) = \begin{cases} 0, & \text{if } \mathbf{v} \in \mathcal{D}, \\ +\infty, & \text{otherwise.} \end{cases}$$
(12)

Then, we obtain the equivalent ADMM reformulation of problem $Q^{(t)}$

minimize
$$\sum_{\{\Gamma, \mathbf{v}, \mathbf{w}\}}^{M} \|\mathbf{w}_{m}\|_{2}^{2} + \mathbb{I}_{\mathcal{C}}(\Gamma) + \mathbb{I}_{\mathcal{D}}(\mathbf{v})$$
subject to (10b), (10c).

The augmented Lagrangian (using the scaled dual variable) of problem (13) is given by

$$\mathcal{L}_{\rho}(\mathbf{\Gamma}, \mathbf{v}, \mathbf{w}, \boldsymbol{\lambda}, \mathbf{z}) = \sum_{m=1}^{M} \|\mathbf{w}_{m}\|_{2}^{2} + \mathbb{I}_{\mathcal{C}}(\mathbf{\Gamma}) + \mathbb{I}_{\mathcal{D}}(\mathbf{v}) + \frac{\rho}{2} \sum_{k=1}^{K} \sum_{m=1}^{M} |\Gamma_{k,m} - \mathbf{h}_{k}^{H} \mathbf{w}_{m} + \lambda_{k,m}|^{2}$$

$$+ \frac{\rho}{2} \sum_{m=1}^{M} \|\mathbf{v}_{m} - \mathbf{w}_{m} + \mathbf{z}_{m}\|_{2}^{2},$$
(14)

where $\rho > 0$ is the penalty parameter, $\lambda \triangleq \{\lambda_{k,m} | k \in \mathcal{K}, m \in \mathcal{M}\}$ and $\mathbf{z} \triangleq \{\mathbf{z}_m \in \mathbb{C}^N | m \in \mathcal{M}\}$ are the scaled dual variables for constraints (10b) and (10c), respectively.

From problem (13), we observe that the variables in the constraints can be split into two blocks, $\{\Gamma, \mathbf{v}\}$ and \mathbf{w} , and that the objective function is also separable across this splitting. Thus, by adopting ADMM, we can minimize $\mathcal{L}_{\rho}(\Gamma, \mathbf{v}, \mathbf{w}, \lambda, \mathbf{z})$ by updating the two blocks of variables, $\{\Gamma, \mathbf{v}\}$ and \mathbf{w} , alternatively. The ADMM procedure is given in Alg. 1.

Algorithm 1 ADMM for solving problem $Q^{(t)}$

Initialization: Initialize $\mathbf{w}_m^0 \leftarrow \mathbf{w}_m^{(t)}, \mathbf{z}_m^0 \leftarrow \mathbf{0}, \lambda_{k,m}^0 \leftarrow 0, \forall m \in \mathcal{M}, k \in \mathcal{K}, \text{ and } j \leftarrow 0.$ Set the penalty parameter ρ .

Repeat

1) Update the first block of variables $\{\Gamma^{j+1}, \mathbf{v}^{j+1}\}$

$$\{\mathbf{\Gamma}^{j+1}, \mathbf{v}^{j+1}\} := \arg\min_{\mathbf{\Gamma}, \mathbf{v}} \mathcal{L}_{\rho}(\mathbf{\Gamma}, \mathbf{v}, \mathbf{w}^{j}, \boldsymbol{\lambda}^{j}, \mathbf{z}^{j}). \tag{15}$$

2) Update the second block of variables \mathbf{w}^{j+1}

$$\mathbf{w}^{j+1} := \arg\min_{\mathbf{w}} \mathcal{L}_{\rho}(\Gamma^{j+1}, \mathbf{v}^{j+1}, \mathbf{w}, \boldsymbol{\lambda}^{j}, \mathbf{z}^{j}). \tag{16}$$

3) Update the dual variables $\{ \boldsymbol{\lambda}^{j+1}, \mathbf{z}^{j+1} \}$

$$\lambda_{k,m}^{j+1} := \lambda_{k,m}^j + (\Gamma_{k,m}^{j+1} - \mathbf{h}_k^H \mathbf{w}_m^{j+1}), \ \forall k \in \mathcal{K}, m \in \mathcal{M},$$

$$(17)$$

$$\mathbf{z}_m^{j+1} := \mathbf{z}_m^j + (\mathbf{v}_m^{j+1} - \mathbf{w}_m^{j+1}), \ \forall m \in \mathcal{M}.$$

4) Set $j \leftarrow j + 1$.

Until convergence criterion is met.

In the following, we elaborate the details of updating the primal variables $\{\Gamma, \mathbf{v}\}$ and \mathbf{w} (the superscript ADMM iteration counter is ignored for simplicity). Note that the update for the first block of variables in (15) can be decomposed into two independent problems, one for each of the two sets Γ and \mathbf{v} , which are given in (19) and (20), respectively.

$$\mathbf{\Gamma}^{j+1} := \arg\min_{\mathbf{\Gamma}} \left\{ \mathbb{I}_{\mathcal{C}}(\mathbf{\Gamma}) + \frac{\rho}{2} \sum_{k=1}^{K} \sum_{m=1}^{M} |\Gamma_{k,m} - \mathbf{h}_{k}^{H} \mathbf{w}_{m}^{j} + \lambda_{k,m}^{j}|^{2} \right\}$$
(19)

$$\mathbf{v}^{j+1} := \arg\min_{\mathbf{v}} \left\{ \mathbb{I}_{\mathcal{D}}(\mathbf{v}) + \frac{\rho}{2} \sum_{m=1}^{M} \|\mathbf{v}_m - \mathbf{w}_m^j + \mathbf{z}_m^j\|_2^2 \right\}$$
(20)

1) Γ Update: The update of Γ in problem (19) is equivalent to solving the problem:

minimize
$$\sum_{k=1}^{K} \sum_{m=1}^{M} |\Gamma_{k,m} - \mathbf{h}_k^H \mathbf{w}_m^j + \lambda_{k,m}^j|^2$$
 subject to (10d). (21)

It is observed that problem (21) can be decomposed into K subproblems, one for each $k \in K$:

minimize
$$\sum_{\{\Gamma_{k,m}\}_{m=1}^{M}}^{M} \sum_{m=1}^{M} |\Gamma_{k,m} - \mathbf{h}_{k}^{H} \mathbf{w}_{m}^{j} + \lambda_{k,m}^{j}|^{2}$$
 (22a)

subject to
$$\gamma_k \left(\sum_{m \neq m_k} |\Gamma_{k,m}|^2 + \sigma_k^2 \right) - 2\Re \left\{ (\mathbf{w}_{m_k}^{(t)})^H \mathbf{h}_k \Gamma_{k,m_k} \right\} + |\mathbf{h}_k^H \mathbf{w}_{m_k}^{(t)}|^2 \le 0,$$
 (22b)

where m_k is the index of the multicast group that user k belongs to. The superscript j and (t) are the ADMM and CCP iteration counters, respectively. We note that each subproblem (22) is a convex QCQP-1 problem. In various cases, QCQP-1 problems can be solved efficiently [29]. In our case, the optimal solution of subproblem (22) is obtained in closed form. The details are given Appendix A.

2) v *Update*: The update of v in problem (20) is equivalent to solving the problem:

minimize
$$\sum_{m=1}^{M} \|\mathbf{v}_m - \mathbf{w}_m^j + \mathbf{z}_m^j\|_2^2$$
 subject to (10e). (23)

By inspection, this problem can be decomposed into N subproblems, one for each antenna $n \in \mathcal{N}$:

minimize
$$\begin{cases} \sum_{\{v_{n,m}\}_{m=1}^{M}}^{M} & \sum_{m=1}^{M} |v_{n,m} - w_{n,m}^{j} + z_{n,m}^{j}|^{2} \\ \text{subject to} & \sum_{m=1}^{M} |v_{n,m}|^{2} \leq P_{n}, \end{cases}$$
 (24)

where $v_{n,m}$ is the *n*-th element of vector \mathbf{v}_m , $w_{n,m}$ and $z_{n,m}$ are defined in a similar manner. Subproblem (24) is also a QCQP-1 problem. We solve this problem optimally in closed form. In specific, we rewrite subproblem (24) as

minimize
$$\|\tilde{\mathbf{v}}_n - \tilde{\mathbf{u}}_n^j\|_2$$
 (25) subject to $\|\tilde{\mathbf{v}}_n\|_2 \le \sqrt{P_n}$,

where $\tilde{\mathbf{v}}_n = [v_{n,1}, \dots, v_{n,M}]^T \in \mathbb{C}^M$ and $\tilde{\mathbf{u}}_n^j = [(w_{n,1}^j - z_{n,1}^j), \dots, (w_{n,M}^j - z_{n,M}^j)]^T \in \mathbb{C}^M$. It is clear that problem (25) can be viewed as the Euclidean projection of the point $\tilde{\mathbf{u}}_n^j$ onto an Euclidean ball, centered at the original point with radius of $\sqrt{P_n}$. The optimal solution is thus given by

$$\tilde{\mathbf{v}}_n = \min\left\{\frac{\sqrt{P_n}}{\|\tilde{\mathbf{u}}_n^j\|_2}, 1\right\} \tilde{\mathbf{u}}_n^j. \tag{26}$$

Here we give an intuitive geometric explanation for the solution (26). If the point $\tilde{\mathbf{u}}_n^j$ is already inside the ball (i.e., $\|\tilde{\mathbf{u}}_n^j\|_2 \leq \sqrt{P_n}$), then the point $\tilde{\mathbf{u}}_n^j$ itself is what we want. Otherwise, we simply scale it to have Euclidean norm equal to $\sqrt{P_n}$.

3) w Update: The problem (20) for updating the second block of variables w can be decomposed into M independent unconstrained quadratic programing problems, one for each group $m \in \mathcal{M}$:

$$\mathbf{w}_{m}^{j+1} := \arg\min_{\mathbf{w}_{m}} \left\{ \|\mathbf{w}_{m}\|_{2}^{2} + \frac{\rho}{2} \sum_{k=1}^{K} |\Gamma_{k,m}^{j+1} - \mathbf{h}_{k}^{H} \mathbf{w}_{m} + \lambda_{k,m}^{j}|^{2} + \frac{\rho}{2} \|\mathbf{v}_{m}^{j+1} - \mathbf{w}_{m} + \mathbf{z}_{m}^{j}\|_{2}^{2} \right\}.$$
(27)

The solution is given in closed form as

$$\mathbf{w}_{m}^{j+1} = \left((2+\rho)\mathbf{I}_{N} + \rho \sum_{k} \mathbf{h}_{k} \mathbf{h}_{k}^{H} \right)^{-1} \left(\rho \sum_{k} \mathbf{h}_{k} (\Gamma_{k,m}^{j+1} + \lambda_{k,m}^{j}) + \rho (\mathbf{v}_{m}^{j+1} + \mathbf{z}_{m}^{j}) \right). \tag{28}$$

Note that the most computational intensive operation in Alg. 1 is the matrix inversion in w update (28), whose complexity is $\mathcal{O}(N^3)$. However, this operation only needs to be computed once for each channel realization and the solution can be readily used in the subsequent iterations.

Up to now, the closed-form expressions for all the updating steps in Alg. 1 have been derived. The ADMM iteration can converge to a global optimum of problem (10) [32, Proposition 4.2]. Note that Alg. 1 is guaranteed to converge for any initial point. We propose to initialize \mathbf{w}_m^0 using the solution obtained in the previous CCP iteration (i.e., $\mathbf{w}_m^{(t)}$) in Alg. 1. It is a warm start and can often speed up the convergence.

Remark 1: The main novelty of our proposed ADMM algorithm lies in the design of the auxiliary variables $\{\Gamma_{k,m}\}$ and $\{\mathbf{v}_m\}$, with which, problem $\mathcal{Q}^{(t)}$ can be reformulated into such a form that each updating step in ADMM is decomposable and thus can be updated through parallelly solving multiple subproblems with much smaller sizes. Specifically, the variables $\{\Gamma_{k,m}\}$ enable that the SINR constraints in (7b) with coupled variables is transformed into the new constraint in (10d) with decoupled variables. As a result, problem (21) can be decomposed into K QCQP-1 subproblems and then solved efficiently in closed form. The introduction of variables $\{\mathbf{v}_m\}$ makes it possible that the w update in (16) can be decomposed into M small-scale unconstrained problems and carried out in closed form. Meanwhile, the update of $\{\mathbf{v}_m\}$ itself can also be efficiently performed through an Euclidean projection onto an Euclidean ball given in (26).

Remark 2: The consensus ADMM proposed in [29] can be applied to optimally solve problem $\mathcal{Q}^{(t)}$ too. The details are given in Appendix B. However, since a local copy of the global optimization variables w is introduced for each single quadratic constraint in (7b) and (7c), yielding a total of $KMN + MN^2$ auxiliary variables, the problem size of each local variable update is still as large as the original problem $\mathcal{Q}^{(t)}$. In comparison, the number of introduced auxiliary variables in our proposed ADMM in Alg. 1 is only KM + NM, which is much smaller. This is especially useful when the number of antennas N is large. Numerical simulations also demonstrate the complexity advantage of our proposed ADMM. One may also apply the consensus ADMM in [29] to directly solve the original QoS problem \mathcal{Q} in (4) without using CCP. The iteration, however, is not guaranteed to converge due to the non-convexity of problem \mathcal{Q} . Numerical results in Section V show that such method performs poorly for the considered multi-group multicast beamforming problem.

C. Initialization of CCP Algorithm

The CCP-based algorithms generally need a feasible starting point, which is difficult to obtain in general. In [18], the authors propose to initialize the CCP with a feasible point found by SDR with Gaussian randomization. This method works well in simulation, though the extra complexity of SDR method is high. In [29], the authors find a feasible point for initialization using the same consensus ADMM for the general QCQP problems. In this work, we find a starting (not necessary feasible due to the per-antenna peak power constraints) point for initialization of the CCP-ADMM algorithm also using ADMM.

We first formulate a feasibility problem that takes into account the non-convex SINR constraints (4b) only and ignores the per-antenna peak power constraints (4c) as:

find
$$\{\Gamma, \mathbf{w}\}$$
 (29a)

such that
$$\Gamma_{k,m} - \mathbf{h}_k^H \mathbf{w}_m = 0, \ \forall k \in \mathcal{K}, m \in \mathcal{M},$$
 (29b)

$$\gamma_k \left(\sum_{j \neq m} |\Gamma_{k,j}|^2 + \sigma_k^2 \right) - |\Gamma_{k,m}|^2 \le 0, \ \forall k \in \mathcal{G}_m, \forall m \in \mathcal{M}.$$
 (29c)

By applying ADMM, the iterative steps are given by

$$\Gamma^{j+1} := \arg\min_{\Gamma} \left\{ \sum_{k=1}^{K} \sum_{m=1}^{M} |\Gamma_{k,m} - \mathbf{h}_k^H \mathbf{w}_m^j + \lambda_{k,m}^j|^2 \right\}$$

$$(30)$$

$$\mathbf{w}^{j+1} := \arg\min_{\mathbf{w}} \left\{ \sum_{k=1}^{K} \sum_{m=1}^{M} |\Gamma_{k,m}^{j+1} - \mathbf{h}_{k}^{H} \mathbf{w}_{m} + \lambda_{k,m}^{j}|^{2} + \sum_{m=1}^{M} ||\mathbf{v}_{m}^{j+1} - \mathbf{w}_{m} + \mathbf{z}_{m}^{j}||_{2}^{2} \right\},$$
(31)

$$\lambda_{k,m}^{j+1} := \lambda_{k,m}^j + (\Gamma_{k,m}^{j+1} - \mathbf{h}_k^H \mathbf{w}_m^{j+1}), \ \forall k \in \mathcal{K}, m \in \mathcal{M},$$
(32)

which are independent of the penalty parameter ρ .

The Γ update in (30) is similar to (19), which can also be decomposed into K subproblems, one for each $k \in \mathcal{K}$:

minimize
$$\sum_{\{\Gamma_{k,m}\}_{m=1}^{M}}^{M} \sum_{m=1}^{M} |\Gamma_{k,m} - \mathbf{h}_{k}^{H} \mathbf{w}_{m}^{j} + \lambda_{k,m}^{j}|^{2}$$
 (33a)

subject to
$$\gamma_k \left(\sum_{m \neq m_k} |\Gamma_{k,m}|^2 + \sigma_k^2 \right) - |\Gamma_{k,m_k}|^2 \le 0.$$
 (33b)

Each subproblem (33) is still a QCQP-1 problem, but non-convex, in contrast to (22) which is convex. Nevertheless, the strong duality still holds, and problem (33) can be solved optimally irrespective of the non-convexity [33, Appendix B]. Its closed-form solution is given in Appendix C.

The update of w in problem (31) is given by

$$\mathbf{w}_{m}^{j+1} = \left(\mathbf{I}_{N} + \sum_{k} \mathbf{h}_{k} \mathbf{h}_{k}^{H}\right)^{-1} \left(\sum_{k} \mathbf{h}_{k} (\Gamma_{k,m}^{j+1} + \lambda_{k,m}^{j}) + (\mathbf{v}_{m}^{j+1} + \mathbf{z}_{m}^{j})\right), \ \forall m \in \mathcal{M}.$$
(34)

With random initialization, the ADMM iteration terminates until a point, denoted as $\mathbf{w}^{(0)}$ that satisfies the SINR constraints is found. This type of iteration usually converges very fast, if such a point $\mathbf{w}^{(0)}$ exists. Note that, with different initializations, the ADMM iteration may find different $\mathbf{w}^{(0)}$. If it fails to find the point even after a large number of trials, then we have, to some extent, numerical evidence that the problem may be infeasible. Once $\mathbf{w}^{(0)}$ is obtained, it can be used as a starting point for the CCP at iteration t=0.

Remark 3: The obtained starting point $\mathbf{w}^{(0)}$ is not necessarily a feasible point of the original QoS problem in (4) since it may not satisfy the per-antenna peak power constraints (4c). However, the per-antenna peak power constraints will automatically be satisfied by the subsequent $\mathbf{w}^{(t)}$ for all iterations t > 1 because of the constraint (7c). One might find a truly feasible starting point using the same ADMM by including the per-antenna peak power constraints (4c) into the above feasibility problem (29). Our numerical results, however, suggest that such ADMM method can hardly find a feasible point with stringent per-antenna power constraints, even though the feasible set is non-empty.

Remark 4: It might be possible that the problem $Q^{(t)}$ for t=0 is infeasible at the obtained starting point $\mathbf{w}^{(0)}$. If this happens, we try multiple $\mathbf{w}^{(0)}$ until this problem is feasible. Note that multiple $\mathbf{w}^{(0)}$ can be obtained by using ADMM iteration (30)-(32) with different initializations.

In a special case, we can find an alternative starting point with closed-form expression, as given in the following lemma.

Lemma 1: If the aggregate channel matrix $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_K] \in \mathbb{C}^{N \times K}$ has full column rank, then there always exists a starting point that satisfies the SINR constraints (4b) and it is given by

$$\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_M] = \mathbf{H}(\mathbf{H}^H \mathbf{H})^{-1} \mathbf{A}, \tag{35}$$

where A is a $K \times M$ complex-valued matrix with the (k, m)-th element defined as

$$A_{k,m} = \begin{cases} \sqrt{\gamma_k \sigma_k^2} e^{j\theta_k}, & \text{if } k \in \mathcal{G}_m, \\ 0, & \text{otherwise.} \end{cases}$$
 (36)

where $\theta_k \in [0, 2\pi]$ is arbitrary. Further, if the starting point in (35) satisfies the per-antenna peak power constraints (4c), it is also a feasible solution of problem Q.

This lemma guarantees that a starting point that satisfies the SINR constraints exists with probability one and has a closed-form expression if $N \geq K$ and the entries of each h_k are independent and identically distributed. The lemma also provides a sufficient condition for checking the feasibility of problem Q.

IV. EXTENSION TO MMF BEAMFORMING DESIGN

In this section, we extend the CCP-ADMM algorithm proposed for the QoS problem to the MMF problem.

A. Duality of MMF Problem \mathcal{F}

In [9], without the per-antenna peak power constraints, the authors prove that the QoS problem and MMF problem are a dual pair. Including the per-antenna peak power constraints, a similar duality is established in [12]. We restate the conclusions in [12] as follows. Consider a per-antenna power minimization problem [12]

$$\mathcal{P}: \underset{\mathbf{w}}{\mathsf{minimize}} \quad r \tag{37a}$$

subject to
$$\frac{|\mathbf{h}_{k}^{H}\mathbf{w}_{m}|^{2}}{\sum_{j\neq m}|\mathbf{h}_{k}^{H}\mathbf{w}_{j}|^{2}+\sigma_{k}^{2}} \geq g_{k}, \ \forall k \in \mathcal{G}_{m}, \forall m \in \mathcal{M},$$

$$1 \quad \frac{M}{2} \quad (37b)$$

$$\frac{1}{P_n} \sum_{m=1}^{M} \mathbf{w}_m^H \mathbf{R}_n \mathbf{w}_m \le r, \ \forall n \in \mathcal{N}.$$
 (37c)

Following the same notation in the literature [9], [12], let $\mathbf{g} = [g_1, \dots, g_K]^T \in \mathbb{R}^K$ be the user target SINR weight vector and $\mathbf{p} = [P_1, \dots, P_N]^T \in \mathbb{R}^N$ be the per-antenna peak power vector. Problems \mathcal{F} and \mathcal{P} receive \mathbf{g} and \mathbf{p} as inputs and the optimal objective values are denoted as $t^* = \mathcal{F}(\mathbf{g}, \mathbf{p})$ and $r^* = \mathcal{P}(\mathbf{g}, \mathbf{p})$, respectively. Then we have the following claims:

Claim 1 (Claim 2 in [12]): Problems \mathcal{F} and \mathcal{P} are related as

$$1 = \mathcal{P}(\mathcal{F}(\mathbf{g}, \mathbf{p}) \cdot \mathbf{g}, \mathbf{p}), \tag{38}$$

$$t = \mathcal{F}(\mathbf{g}, \mathcal{P}(t \cdot \mathbf{g}, \mathbf{p}) \cdot \mathbf{p}). \tag{39}$$

Claim 2: The optimum objective value of problem $\mathcal{P}(t \cdot \mathbf{g}, \mathbf{p})$ is monotonically nondecreasing in t, for a given \mathbf{g} and \mathbf{p} .

Proof: The feasible region of problem $\mathcal{P}(t \cdot \mathbf{g}, \mathbf{p})$ becomes smaller as t increases, thus completing the proof.

By the above claims, a solution to the MMF problem $\mathcal{F}(\mathbf{g},\mathbf{p})$ can be found by iteratively solving $\mathcal{P}(t \cdot \mathbf{g},\mathbf{p})$ through adjusting the value of t. Claim 1 guarantees the optimality of the solution for $1 = \mathcal{P}(t^* \cdot \mathbf{g},\mathbf{p})$ and Claim 2 enables the use of a simple one-dimensional bisection search for the sought t. Given the non-negativity of t, a lower bound of t can be set as L = 0. An upper bound of t can be obtained by transmitting all the available power $P_{\text{all}} = \sum_{n=1}^{N} P_n$ towards the user with the best channel condition, namely, $U = \max_{k \in \mathcal{K}} \{(P_{\text{all}} || \mathbf{h}_k ||_2^2)/(g_k \sigma_k^2)\}$.

B. Fast Algorithm for Solving Problem \mathcal{P}

We first introduce the following auxiliary variables:

$$\Gamma_{k,m} = \mathbf{h}_k^H \mathbf{w}_m, \ \forall k \in \mathcal{K}, m \in \mathcal{M},$$
 (40)

$$\mathbf{v}_m = \mathbf{w}_m, \ \forall m \in \mathcal{M},\tag{41}$$

$$\alpha_n = r, \ \forall n \in \mathcal{N},\tag{42}$$

where α_n represents the local copy of power weight r for the n-th antenna. To facilitate the notation, we define $\alpha \triangleq \{\alpha_n | n \in \mathcal{N}\}.$

By adopting CCP, in each iteration, we solve the following problem $\mathcal{P}^{(t)}$ with ADMM reformulation

$$\mathcal{P}^{(t)}: \underset{\{\Gamma, \mathbf{v}, \mathbf{w}, \alpha, r\}}{\text{minimize}} \quad r \tag{43a}$$

subject to
$$\Gamma_{k,m} - \mathbf{h}_k^H \mathbf{w}_m = 0, \ \forall k \in \mathcal{K}, m \in \mathcal{M},$$
 (43b)

$$\mathbf{v}_m - \mathbf{w}_m = 0, \ \forall m \in \mathcal{M}, \tag{43c}$$

$$\alpha_n - r = 0, \ \forall n \in \mathcal{N},\tag{43d}$$

$$g_k \left(\sum_{j \neq m} |\Gamma_{k,j}|^2 + \sigma_k^2 \right) - 2\mathfrak{Re} \left\{ (\mathbf{w}_m^{(t)})^H \mathbf{h}_k \Gamma_{k,m} \right\}$$

$$+ |\mathbf{h}_k^H \mathbf{w}_m^{(t)}|^2 \le 0, \ \forall k \in \mathcal{G}_m, \forall m \in \mathcal{M},$$
(43e)

$$\frac{1}{P_n} \sum_{m=1}^{M} \mathbf{v}_m^H \mathbf{R}_n \mathbf{v}_m \le \alpha_n, \ \forall n \in \mathcal{N}.$$
(43f)

To solve problem $\mathcal{P}^{(t)}$ using ADMM, the two blocks of variables $\{\Gamma, \mathbf{v}, \alpha\}$ and $\{\mathbf{w}, r\}$ are alternatively updated. In particular, the ADMM updates (scaled form) for problem (43) are given by Alg. 2, where $\boldsymbol{\mu} \triangleq \{\mu_n | n \in \mathcal{N}\}$ is the dual variable for constraints (43d).

Note that the updates of Γ and $\{v, \alpha\}$ can be separated, both of which can be further decomposed into multiple QCQP-1 subproblems with much smaller size, and performed in closed form. The procedure for $\{v, \alpha\}$ update is similar to Γ update in Section III-B1, thus is omitted here. The updates of w in (46) and r in (47) can be obtained in closed form through solving two unconstrained quadratic programs, which are given by (34) and (51), respectively.

$$r = \frac{1}{N} \sum_{n=1}^{N} (\alpha_n^{j+1} + \mu_n^j) - \frac{1}{N\rho}.$$
 (51)

Algorithm 2 ADMM for solving problem $\mathcal{P}^{(t)}$

Initialization: Initialize $\mathbf{w}_m^0 \leftarrow \mathbf{w}_m^{(t)}, \mathbf{z}_m^0 \leftarrow \mathbf{0}, \lambda_{k,m}^0 \leftarrow 0, \mu_n^0 \leftarrow 0, \forall m \in \mathcal{M}, k \in \mathcal{K}, n \in \mathcal{N}, r^0 \leftarrow 0, m \in \mathcal{M}, k \in \mathcal{K}, k \in \mathcal{$ $\max_{n \in \mathcal{N}} \left\{ \frac{1}{P_n} \sum_{m=1}^{M} (\mathbf{w}_m^{(t)})^H \mathbf{R}_n \mathbf{w}_m^{(t)} \right\}$, and $j \leftarrow 0$. Set the penalty parameter ρ .

1) Update the first block of variables $\{\Gamma^{j+1}, \mathbf{v}^{j+1}, \boldsymbol{\alpha}^{j+1}\}$

$$\Gamma^{j+1} := \arg\min_{\Gamma} \left\{ \sum_{k=1}^{K} \sum_{m=1}^{M} |\Gamma_{k,m} - \mathbf{h}_k^H \mathbf{w}_m^j + \lambda_{k,m}^j|^2 \right\}$$
s.t. (43e),

$$\{\mathbf{v}^{j+1}, \boldsymbol{\alpha}^{j+1}\} := \arg\min_{\mathbf{v}, \boldsymbol{\alpha}} \sum_{m=1}^{M} \|\mathbf{v}_m - \mathbf{w}_m^j + \mathbf{z}_m^j\|_2^2 + \frac{\rho}{2} \sum_{n=1}^{N} |\alpha_n - r^j + \mu_n^j|^2$$
s.t. (43f).

2) Update the second block of variables $\{\mathbf{w}^{j+1}, r^{j+1}\}$

$$\mathbf{w}^{j+1} := \arg\min_{\mathbf{w}} \left\{ \sum_{k=1}^{K} \sum_{m=1}^{M} |\Gamma_{k,m}^{j+1} - \mathbf{h}_{k}^{H} \mathbf{w}_{m} + \lambda_{k,m}^{j}|^{2} + \sum_{m=1}^{M} ||\mathbf{v}_{m}^{j+1} - \mathbf{w}_{m} + \mathbf{z}_{m}^{j}||_{2}^{2} \right\},$$
(46)

$$r^{j+1} := \arg\min_{r} \left\{ r + \frac{\rho}{2} \sum_{n=1}^{N} |\alpha_n^{j+1} - r + \mu_n^j|^2 \right\}.$$
 (47)

3) Update the dual variables $\{\lambda^{j+1}, \mathbf{z}^{j+1}, \boldsymbol{\mu}^{j+1}\}$

$$\lambda_{k,m}^{j+1} := \lambda_{k,m}^j + (\Gamma_{k,m}^{j+1} - \mathbf{h}_k^H \mathbf{w}_m^{j+1}), \ \forall k \in \mathcal{K}, m \in \mathcal{M},$$

$$(48)$$

$$\mathbf{z}_{m}^{j+1} := \mathbf{z}_{m}^{j} + (\mathbf{v}_{m}^{j+1} - \mathbf{w}_{m}^{j+1}), \ \forall m \in \mathcal{M},$$

$$\mu_{n}^{j+1} := \mu_{n}^{j} + (\alpha_{n}^{j+1} - r^{j+1}), \ \forall n \in \mathcal{N}.$$
(49)

$$\mu_n^{j+1} := \mu_n^j + (\alpha_n^{j+1} - r^{j+1}), \ \forall n \in \mathcal{N}.$$
 (50)

4) Set $j \leftarrow j + 1$.

Until convergence criterion is met.

A starting point $\mathbf{w}^{(0)}$ satisfying the SINR constraints can be obtained as suggested in Section III-C. Together with $r^{(0)} = \max_{n \in \mathcal{N}} \frac{1}{P_n} \sum_{m=1}^{M} (\mathbf{w}_m^{(0)})^H \mathbf{R}_n \mathbf{w}_m^{(0)}$, $\{\mathbf{w}^{(0)}, r^{(0)}\}$ is a feasible point of the original problem \mathcal{P} .

C. On the Extension to NUM Problems

Our proposed optimization framework can be extended to general network utility maximization (NUM) problems in the multicast transmission scenario but not necessarily in an optimal way. Similar to the unicast scenario as in [27, Section II-C], the NUM problem in multicast transmission is also a monotonic optimization problem for a general strictly increasing utility (e.g., weighted sum-rate) and thus can be solved using the polyblock outer approximation algorithm or the branch-reduce-and-bound algorithm through solving a series of feasibility subproblems with given SINR targets [34]. Different from unicast transmission, the feasibility subproblems with given SINR targets in the multicast transmission are nonconvex in general. However, we can still apply the proposed ADMM to determine the feasibility of the subproblems numerically and approximately.

V. SIMULATION RESULTS

In this section, the performance and complexity of our fast algorithm CCP-ADMM are demonstrated via numerical simulations. The following baselines are selected,

- **SDR-GauRan:** the SDR method with Gaussian randomization proposed in [9], but modified slightly to cope with the per-antenna peak power constraints. The relaxed SDP problem is solved using the CVX package via interior-point solver SDPT3. If the solution is not rank-one, 200 Gaussian randomization samples are adopted.
- **FPP-SCA:** the feasible point pursuit SCA algorithm proposed in [16], [17], which is equivalent to the CCP algorithm in [18] except the initialization. The convex subproblems are solved using CVX via SDPT3 solver.
- ConADMM: the consensus ADMM in [29] that is directly applied to solve the original non-convex QCQP problem without using CCP as mentioned in Remark 2. The penalty parameter is set as $\rho = 30$ to prevent $\{\mathbf{w}^j\}$ from diverging towards infeasibility. If the ADMM iteration cannot converge to a feasible point within 3000 iterations, a scaling operation similar to [9] is then followed to refine the ADMM solution. If the scaling problem is still infeasible, we claim that ConADMM fails.
- CCP-ConADMM: the CCP method with each subproblem solved using the consensus ADMM in [29] as mentioned in Remark 2 with details given in Alg. 3.
- CCP-CVX-SCS: the CCP method with each subproblem first transformed into the standard cone programming form via CVX then solved using the general ADMM proposed in [27], [28] via SCS solver. The maximum iteration and the convergence tolerance are set to be 3000 and 10⁻⁶, respectively. Other settings are the same as in [28].

For the CCP (or SCA)-based algorithms (namely, FPP-SCA, CCP-ConADMM, CCP-CVX-SCS, CCP-ADMM), the iteration stops when the relative decrease of the objective $\frac{|p^{(t+1)}-p^{(t)}|}{p^{(t)}}$ is less than 10^{-3} or a maximum of 30 iterations is reached, where $p^{(t)}$ denotes the objective value of the t-th iteration. For CCP-ADMM and CCP-ConADMM, the penalty parameters in Alg. 1 and Alg. 3 are set as $\rho = \frac{2}{\sqrt{N}}$ and $\rho = \frac{10}{\sqrt{N}}$, respectively, which are empirically found to work very well (converge fast). The convergence criterion of the ADMM is set as suggested in [26, Section 3.3] with the absolute tolerance $\epsilon^{abs} = 10^{-6}$ and relative tolerance $\epsilon^{rel} = 10^{-6}$. If the ADMM fails to converge within 3000 iterations, we claim the problem is infeasible.

The downlink channels of all users are assumed to be independent and follow the standard complex Gaussian distribution $\mathcal{CN}(0,1)$. The noise variance is set to $\sigma_k^2=1, \forall k$. The SINR target for each user is $\gamma_k=10\text{dB}, \forall k$. All experiments are carried out on a Windows x64 machine with 3.3 GHz CPU and 24 GB of RAM. The plots are obtained after averaging over 100 channel realizations. The performance of our fast algorithm is mainly validated for the QoS problem as similar conclusions can be extended to the MMF problem.

A. Convergence Analysis

We first demonstrate the convergence behavior of the proposed CCP-ADMM algorithm. The algorithm involves an outer-loop iteration for CCP and an inner-loop iteration for ADMM, whose convergences are illustrated in Fig. 1 and Fig. 2, respectively. Here, in Fig. 2, the relative error in each iteration is defined as $\frac{|a^j-a^*|}{a^*}$, where a^j is the objective value of the j-th iteration and a^* is the optimal objective value of the given problem $\mathcal{Q}^{(t)}$ ². The number of transmit antennas is fixed as N=100, and the number of users K varies from 60 to 140. All the users are equally divided into M=4 multicast groups. The maximum transmit power for each antenna is set to $P_n=40 \mathrm{dBm}, \forall n \in \mathcal{N}$. From Fig. 1, it is observed that the CCP algorithm converges to modest accuracy, e.g., relative decrease of 10^{-2} within 15 iterations and then higher accuracy as the iterations progress. From Fig. 2, it is seen that the proposed ADMM algorithm converges

² The optimal objective value a^* can be obtained using interior-point methods.

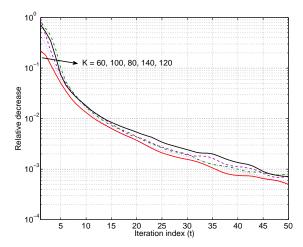


Fig. 1: Convergence behavior of the CCP algorithm in the outer loop (N = 100).

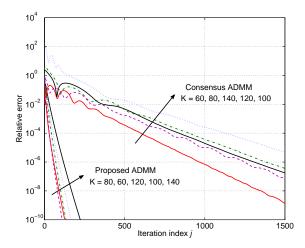


Fig. 2: Convergence behavior of the ADMM algorithm in the inner loop (N = 100).

to a high accuracy solution, e.g., relative error of 10^{-4} , within 100 iterations for various K, while the existing consensus ADMM needs about 700-1200 iterations to achieve the same level of accuracy.

B. Performance and Complexity Comparison

1) Comparison with varying number of users: With the same scenario as in Section V-A, Fig. 3 and Fig. 4 compare the performance in terms of the actual transmit power and the computational complexity in terms of the actual simulation running time, respectively. We first observe that the conventional SDR-GauRan method can hardly find a feasible solution where K>100, and its gap to the SDR lower bound is as large as 12dB when K=100. The reason is that the rank-one probability of SDP solutions is very low in the considered scenario. We also see that ConADMM performs much worse than our proposed CCP-ADMM in terms of power efficiency and simulation running time. This is mainly due to that ConADMM fails to converge to a feasible solution and the refinement must be performed most of the time.

From Fig. 3, it is also seen that the three CCP-based algorithms, namely, CCP-ADMM, CCP-ConADMM, and CCP-CVX-SCS, achieve the same favorable performance as FPP-SCA, as expected, and their solutions are within 1dB close to the SDR lower bound. Their computational complexities are however different. It is seen from Fig. 4 that the running time of CCP-ADMM is about $70 \sim 170$ times faster than FPP-SCA and CCP-ConADMM. It is also $30 \sim 100$ times faster than CCP-CVX-SCS ³. This is mainly because

³ The running time of CCP-CVX-SCS includes the modeling time for transforming the original problem instance into the standard form using CVX and the solving time for calling SCS solver.

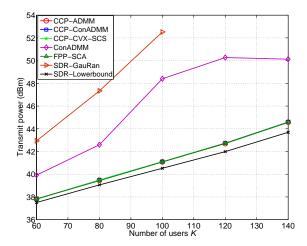


Fig. 3: Transmit power versus the number of users K (N = 100).

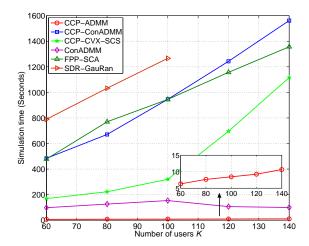


Fig. 4: Simulation time versus the number of users K (N = 100).

the consensus ADMM in [29] and the general ADMM in [27], [28] need to introduce a large number of auxiliary variables to transform the original problem instance into certain standard forms, the dimension of which are much larger than the original problem. In addition, the consensus ADMM numerically converges much slower than our proposed ADMM, as shown in Fig. 2. The modeling time for transforming the original problem instance into the standard form in CCP-CVX-SCS is also non-negligible.

We would like to remark that since the consensus ADMM in [29] and the general ADMM in [27], [28] deal with general frameworks, they are not really comparable to our proposed tailor-made algorithm. We mainly use them to verify the high-performance and low-complexity of our customized algorithm in terms of the solution quality and the running time via numerical simulation.

2) Comparison with varying number of transmit antennas: The performance and complexity comparison at K=50 users with varying number of transmit antennas are illustrated in Fig. 5 and Fig. 6, respectively. Here, the users are equally divided into M=5 multicast groups. The maximum transmit power for each antenna is set as $P_n=P_{\rm all}/N, \forall n\in\mathcal{N}$, where $P_{\rm all}=57{\rm dBm}$ (500 W). This ensures that the total transmit power keeps constant when the number of antennas increases.

Similar observations as in the previous two figures can be observed. In particular, the CCP-ADMM achieves the same performance as FPP-SCA, which is within 0.3dB close to the SDR lower bound at several orders of magnitude reduction in complexity. It is also observed from Fig. 6 that as the number of transmit antennas increases, the timing curves of CCP-ADMM and CCP-CVX-SCS grow almost linearly, and the curves of the other four baselines grow almost exponentially. Compared with CCP-CVX-SCS, our

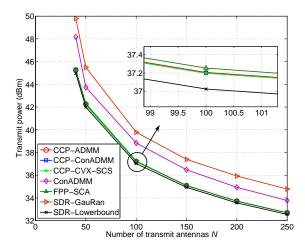


Fig. 5: Transmit power versus the number of transmit antennas N (K = 50).

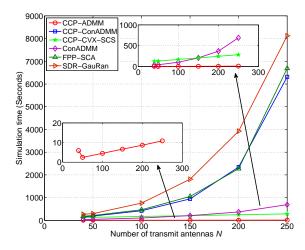


Fig. 6: Simulation time versus the number of transmit antennas N (K = 50).

proposed CCP-ADMM can speedup the running time about $30 \sim 50$ times. Note that the discontinuity of the running time of CCP-ADMM at N=40 is due to that the closed-form starting point (35) in Lemma 1 only holds for $N \geq K (=50)$ and hence the starting point at N=40 has to be found using the ADMM iterations (30)-(32), which is more time-consuming.

3) Comparison for the MMF problem: We finally demonstrate the performance of our fast algorithm for the MMF problem. The obtained minimum SINR (with equal weights) at different per-antenna peak power levels is illustrated in Fig. 7. The corresponding average simulation running time is shown in Table I. Here, N=100 transmit antennas and K=50 users evenly divided into M=5 groups. It is seen that our proposed CCP-ADMM achieves the same performance as CCP-CVX-SCS and FPP-SCA, which is within $0.5 \, \mathrm{dB}$ close to the SDR upper bound, but the running time is more than 10 times faster.

The above numerical results verify that the proposed CCP-ADMM has high performance and extremely low complexity compared with the state-of-the-art algorithms in the literature.

TABLE I: Average simulation time of the MMF problem.

Algorithm	CCP-ADMM	CCP-CVX-SCS	FPP-SCA
Time (sec)	71.53	898.23	905.39

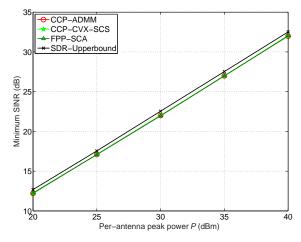


Fig. 7: Minimum SINR versus the per-antenna peak power P ($P_n = P$, $\forall n$).

VI. CONCLUSION

This paper develops a fast algorithm for the multi-group multicast beamforming problem in large-scale wireless systems. The algorithm adopts the CCP method to maintain near-optimal performance and utilizes the ADMM to realize low-complex and parallel implementation. In addition, we propose an efficient ADMM-based method to obtain a starting point for the algorithm. Numerical results verify that the proposed fast algorithm reduces the complexity by multiple orders of magnitude compared with the state-of-the-art algorithms and is very suitable for large-scale wireless systems.

APPENDIX A

FINDING A CLOSED-FORM SOLUTION OF PROBLEM (22)

The Lagrangian of problem (22) is given by

$$\mathcal{L}_{\pi} = \sum_{m=1}^{M} |\Gamma_{k,m} - \mathbf{h}_{k}^{H} \mathbf{w}_{m}^{j} + \lambda_{k,m}^{j}|_{2}^{2}$$

$$+ \pi \left[\gamma_{k} \left(\sum_{m \neq m_{k}} |\Gamma_{k,m}|^{2} + \sigma_{k}^{2} \right) - 2 \mathfrak{Re} \left\{ (\mathbf{w}_{m_{k}}^{(t)})^{H} \mathbf{h}_{k} \Gamma_{k,m_{k}} \right\} + |\mathbf{h}_{k}^{H} \mathbf{w}_{m_{k}}^{(t)}|^{2} \right],$$

$$(52)$$

where $\pi \geq 0$ is the dual variable associated with the inequality constraint (22b). By setting the derivative of \mathcal{L}_{π} with respect to $\{\Gamma_{k,m}\}_{m=1}^{M}$ to zero, i.e., $\partial \mathcal{L}_{\pi}/\partial \Gamma_{k,m} = 0, \forall m \in \mathcal{M}$, the optimal primal variables $\{\Gamma_{k,m}\}_{m=1}^{M}$ can be expressed as

$$\Gamma_{k,m} = \begin{cases} (\mathbf{h}_k^H \mathbf{w}_m^j - \lambda_{k,m}^j) / (\pi \gamma_k + 1), & \text{if } m \neq m_k, m \in \mathcal{M}, \\ \pi \mathbf{h}_k^H \mathbf{w}_m^{(t)} + \mathbf{h}_k^H \mathbf{w}_m^j - \lambda_{k,m}^j, & \text{if } m = m_k. \end{cases}$$
(53)

Next, we find the optimal dual variable π^* . Substituting the solution (53) back into the inequality constraint (22b), we have

$$g(\pi) = \frac{\gamma_k \sum_{m \neq m_k} |\mathbf{h}_k^H \mathbf{w}_m^j - \lambda_{k,m}^j|^2}{(\pi \gamma_k + 1)^2} - 2\pi |\mathbf{h}_k^H \mathbf{w}_{m_k}^{(t)}|^2 + \gamma_k \sigma_k^2 - 2\Re \left\{ (\mathbf{w}_{m_k}^{(t)})^H \mathbf{h}_k (\mathbf{h}_k^H \mathbf{w}_{m_k}^j - \lambda_{k,m_k}^j) \right\} + |\mathbf{h}_k^H \mathbf{w}_{m_k}^{(t)}|^2 \le 0.$$

Taking a closer look at the function $g(\pi)$, we obtain its derivative

$$\nabla g(\pi) = -\frac{2\gamma_k^2 \sum_{m \neq m_k} |\mathbf{h}_k^H \mathbf{w}_m^j - \lambda_{k,m}^j|^2}{(\pi \gamma_k + 1)^3} - 2|\mathbf{h}_k^H \mathbf{w}_{m_k}^{(t)}|^2 < 0.$$

Obviously, there is $\nabla g(\pi) < 0$ (Note that when $\nabla g(\pi) = 0$, problem (22) is infeasible.). Thus, $g(\pi)$ is monotonically decreasing in the region $\pi \geq 0$. According to the complementary slackness condition, we have $\pi^*g(\pi^*)=0$. Therefore, if g(0)<0, we have $\pi^*=0$. Otherwise, $g(\pi)$ has a zero-crossing in the region $\pi \geq 0$, and we have $g(\pi^*)=0$. Thus, π^* can be found efficiently using a root finding method, such as bisection search or Newtons method [29]. In our context, closed-form solution of π^* can be derived. For ease of notation, we let

$$\begin{split} a &\triangleq \gamma_k \sum_{m \neq m_k} |\mathbf{h}_k^H \mathbf{w}_m^j - \lambda_{k,m}^j|^2, \\ b &\triangleq -2|\mathbf{h}_k^H \mathbf{w}_{m_k}^{(t)}|^2, \\ c &\triangleq \gamma_k \sigma_k^2 - 2 \mathfrak{Re} \left\{ (\mathbf{w}_{m_k}^{(t)})^H \mathbf{h}_k (\mathbf{h}_k^H \mathbf{w}_{m_k}^j - \lambda_{k,m_k}^j) \right\} + |\mathbf{h}_k^H \mathbf{w}_{m_k}^{(t)}|^2. \end{split}$$

Then solving $g(\pi) = 0$ is equivalent to solving the cubic equation

$$h(\pi) = (b\gamma_k^2)\pi^3 + (2b\gamma_k + c\gamma_k^2)\pi^2 + (b + 2c\gamma_k)\pi + (c + a) = 0,$$
(54)

of which the closed-form solution can be expressed using the cubic formula. The real non-negative root is guaranteed to be the unique (thus correct) solution for π^* . Then, the optimal $\{\Gamma_{k,m}\}_{m=1}^M$ can be recovered by substituting π^* into equation (53).

APPENDIX B

Consensus-ADMM Algorithm for Problem $\mathcal{Q}^{(t)}$

By adopting the consensus form of ADMM in [29], we introduce a auxiliary variable copy of w for every single quadratic constraint, and problem $Q^{(t)}$ can be equivalently expressed as

$$\underset{\{\mathbf{x}_k, \mathbf{y}_n, \mathbf{w}\}}{\text{minimize}} \quad \sum_{m=1}^M \|\mathbf{w}_m\|_2^2 \tag{55a}$$

subject to
$$\mathbf{x}_{m,k} - \mathbf{w}_m = 0, \ \forall k \in \mathcal{K}, m \in \mathcal{M},$$
 (55b)

$$\mathbf{y}_{m,n} - \mathbf{w}_m = 0, \ \forall n \in \mathcal{N}, m \in \mathcal{M},$$
 (55c)

$$\gamma_k \left(\sum_{j \neq m} |\mathbf{h}_k^H \mathbf{x}_{j,k}|^2 + \sigma_k^2 \right) - 2 \mathfrak{Re} \left\{ (\mathbf{w}_m^{(t)})^H \mathbf{h}_k \mathbf{h}_k^H \mathbf{x}_{m,k} \right\}$$

$$+ |\mathbf{h}_k^H \mathbf{w}_m^{(t)}|^2 \le 0, \ \forall k \in \mathcal{G}_m, \forall m \in \mathcal{M},$$
 (55d)

$$\sum_{m=1}^{M} \mathbf{y}_{m,n}^{H} \mathbf{R}_{n} \mathbf{y}_{m,n} \leq P_{n}, \ \forall n \in \mathcal{N}.$$
 (55e)

where $\mathbf{x}_k \triangleq \{\mathbf{x}_{m,k} \in \mathbb{C}^N | m \in \mathcal{M}\}$ is a set of auxiliary variables for each SINR constraint $k \in \mathcal{K}$ and $\mathbf{y}_n \triangleq \{\mathbf{y}_{m,n} \in \mathbb{C}^N | m \in \mathcal{M}\}$ is a set of auxiliary variables for each per-antenna power constraint $n \in \mathcal{N}$. Then the corresponding ADMM algorithm is summarized in Alg. 3.

The update of variable \mathbf{x}_k for each $k \in \mathcal{K}$ in problem (56) and \mathbf{y}_n for each $n \in \mathcal{N}$ in problem (57) are convex QCQP-1 problems, which are similar to the Γ update in Section III-B1 and \mathbf{v} update in Section III-B2, respectively. The details are therefore omitted here.

APPENDIX C

FINDING A CLOSED-FORM SOLUTION OF PROBLEM (33)

The Lagrangian of problem (33) is given by

$$\mathcal{L}_{\pi} = \sum_{m=1}^{M} |\Gamma_{k,m} - \mathbf{h}_{k}^{H} \mathbf{w}_{m}^{j} + \lambda_{k,m}^{j}|_{2}^{2} + \pi \left[\gamma_{k} \left(\sum_{m \neq m_{k}} |\Gamma_{k,m}|^{2} + \sigma_{k}^{2} \right) - |\Gamma_{k,m_{k}}|^{2} \right],$$
 (61)

Algorithm 3 Consensus-ADMM for solving problem $Q^{(t)}$

Initialization: Initialize $\mathbf{w}_m^0 \leftarrow \mathbf{w}_m^{(t)}, \mathbf{u}_{m,k}^0 \leftarrow \mathbf{0}, \mathbf{v}_{m,n}^0 \leftarrow 0, \forall m \in \mathcal{M}, k \in \mathcal{K}, n \in \mathcal{N}, \text{ and } j \leftarrow 0.$ Set the penalty parameter ρ .

Repeat

1) Update the first block of variables $\{\mathbf{x}_k^{j+1}, \mathbf{y}_n^{j+1}, \forall k \in \mathcal{K}, n \in \mathcal{N}\}$

$$\mathbf{x}_{k}^{j+1} := \arg\min_{\mathbf{x}_{k}} \sum_{m=1}^{M} \|\mathbf{x}_{m,k} - \mathbf{w}_{m}^{j} + \mathbf{u}_{m,k}^{j}\|^{2}$$

$$\text{s.t. } \gamma_{k} \left(\sum_{m \neq m_{k}} |\mathbf{h}_{k}^{H} \mathbf{x}_{m,k}|^{2} + \sigma_{k}^{2} \right) - 2\mathfrak{Re} \left\{ (\mathbf{w}_{m_{k}}^{(t)})^{H} \mathbf{h}_{k} \mathbf{h}_{k}^{H} \mathbf{x}_{m_{k},k} \right\}$$

$$+ |\mathbf{h}_{k}^{H} \mathbf{w}_{m_{k}}^{(t)}|^{2} \leq 0,$$

$$\mathbf{y}_{n}^{j+1} := \arg\min_{\mathbf{y}_{n}} \sum_{m=1}^{M} \|\mathbf{y}_{m,n} - \mathbf{w}_{m}^{j} + \mathbf{v}_{m,n}^{j}\|_{2}^{2}$$

$$\text{s.t. } \sum_{m=1}^{M} \mathbf{y}_{m,n}^{H} \mathbf{R}_{n} \mathbf{y}_{m,n} \leq P_{n}.$$

$$(56)$$

2) Update the second block of variables \mathbf{w}^{j+1}

$$\mathbf{w}_{m}^{j+1} := \frac{\rho}{2 + \rho(K + N)} \left(\sum_{k} (\mathbf{x}_{m,k}^{j+1} + \mathbf{u}_{m,k}^{j}) + \sum_{n} (\mathbf{y}_{m,n}^{j+1} + \mathbf{v}_{m,n}^{j}) \right), \forall m \in \mathcal{M}.$$
 (58)

3) Update the dual variables $\{\mathbf{u}_k^{j+1}, \mathbf{v}_n^{j+1}, \forall k \in \mathcal{K}, n \in \mathcal{N}\}$

$$\mathbf{u}_{m,k}^{j+1} := \mathbf{u}_{m,k}^{j} + \mathbf{x}_{m,k}^{j+1} - \mathbf{w}_{m}^{j+1}, \ \forall m \in \mathcal{M},$$
 (59)

$$\mathbf{v}_{m,n}^{j+1} := \mathbf{v}_{m,n}^{j} + \mathbf{y}_{m,n}^{j+1} - \mathbf{w}_{m}^{j+1}, \ \forall m \in \mathcal{M}.$$
 (60)

4) Set $j \leftarrow j + 1$.

Until convergence criterion is met.

where $\pi \geq 0$ is the dual variable associated with the inequality constraint (33b). If problem (33) is feasible, then from the dual of the problem, there should be $\pi \leq 1$ [33, Appendix B]. Setting the derivative of \mathcal{L}_{π} with respect to $\{\Gamma_{k,m}\}_{m=1}^{M}$ to zero, we have

$$\begin{cases}
(\pi \gamma_k + 1) \Gamma_{k,m} = \mathbf{h}_k^H \mathbf{w}_m^j - \lambda_{k,m}^j, & \text{if } m \neq m_k, m \in \mathcal{M}, \\
(1 - \pi) \Gamma_{k,m} = \mathbf{h}_k^H \mathbf{w}_m^j - \lambda_{k,m}^j, & \text{if } m = m_k.
\end{cases}$$
(62)

If the optimal dual variable π^* satisfies $\pi^* = 1$, we must have $\mathbf{h}_k^H \mathbf{w}_{m_k}^j - \lambda_{k,m_k}^j = 0$. Then $\{\Gamma_{k,m}\}_{m=1}^M$ can be obtained from (62). Otherwise, we have

$$\Gamma_{k,m} = \begin{cases} \left(\mathbf{h}_k^H \mathbf{w}_m^j + \lambda_{k,m}^j\right) / (\pi \gamma_k + 1), & \text{if } m \neq m_k, m \in \mathcal{M}, \\ \left(\mathbf{h}_k^H \mathbf{w}_m^j + \lambda_{k,m}^j\right) / (1 - \pi), & \text{if } m = m_k. \end{cases}$$
(63)

Substituting the solution (63) back into the inequality constraint (33b), we have

$$\tilde{g}(\pi) = \frac{\gamma_k \sum_{m \neq m_k} |\mathbf{h}_k^H \mathbf{w}_m^j - \lambda_{k,m}^j|^2}{(\pi \gamma_k + 1)^2} + \gamma_k \sigma_k^2 - \frac{|\mathbf{h}_k^H \mathbf{w}_{m_k}^j - \lambda_{k,m_k}^j|^2}{(1 - \pi)^2} \le 0.$$

It is easy to verify that $\nabla \tilde{g}(\pi) < 0$ and $\tilde{g}(\pi)$ is monotonically decreasing in the region $0 \le \pi < 1$. According to the complementary slackness condition, we have $\pi^* \tilde{g}(\pi^*) = 0$. Therefore, if $\tilde{g}(0) < 0$, we have $\pi^* = 0$. Otherwise, $\tilde{g}(\pi)$ has a zero-crossing in the region $0 \le \pi < 1$ and we have $\tilde{g}(\pi^*) = 0$. For ease of notation, we let

$$\tilde{a} \triangleq \gamma_k \sum_{m \neq m_k} |\mathbf{h}_k^H \mathbf{w}_m^j - \lambda_{k,m}^j|^2,$$

$$\tilde{b} \triangleq \gamma_k \sigma_k^2,$$

$$\tilde{c} \triangleq |\mathbf{h}_k^H \mathbf{w}_{m_k}^j - \lambda_{k,m_k}^j|^2.$$

Note that $\tilde{g}(\pi) = 0$ is equivalent to solving the quartic equation

$$\tilde{h}(\pi) = (\tilde{b}\gamma_k^2)\pi^4 + (2\tilde{b}\gamma_k - 2\tilde{b}\gamma_k^2)\pi^3 + (\tilde{b}\gamma_k^2 + \tilde{b} - 4\tilde{b}\gamma_k + \tilde{a} - \tilde{c}\gamma_k^2)\pi^2 + (2\tilde{b}\gamma_k - 2\tilde{b} - 2\tilde{a} - 2\tilde{c}\gamma_k)\pi + (\tilde{b} + \tilde{a} - \tilde{c}) = 0,$$

of which the closed-form solution can be expressed using the quartic formula. The root satisfying $0 \le \pi < 1$ is guaranteed to be the unique (thus correct) solution for π^* . Then, the optimal $\{\Gamma_{k,m}\}_{m=1}^M$ can be recovered through equation (63) with the obtained π^* .

APPENDIX D PROOF OF LEMMA 1

To find a point that satisfies the SINR constraints, we propose to solve the following problem with the objective of minimizing $\sum_{k=1}^K \sum_{m=1}^M |\mathbf{h}_k^H \mathbf{w}_m|^2$ subject to the SINR constraints

minimize
$$\sum_{k=1}^{K} \sum_{m=1}^{M} |\mathbf{h}_{k}^{H} \mathbf{w}_{m}|^{2}$$
subject to
$$\gamma_{k} \left(\sum_{j \neq m} |\mathbf{h}_{k}^{H} \mathbf{w}_{j}|^{2} + \sigma_{k}^{2} \right) - |\mathbf{h}_{k}^{H} \mathbf{w}_{m}|^{2} \leq 0, \ \forall k \in \mathcal{G}_{m}, \forall m \in \mathcal{M}.$$
(64)

Let $A_{k,m} = \mathbf{h}_k^H \mathbf{w}_m$, $\forall k \in \mathcal{K}, m \in \mathcal{M}$, then problem (64) can be decomposed into K subproblems, one for each $k \in \mathcal{K}$:

minimize
$$\sum_{\{A_{k,m}\}_{m=1}^{M}}^{M} \sum_{m=1}^{M} |A_{k,m}|^2$$
 (65) subject to $\gamma_k \left(\sum_{m \neq m_k} |A_{k,m}|^2 + \sigma_k^2 \right) - |A_{k,m_k}|^2 \le 0.$

where m_k is the index of the multicast group that user k belongs to. Each subproblem (65) is a QCQP-1 problem, the optimal solution of which is given by

$$A_{k,m} = \begin{cases} \sqrt{\gamma_k \sigma_k^2} e^{j\theta_k}, & \text{if } m = m_k, \\ 0, & \text{otherwise,} \end{cases}$$
 (66)

where $\theta_k \in [0, 2\pi]$ is arbitrary. Gathering the solution (66) for all $k \in \mathcal{K}$, we obtain the equation (36).

Since the channel matrix \mathbf{H} has full column rank, the equation $\mathbf{h}_k^H \mathbf{w}_m = A_{k,m}, \forall k, m$, or equivalently, $\mathbf{H}^H \mathbf{W} = \mathbf{A}$ always holds with $\mathbf{W} = \mathbf{H}(\mathbf{H}^H \mathbf{H})^{-1} \mathbf{A}$, where \mathbf{A} is a $K \times M$ matrix with the (k, m)-th element defined as $A_{k,m}$. Note that the objective function $\sum_{k=1}^K \sum_{m=1}^M |\mathbf{h}_k^H \mathbf{w}_m|^2$ in problem (64) is a key to enable the simple and closed-form solution (66) for each subproblem (65). Replacing it with a different objective function may not result in closed-form solution.

REFERENCES

- [1] E. Chen and M. Tao, "A fast algorithm for multi-group multicast beamforming in large-scale wireless systems," in (to appear) Proc. IEEE International Conference on Communications (ICC), May 2017.
- [2] A. Lee, "eMBMS delivers mobile video to the mass audience," https://www.itu.int/en/ITU-D/Regional-Presence/AsiaPacific/Documents/Events/2015/A 2015.
- [3] T. Lohmar, M. Slssingar, V. Kenehan, and S. Puustinen, "Delivering content with LTE broadcast," *Ericsson Review*, vol. 1, no. 11, Feb. 2013.
- [4] Qualcomm, "LTE broadcast," https://www.qualcomm.com/invention/technologies/lte/broadcast, 2016.
- [5] B. Hu, C. Hua, C. Chen, and X. Guan, "Multicast beamforming for wireless backhaul with user-centric clustering in cloud-RANs," in *Proc. IEEE International Conference on Communications (ICC)*, May 2016, pp. 1–6.
- [6] M. A. Maddah-Ali and U. Niesen, "Fundamental limits of caching," *IEEE Trans. Inf. Theory*, vol. 60, no. 5, pp. 2856–2867, May 2014.
- [7] F. Xu, K. Liu, and M. Tao, "Cooperative Tx/Rx caching in interference channels: A storage-latency tradeoff study," in *Proc. IEEE International Symposium on Information Theory (ISIT)*, Jul. 2016, pp. 2034–2038.
- [8] 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.
- [9] E. Karipidis, N. 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.
- [10] 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.
- [11] ——, "Massive MIMO multicasting in noncooperative cellular networks," *IEEE J. Sel. Areas Commun.*, vol. 32, no. 6, pp. 1180–1193, Jun. 2014.
- [12] D. Christopoulos, S. Chatzinotas, and B. Ottersten, "Weighted fair multicast multigroup beamforming under per-antenna power constraints," *IEEE Trans. Signal Process.*, vol. 62, no. 19, pp. 5132–5142, Oct. 2014.
- [13] B. R. Marks and G. P. Wright, "A general inner approximation algorithm for nonconvex mathematical programs," *Operations Research*, vol. 26, no. 4, pp. 681–683, 1978.
- [14] G. R. Lanckriet and B. K. Sriperumbudur, "On the convergence of the concave-convex procedure," in *Proc. Advances in neural information processing systems*, 2009, pp. 1759–1767.
- [15] L.-N. Tran, M. F. Hanif, and M. Juntti, "A conic quadratic programming approach to physical layer multicasting for large-scale antenna arrays," *IEEE Signal Process. Lett.*, vol. 21, no. 1, pp. 114–117, Jan. 2014.
- [16] O. Mehanna, K. Huang, B. Gopalakrishnan, A. Konar, and N. D. Sidiropoulos, "Feasible point pursuit and successive approximation of non-convex QCQPs," *IEEE Signal Process. Lett.*, vol. 22, no. 7, pp. 804–808, July 2015.
- [17] D. Christopoulos, S. Chatzinotas, and B. Ottersten, "Multicast multigroup beamforming for per-antenna power constrained large-scale arrays," in *Proc. IEEE 16th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, Jun. 2015, pp. 271–275.
- [18] M. Tao, E. Chen, H. Zhou, and W. Yu, "Content-centric sparse multicast beamforming for cache-enabled cloud RAN," *IEEE Trans. Wireless Commun.*, vol. 15, no. 9, pp. 6118–6131, Sept. 2016.
- [19] T. X. Tran and D. Pompili, "Dynamic radio cooperation for downlink cloud-RANs with computing resource sharing," in *Proc. IEEE 12th International Conference on Mobile Ad Hoc and Sensor Systems (MASS)*, Oct. 2015, pp. 118–126.
- [20] S. Mosleh, L. Liu, and J. Zhang, "Proportional-fair resource allocation for coordinated multi-point transmission in LTE-advanced," *IEEE Trans. Wireless Commun.*, vol. 15, no. 8, pp. 5355–5367, Aug. 2016.
- [21] M. Grant and S. Boyd, "CVX: Matlab software for disciplined convex programming, version 2.1," http://cvxr.com/cvx, Mar. 2014.
- [22] P. Rost, C. Bernardos, A. Domenico, M. Girolamo, M. Lalam, A. Maeder, D. Sabella, and D. Wübben, "Cloud technologies for flexible 5G radio access networks," *IEEE Commun. Mag.*, vol. 52, no. 5, pp. 68–76, May 2014.
- [23] F. Rusek, D. Persson, B. K. Lau, E. G. Larsson, T. L. Marzetta, O. Edfors, and F. Tufvesson, "Scaling up MIMO: Opportunities and challenges with very large arrays," *IEEE Signal Process. Mag.*, vol. 30, no. 1, pp. 40–60, Jan. 2013.
- [24] J. G. Andrews, S. Buzzi, W. Choi, S. V. Hanly, A. Lozano, A. C. K. Soong, and J. C. Zhang, "What will 5G be?" *IEEE J. Sel. Areas Commun.*, vol. 32, no. 6, pp. 1065–1082, Jun. 2014.
- [25] B. Gopalakrishnan and N. D. Sidiropoulos, "High performance adaptive algorithms for single-group multicast beamforming," *IEEE Trans. Signal Process.*, vol. 63, no. 16, pp. 4373–4384, Aug. 2015.
- [26] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends* in *Machine Learning*, vol. 3, no. 1, pp. 1–122, Jan. 2011.
- [27] Y. Shi, J. Zhang, B. O'Donoghue, and K. B. Letaief, "Large-scale convex optimization for dense wireless cooperative networks," *IEEE Trans. Signal Process.*, vol. 63, no. 18, pp. 4729–4743, Sept. 2015.
- [28] B. O'Donoghue, E. Chu, N. Parikh, and S. Boyd, "Conic optimization via operator splitting and homogeneous self-dual embedding," *Journal of Optimization Theory and Applications*, vol. 169, no. 3, pp. 1042–1068, Jun. 2016.
- [29] K. Huang and N. D. Sidiropoulos, "Consensus-ADMM for general quadratically constrained quadratic programming," *IEEE Trans. Signal Process.*, vol. 64, no. 20, pp. 5297–5310, Oct. 2016.
- [30] C. Shen, T. H. Chang, K. Y. Wang, Z. Qiu, and C. Y. Chi, "Distributed robust multicell coordinated beamforming with imperfect CSI: An ADMM approach," *IEEE Trans. Signal Process.*, vol. 60, no. 6, pp. 2988–3003, Jun. 2012.
- [31] W. C. Liao, M. Hong, Y. F. Liu, and Z. Q. Luo, "Base station activation and linear transceiver design for optimal resource management in heterogeneous networks," *IEEE Trans. Signal Process.*, vol. 62, no. 15, pp. 3939–3952, Aug. 2014.
- [32] D. P. Bertsekas and J. N. Tsitsiklis, Parallel and Distributed Computation: Numerical Methods. New York, NY, USA: Athena Scientific, 1997.
- [33] S. Boyd and L. Vandenberghe, Convex Optimization. Cambridge University Press, 2004.

[34] E. Björnson and E. Jorswieck, "Optimal resource allocation in coordinated multi-cell systems," *Foundations and Trends*® in *Communications and Information Theory*, vol. 9, no. 2-3, pp. 113–381, Jan. 2013.