Ultra-Low-Complexity Algorithms with Structurally Optimal Multi-Group Multicast Beamforming in Large-Scale Systems

Chong Zhang, Student Member, IEEE, Min Dong, Senior Member, IEEE, and Ben Liang, Fellow, IEEE

Abstract-In this work, we propose ultra-low-complexity design solutions for multi-group multicast beamforming in large-scale systems. For the quality-of-service (QoS) problem, by utilizing the optimal multicast beamforming structure obtained recently in [2], we convert the original problem into a non-convex weight optimization problem of a lower dimension and propose two fast first-order algorithms to solve it. Both algorithms are based on successive convex approximation (SCA) and provide fast iterative updates to solve each SCA subproblem. The first algorithm uses a saddle point reformulation in the dual domain and applies the extragradient method with an adaptive step-size procedure to find the saddle point with simple closed-form updates. The second algorithm adopts the alternating direction method of multipliers (ADMM) method by converting each SCA subproblem into a favorable ADMM structure. The structure leads to simple closedform ADMM updates, where the problem in each update block can be further decomposed into parallel subproblems of small sizes, for which closed-form solutions are obtained. We also propose efficient initialization methods to obtain favorable initial points that facilitate fast convergence. Furthermore, taking advantage of the proposed fast algorithms, for the max-min fair (MMF) problem, we propose a simple closed-form scaling scheme that directly uses the solution obtained from the QoS problem, avoiding the conventional computationally expensive method that iteratively solves the inverse QoS problem. We further develop lower and upper bounds on the performance of this scaling scheme. Simulation results show that the proposed algorithms offer near-optimal performance with substantially lower computational complexity than the state-of-theart algorithms for large-scale systems.

Index Terms—Multicast beamforming, optimal beamforming structure, large-scale optimization, extragradient algorithm, alternating direction method of multipliers, low complexity.

I. INTRODUCTION

Content distribution through wireless multicasting has been increasingly popular among new wireless services and applications and is expected to dominate future wireless traffic. Multiantenna multicast beamforming is an efficient transmission technique to support high-speed content distribution to multiple user groups simultaneously. With massive multiple-input multipleoutput (MIMO) being the essential technology for future networks [3], it is critical for multicast beamforming solutions to be scalable to meet the ultra-low complexity requirement for large-scale systems.

Downlink multicast beamforming has been studied for traditional multi-antenna systems in various scenarios, such as single-group or multi-group multicasting [4]–[6], multi-cell networks [7], [8], relay networks [9], [10], and cognitive radio networks [11], [12]. Multicast beamforming problems are challenging to solve, as they are generally non-convex and NP-hard [4]. Semi-definite relaxation (SDR) has been a prior state-of-the-art method considered in the existing works [13]. For traditional multi-antenna systems with relatively small problem sizes, SDR provides a good approximate solution [4]– [6], [8]. However, SDR is not a scalable method for largescale problems, where the computational complexity becomes very high and the performance deteriorates significantly as the problem size grows.

With the increasing number of transmit antennas, successive convex approximation (SCA) [14] has become a more attractive approach for solving the multi-group multicast beamforming problems due to its computational and performance advantages over SDR [15]–[17]. SCA-based methods convexify the nonconvex problem into a sequence of convex approximation subproblems, where the convex subproblems are typically solved by the interior-point method (IPM) [18]. However, IPM is a second-order algorithm. For computing a multicast beamforming solution in large-scale massive MIMO systems, using IPM still results in a relatively high computational complexity.

Following this, several methods have been proposed to further reduce the computational complexity at each SCA iteration. For the single-group case, first-order methods have been proposed to solve the convex subproblems [19], [20]. However, these methods are not directly applicable to multiple groups with inter-group interference. For the multi-group scenario with perantenna power constraints, the alternating direction method of multipliers (ADMM) [21] has been considered for solving the subproblem in each SCA iteration. In [22], zero-forcing pre-processing for interference elimination has been proposed to reduce the multi-group case to a single-group equivalent problem, for which SCA is applied. These methods provide much lower complexity than the original SCA-based method. However, since the dimension of beamforming vectors is dictated by the number of antennas, the computational complexity of these methods still grows with the number of antennas in polynomial time. This renders these methods still computationally costly for massive MIMO systems. Alternatively, for multi-cell systems, low-complexity beamforming schemes using weighted maximum ratio transmission (MRT) in combination with SDR have been developed, where only the weights, one for each user, need to be optimized, and thus the size of the optimization problem is reduced [23], [24]. Despite all the above advancements in computational algorithms, they do not optimally utilize the multicast beamforming structure.

The optimal multi-group multicast beamforming structure has been recently obtained in [2]. It is shown that the optimal solution is a weighted minimum mean-square error (MMSE) filter with an inherent low-dimensional structure for the unknown weights to be computed. With this structure, the multicast beamforming problem can be transformed into a weight optimization problem of a much lower dimension, independent

Part of this work was presented at the IEEE International Conference on Acoustics, Speech and Signal Processing, 2021 [1].

Chong Zhang and Ben Liang are with the Department of Electrical and Computer Engineering, University of Toronto, Toronto, ON M5S 1A1, Canada (e-mail: chongzhang@ece.utoronto.ca; liang@ece.utoronto.ca).

Min Dong is with the Department of Electrical, Computer and Software Engineering, Ontario Tech University, Oshawa, ON L1G 0C5, Canada (e-mail: min.dong@ontariotechu.ca).

of the number of antennas. As a result, the solution can be computed with significantly lower computational complexity, no longer growing with the number of antennas [2]. Thus, it offers design opportunities for computationally efficient algorithms for large-scale massive MIMO systems. However, [2] still adopted the conventional IPM-based SCA method for the weight optimization, whose computational complexity still does not scale well with the total number of users. Our goal in this work is to develop scalable first-order fast algorithms for large-scale systems that exploit both the optimal beamforming structure and optimization techniques.

Multicast beamforming problems considered in all the above-mentioned works typically are cast into two problem formulations: a quality-of-service (QoS) problem for transmit power minimization with signal-to-interference-and-noise (SINR) guarantees for all users, or a max-min fair (MMF) problem for maximizing the minimum SINR among all users subject to a transmit power budget. Although both types of problems are non-convex and NP-hard, the MMF problem is more complicated to solve than the QoS problem. Typically, the solution to the MMF problem is obtained through iteratively solving its inverse QoS problem along with a bi-section search over the value of minimum SINR [2], [5], [6], [17], [21]. This additional layer of bi-section iterations results in high computational complexity for the MMF problem in large-scale systems. Therefore, developing a low-complexity method to obtain a good solution for the MMF problem is also critically important.

Besides the above-mentioned works on algorithm development for multicast beamforming design, asymptotic multicast beamforming in massive MIMO systems has been analyzed in [25], [26] without inter-group interference consideration, and in [2] with inter-group interference consideration. Multicast beamforming has also been investigated for energy efficiency maximization [27] and for joint unicast and multicast transmission in massive MIMO systems [28], [29]. For overloaded systems with fewer antennas than the users, the rate-splittingbased multicast beamforming strategies have been proposed [30]–[32]. These studies address different problems from the one considered in our work.

A. Contributions

In this paper, for downlink multi-group multicast beamforming in large-scale massive MIMO systems, we propose two fast first-order algorithms for the QoS problem, which are scalable in both antenna and user dimensions. Utilizing the optimal multicast beamforming structure, we convert the original QoS problem into a non-convex weight optimization problem of a much lower dimension. We then develop two fast algorithms to solve this weight optimization problem to obtain our multicast beamforming solution. The two algorithms are SCA-based methods, referred to as the extragradient-based SCA (ESCA) and ADMM-based SCA (ASCA). They are developed using two different first-order optimization techniques. Using these algorithms, we also propose a simple closed-form scaling scheme for solving the MMF problem. The main novelty and contribution of this work are summarized below:

• We propose ESCA to solve each SCA subproblem in the dual domain. In particular, we construct a saddle point reformulation of the dual problem, which can be further

cast as a variational inequality problem for us to apply the extragradient method [33] to find the saddle point. Instead of directly considering the primal problem, our approach explores the dual problem where the extragradient method can be used efficiently, and we obtain simple closed-form gradient updates in each updating step, which is the key for our algorithm to compute the solution with low complexity. Furthermore, to avoid finding the Lipschitz constant of the gradient function to set the step size, which is generally difficult to obtain, we adopt a prediction-correction procedure for an adaptive step size in each update to ensure convergence to the saddle point for the optimal solution. We also consider two initialization methods, a fast extragradient-based initialization method and an SDR-based method, for generating favorable initial feasible points for ESCA to facilitate fast convergence.

- We also propose a fast ADMM-based algorithm for the SCA subproblems, referred to as ASCA. We transform each SCA subproblem into a favorable form to construct the ADMM procedure with two ADMM updating blocks. In particular, the structure of the transformed problem leads to simple updates in each ADMM update block, where the problem in each update block can be further decomposed into parallel subproblems of small sizes, each of which yields a closed-form solution. Thus, ASCA takes advantage of the closed-form updates in the ADMM procedure for fast computation and is guaranteed to converge to the optimal solution. The similar procedure is used to provide a ADMM-based fast initialization method to facilitate fast convergence of the algorithm.
- Taking advantage of the proposed fast algorithms for the QoS problems, we propose a simple closed-form scaling scheme to obtain a multicast beamforming solution for the MMF problem. The scheme directly scales the beamforming solution obtained from the QoS problem to meet the transmit power budget. It thus has substantially lower computational complexity than the conventional method of iteratively solving the inverse QoS problem with bi-section search. We also provide lower and upper bounds on the performance of the scaling scheme.
- Simulation results demonstrate that both ESCA and ASCA with their initialization methods provide near-optimal performance for the QoS problem. Their computational complexity is substantially lower than the state-of-the-art algorithms for large-scale systems as the number of antennas and users increases, demonstrating the algorithm scalability in large-scale systems. Between the two algorithms, ASCA is preferable with faster computation for small to moderately large systems where the number of antennas is less than 300 and the number of users per group is less than 20, while ESCA provides faster computation for larger systems. In addition, the proposed scaling scheme for the MMF problem is shown to result in only a mild loss to the optimal performance as the number of antennas becomes large but substantially faster to compute the solution.

We note that a multicast beamforming problem under perantenna power constraints has been considered in [21], where an ADMM-based algorithm has been proposed. Besides the problem being different from ours, the optimal beamforming structure is not considered there, and the iterative algorithm targets directly computing the beamforming vector. The resulting algorithm in [21] has three ADMM updating blocks, while our ASCA contains two ADMM blocks. In general, different from the existing algorithms [19]–[24], our proposed two algorithms, ESCA and ASCA, exploit both optimal multicast beamforming structure and the numerical optimization techniques, which lead to the simple closed-form updating steps to compute the multicast beamforming solution with ultra-low computational complexity for large-scale systems.

B. Organization and Notations

The rest of this paper is organized as follows. Section II presents the system model and problem formulation for multigroup multicast beamforming. Section III reviews the optimal multicast beamforming structure and the SCA method. In Sections IV and V, we present two fast algorithms, ESCA and ASCA, for the QoS problem. In Section VI, we propose a simple closed-form scaling scheme to find the solution for the MMF problem. Simulation results are provided in Section VII, and the conclusion is presented in Section VIII.

Notations: Hermitian, transpose, and conjugate are denoted as $(\cdot)^H$, $(\cdot)^T$, and $(\cdot)^*$, respectively. The real and imaginary parts of a complex number are denoted as $\Re\{\cdot\}$ and $\Im\{\cdot\}$, respectively. The Euclidean norm of a vector is denoted as $\|\cdot\|$. The notation $\mathbf{x} \succeq \mathbf{0}$ indicates element-wise non-negative. The identity matrix is denoted as I. The notation $\mathbf{z} \sim \mathcal{CN}(\mathbf{0}, \mathbf{C})$ means z is a complex Gaussian random vector with zero mean and covariance C. The non-negative real Euclidean space is denoted as \mathbb{R}_+ .

II. SYSTEM MODEL AND PROBLEM FORMULATION

We consider a downlink multi-group multicast beamforming scenario, where the base station (BS) equipped with N antennas provides service for G user groups. Each group receives a common message that is independent of the messages to other groups. Denote the set of group indices by $\mathcal{G} \triangleq \{1, \dots, G\}$. Assume that there are K_i single-antenna users in group i, and the set of user indices in the group is denoted by $\mathcal{K}_i \triangleq$ and the set of user indices in the group is denoted by $\kappa_i = \{1, \dots, K_i\}, i \in \mathcal{G}$. Define the total number of users in all groups as $K_{\text{tot}} \triangleq \sum_{i=1}^{G} K_i$. Let $\mathbf{h}_{ik} \in \mathbb{C}^{N \times 1}$ be the channel vector from the BS to user k in group i, for $k \in \mathcal{K}_i, i \in \mathcal{G}$. Let $\mathbf{w}_i \in \mathbb{C}^{N \times 1}$ be the multicast

beamforming vector for group $i \in \mathcal{G}$. The received signal at user k in group i is given by

$$y_{ik} = \mathbf{w}_i^H \mathbf{h}_{ik} s_i + \sum_{j \neq i} \mathbf{w}_j^H \mathbf{h}_{ik} s_j + n_{ik}$$

where s_i is the data symbol transmitted to group i with $E(|s_i|^2) = 1$, and n_{ik} is the receiver additive white Gaussian noise with zero mean and variance σ^2 . The received SINR at user k in group i is expressed as

$$\operatorname{SINR}_{ik} = \frac{|\mathbf{w}_i^H \mathbf{h}_{ik}|^2}{\sum_{j \neq i} |\mathbf{w}_j^H \mathbf{h}_{ik}|^2 + \sigma^2}.$$

The transmit power at the BS is given by $P_{t} \triangleq \sum_{i=1}^{G} \|\mathbf{w}_{i}\|^{2}$.

Two types of problem formulations are typically considered for the multicast beamforming design. The QoS problem is to minimize the transmit power while meeting the received SINR target at each user, which is formulated as

$$\mathcal{P}_{o}: \min_{\mathbf{w}} \sum_{i=1}^{G} \|\mathbf{w}_{i}\|^{2} \quad \text{s.t. SINR}_{ik} \ge \gamma_{ik}, \ k \in \mathcal{K}_{i}, i \in \mathcal{G} \quad (1)$$

where $\mathbf{w} \triangleq [\mathbf{w}_1^H, \cdots, \mathbf{w}_G^H]^H$, and γ_{ik} is the SINR target for user k in group i. The other problem formulation is the (weighted) MMF problem, which is to maximize the minimum (weighted) SINR among all users subject to the transmit power constraint at the BS:

$$S_o: \max_{\mathbf{w}} \min_{i,k} \frac{\mathrm{SINR}_{ik}}{\gamma_{ik}} \qquad \text{s.t.} \quad \sum_{i=1}^G \|\mathbf{w}_i\|^2 \le P \quad (2)$$

where P is the transmit power budget at the BS, and $\gamma_{ik} > 0$, $\forall i, k$, here serves as the weight to control the grade of service or fairness among users. It is well-known that both \mathcal{P}_o and S_o are non-convex and NP-hard. The existing works have proposed various computational optimization methods to find good suboptimal solutions. We will first focus on the OoS problem \mathcal{P}_o and develop two fast first-order algorithms to obtain the solutions for \mathcal{P}_o . Then, we discuss how to use the proposed algorithms to solve the MMF problem efficiently.

III. OPTIMAL MULTICAST BEAMFORMING STRUCTURE AND THE SCA METHOD

The optimal multicast beamforming structure has been obtained recently in [2]. Under this structure, problem \mathcal{P}_o is transformed into an equivalent weight optimization problem of a much lower dimension that is independent of the number of antennas. This leads to a significant computational saving and provides opportunities for efficient algorithm designs for massive MIMO systems. To facilitate our algorithm development later, we briefly describe the optimal multicast beamforming structure obtained in [2], the transformed weight optimization problem, and the SCA method for the problem.

A. Optimal Multicast Beamforming Structure

It is shown in [2] that the optimal solution to \mathcal{P}_o is a weighted MMSE filter given by

$$\mathbf{w}_i^o = \mathbf{R}^{-1}(\boldsymbol{\lambda}^o) \mathbf{H}_i \mathbf{a}_i^o, \ i \in \mathcal{G}$$
(3)

where $\mathbf{H}_i \triangleq [\mathbf{h}_{i1}, \cdots, \mathbf{h}_{iK_i}] \in \mathbb{C}^{N \times K_i}$ is the channel matrix for group i, $\mathbf{a}_i^o \in \mathbb{C}^{K_i \times 1}$ is the optimal weight vector for group i, and $\mathbf{R}(\boldsymbol{\lambda}^o) \triangleq \mathbf{I} + \sum_{i=1}^G \sum_{k=1}^{K_i} \lambda_{ik}^o \gamma_{ik} \mathbf{h}_{ik} \mathbf{h}_{ik}^H \in \mathbb{C}^{N \times N}$ is the (normalized) noise plus weighted channel covariance matrix, with $oldsymbol{\lambda}^o \in \mathbb{R}^{K_{ ext{tot}} imes 1}$ containing the optimal Lagrangian multipliers $\{\lambda_{ik}^o\}$ associated with the SINR constraints in (1). The kth weight element a_{ik}^o in \mathbf{a}_i^o indicates the relative significance of user k's channel h_{ik} in the overall group-channel direction $\mathbf{H}_i \mathbf{a}_i^o$.

To determine \mathbf{w}^{o} in (3), we need to numerically compute both λ^{o} and $\{\mathbf{a}_{i}^{o}\}$. The parameter λ^{o} can be approximately computed by the simple fixed-point iterative method proposed in [2]. Given λ , based on the optimal solution structure \mathbf{w}_i^o in (3), the original problem \mathcal{P}_o is transformed into the following equivalent weight optimization problem for $\{a_i\}$

$$\mathcal{P}_1: \min_{\mathbf{a}} \sum_{i=1}^G \|\mathbf{C}_i \mathbf{a}_i\|^2$$

s.t.
$$\frac{|\mathbf{a}_{i}^{H}\mathbf{f}_{iik}|^{2}}{\sum_{j\neq i}|\mathbf{a}_{j}^{H}\mathbf{f}_{jik}|^{2}+\sigma^{2}} \geq \gamma_{ik}, \ k \in \mathcal{K}_{i}, i \in \mathcal{G}$$
(4)

where $\mathbf{a} \triangleq [\mathbf{a}_1^H, \cdots, \mathbf{a}_G^H]^H$, $\mathbf{C}_i \triangleq \mathbf{R}^{-1}(\boldsymbol{\lambda})\mathbf{H}_i \in \mathbb{C}^{N \times K_i}$, $\mathbf{f}_{jik} \triangleq \mathbf{C}_j^H \mathbf{h}_{ik} \in \mathbb{C}^{K_j \times 1}$, $k \in \mathcal{K}_i, i, j \in \mathcal{G}$. Different from the original problem \mathcal{P}_o with GN variables, the converted problem \mathcal{P}_1 has K_{tot} variables, which does not depend on the number of antennas N. The problem dimension of \mathcal{P}_1 is much smaller than that of \mathcal{P}_o in massive MIMO systems with $K_i \ll N$. The inherent structure of the optimal multicast beamforming solution in (3) paves the way for developing low-complexity fast algorithms for multicast beamforming design in large-scale massive MIMO systems.

B. Obtaining Weights $\{a_i\}$ via SCA

The weight optimization problem \mathcal{P}_1 is still a non-convex NPhard problem. The SCA method can be adopted to solve \mathcal{P}_1 , which is guaranteed to converge to the local minimum [14]. Specifically, given any auxiliary vector $\mathbf{v}_i \in \mathbb{C}^{K_i \times 1}$, $i \in \mathcal{G}$, based on the inequality $(\mathbf{a}_i - \mathbf{v}_i)^H \mathbf{f}_{iik} \mathbf{f}_{iik}^H (\mathbf{a}_i - \mathbf{v}_i) \ge 0$, we have $|\mathbf{a}_i^H \mathbf{f}_{iik}|^2 \ge 2\mathfrak{Re}\{\mathbf{a}_i^H \mathbf{f}_{iik} \mathbf{f}_{iik}^H \mathbf{v}_i\} - |\mathbf{v}_i^H \mathbf{f}_{iik}|^2$. Replacing the numerator of the SINR expression in (4) by the right-hand side (RHS) of the above inequality, we convexify the SINR constraint and change \mathcal{P}_1 to the following convex optimization problem

$$\begin{aligned} \mathcal{P}_{1\text{SCA}}(\mathbf{v}) : & \min_{\mathbf{a}} \sum_{i=1}^{G} \|\mathbf{C}_{i}\mathbf{a}_{i}\|^{2} \\ \text{s.t. } \gamma_{ik} \sum_{j \neq i} |\mathbf{a}_{j}^{H}\mathbf{f}_{jik}|^{2} + |\mathbf{v}_{i}^{H}\mathbf{f}_{iik}|^{2} + \gamma_{ik}\sigma^{2} \\ &- 2\Re \mathfrak{e}\{\mathbf{a}_{i}^{H}\mathbf{f}_{iik}\mathbf{f}_{iik}^{H}\mathbf{v}_{i}\} \leq 0, \ k \in \mathcal{K}_{i}, i \in \mathcal{G} \end{aligned}$$

where $\mathbf{v} \triangleq [\mathbf{v}_1^H, \cdots, \mathbf{v}_G^H]^H$. Note that the solution to $\mathcal{P}_{1\text{SCA}}(\mathbf{v})$ is always feasible to \mathcal{P}_1 . By updating \mathbf{v} with the solution \mathbf{a} to $\mathcal{P}_{1\text{SCA}}(\mathbf{v})$, we iteratively solve $\mathcal{P}_{1\text{SCA}}(\mathbf{v})$ until convergence.

Since $\mathcal{P}_{1\text{SCA}}(\mathbf{v})$ is convex, it can be solved by IPM available through standard convex solvers. However, IPM is a secondorder algorithm (*i.e.*, based on the Hessian matrix of the objective function), whose best computational complexity is $O(K_{\text{tot}}^{3.5})$ and worst is $O(K_{\text{tot}}^4)$. Thus, iteratively solving $\mathcal{P}_{1\text{SCA}}(\mathbf{v})$ via IPM still incurs relatively high computational complexity for large-scale systems when K_i is large. To address this, for the rest of this paper, we develop two fast first-order algorithms to solve $\mathcal{P}_{1\text{SCA}}(\mathbf{v})$ that maintain a low complexity in computing the solution for large-scale systems.

IV. EXTRAGRADIENT-BASED SCA ALGORITHM

A. Dual Saddle Point Reformulation

For the purpose of computation, we rewrite problem $\mathcal{P}_{1SCA}(\mathbf{v})$ using the real and imaginary parts of each complex quantity. Define $\mathbf{x}_i \triangleq [\mathfrak{Re}\{\mathbf{a}_i\}^T, \mathfrak{Im}\{\mathbf{a}_i\}^T]^T, \mathbf{y}_i \triangleq [\mathfrak{Re}\{\mathbf{v}_i\}^T, \mathfrak{Im}\{\mathbf{v}_i\}^T]^T$. Also, define

$$\mathbf{A}_{i} \triangleq \begin{bmatrix} \mathfrak{Re}\{\mathbf{C}_{i}\} - \mathfrak{Im}\{\mathbf{C}_{i}\}\\ \mathfrak{Im}\{\mathbf{C}_{i}\} & \mathfrak{Re}\{\mathbf{C}_{i}\} \end{bmatrix}, \widetilde{\mathbf{F}}_{jik} \triangleq \begin{bmatrix} \mathfrak{Re}\{\mathbf{f}_{jik}\} - \mathfrak{Im}\{\mathbf{f}_{jik}\}\\ \mathfrak{Im}\{\mathbf{f}_{jik}\} & \mathfrak{Re}\{\mathbf{f}_{jik}\} \end{bmatrix}$$

for $k \in \mathcal{K}_i, i, j \in \mathcal{G}$. Then, we have $\|\mathbf{C}_i \mathbf{a}_i\|^2 = \|\mathbf{A}_i \mathbf{x}_i\|^2$ and $|\mathbf{a}_j^H \mathbf{f}_{jik}|^2 = \|\mathbf{x}_j^T \widetilde{\mathbf{F}}_{jik}\|^2 = \mathbf{x}_j^T \mathbf{F}_{jik} \mathbf{x}_j$, where $\mathbf{F}_{jik} \triangleq \widetilde{\mathbf{F}}_{jik} \widetilde{\mathbf{F}}_{jik}^T$, for $k \in \mathcal{K}_i$ and $j, i \in \mathcal{G}$. Define $\mathbf{x} \triangleq [\mathbf{x}_1^T, \cdots, \mathbf{x}_G^T]^T$ and

 $\mathbf{y} \triangleq [\mathbf{y}_1^T, \cdots, \mathbf{y}_G^T]^T$. With these new variables, $\mathcal{P}_{1\text{SCA}}(\mathbf{v})$ can be equivalently expressed as

$$\mathcal{P}_{1\text{SCA}}^{r}(\mathbf{y}): \min_{\mathbf{x}} \sum_{i=1}^{G} \|\mathbf{A}_{i}\mathbf{x}_{i}\|^{2}$$

s.t. $\mathbf{y}_{i}^{T}\mathbf{F}_{iik}\mathbf{y}_{i} + \gamma_{ik}\sum_{j\neq i} \mathbf{x}_{j}^{T}\mathbf{F}_{jik}\mathbf{x}_{j} - 2\mathbf{y}_{i}^{T}\mathbf{F}_{iik}\mathbf{x}_{i} + \gamma_{ik}\sigma^{2} \leq 0,$
 $k \in \mathcal{K}_{i}, i \in \mathcal{G}.$ (5)

We first describe the class of variational inequality problems [34] below.

Definition 1 (Variational Inequality). Given $\mathcal{Z} \subseteq \mathbb{R}^n$ and a mapping $\psi : \mathcal{Z} \to \mathbb{R}^n$, the variational inequality is to find $\mathbf{z} \in \mathcal{Z}$ satisfying $\psi(\mathbf{z})^T(\mathbf{z}'-\mathbf{z}) \ge 0, \forall \mathbf{z}' \in \mathcal{Z}$. Operator $\psi(\cdot)$ is said to be monotone on \mathcal{Z} if $[\psi(\mathbf{z}) - \psi(\mathbf{z}')]^T(\mathbf{z} - \mathbf{z}') \ge 0, \forall \mathbf{z}, \mathbf{z}' \in \mathcal{Z}$. The problem is monotone if operator $\psi(\cdot)$ is monotone.

The projection methods belong to a class of iterative algorithms that solve the monotone variational inequality problems [34]. At each iteration, a projection method uses the updating step to compute the point (*i.e.*, the value of the optimization variable) and then projects it onto the feasible set of the problem to ensure the updated point is feasible. Note that the projection method may not be an efficient method. In general, the projection methods are only computationally cheap when the projection is easy to compute.

Note that $\mathcal{P}_{1SCA}^{r}(\mathbf{y})$ is convex, and the objective function is differentiable. Let operator $\psi(\mathbf{x})$ be the gradient of the objective function of $\mathcal{P}_{1SCA}^{r}(\mathbf{y})$. Following the optimality criterion for a convex optimization problem, $\psi(\mathbf{x})$ is monotone, and thus $\mathcal{P}_{1SCA}^{r}(\mathbf{y})$ is a monotone variational inequality problem. However, it is difficult to find a closed-form expression for the projection onto the feasible set of $\mathcal{P}_{1SCA}^{r}(\mathbf{y})$. Thus, directly applying the projection method to solve $\mathcal{P}_{1SCA}^{r}(\mathbf{y})$ is not computationally attractive. To overcome this difficulty, we resort to the Lagrange dual problem of $\mathcal{P}_{1SCA}^{r}(\mathbf{y})$.

Define $\phi_{ik}(\mathbf{x}) \triangleq \mathbf{y}_i^T \mathbf{F}_{iik} \mathbf{y}_i + \gamma_{ik} \sum_{j \neq i} \mathbf{x}_j^T \mathbf{F}_{jik} \mathbf{x}_j - 2\mathbf{y}_i^T \mathbf{F}_{iik} \mathbf{x}_i + \gamma_{ik} \sigma^2$ and $\phi_i(\mathbf{x}) \triangleq [\phi_{i1}(\mathbf{x}), \dots, \phi_{iK_i}(\mathbf{x})]^T$, for $k \in \mathcal{K}_i, i \in \mathcal{G}$. The Lagrangian of $\mathcal{P}_{1SCA}^{r}(\mathbf{y})$ is given by

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\eta}) = \sum_{i=1}^{G} \left(\|\mathbf{A}_i \mathbf{x}_i\|^2 + \boldsymbol{\eta}_i^T \boldsymbol{\phi}_i(\mathbf{x}) \right)$$
(6)

where $\eta_{ik} \geq 0$ is the dual variable associated with constraint (5), and $\boldsymbol{\eta} \triangleq \left[\boldsymbol{\eta}_{1}^{T}, \cdots, \boldsymbol{\eta}_{G}^{T}\right]^{T}$ with $\boldsymbol{\eta}_{i} \triangleq \left[\eta_{i1}, \cdots, \eta_{iK_{i}}\right]^{T}$. The Lagrange dual problem of $\mathcal{P}_{1\text{SCA}}^{r}(\mathbf{y})$ is given by

$$\mathcal{D}_{1 ext{SCA}}^{ ext{r}}(\mathbf{y}): \max_{oldsymbol{\eta} \succcurlyeq oldsymbol{0}} \min_{\mathbf{x}} \sum_{i=1}^{G} \left(\|\mathbf{A}_i \mathbf{x}_i\|^2 + oldsymbol{\eta}_i^T \phi_i(\mathbf{x})
ight).$$

Let \mathbf{x}^{o} and $\boldsymbol{\eta}^{o}$ be the primal and dual optimal solutions for $\mathcal{P}_{1\text{SCA}}^{r}(\mathbf{y})$. Since $\mathcal{P}_{1\text{SCA}}^{r}(\mathbf{y})$ is convex and Slater's condition holds, the strong duality holds. It follows that $\mathbf{u}^{o} \triangleq (\mathbf{x}^{o}, \boldsymbol{\eta}^{o})$ is a saddle point of the Lagrangian $\mathcal{L}(\mathbf{x}, \boldsymbol{\eta})$ [18]. Define operator $g(\mathbf{u})$ as

$$g(\mathbf{u}) = g(\mathbf{x}, \boldsymbol{\eta}) \triangleq \begin{bmatrix} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\eta}) \\ -\nabla_{\boldsymbol{\eta}} \mathcal{L}(\mathbf{x}, \boldsymbol{\eta}) \end{bmatrix}, \quad \mathbf{u} \in \mathcal{U}$$
(7)

where $\mathcal{U} \triangleq \mathbb{R}^{2K_{\text{tot}}} \times \mathbb{R}^{K_{\text{tot}}}_+$ is a closed convex set. Then, $\mathcal{D}^{r}_{1\text{SCA}}(\mathbf{y})$ can be interpreted as finding the saddle point \mathbf{u}^o . It is shown in [34] that the problem of finding the saddle point \mathbf{u}^o can be cast

as the variational inequality problem: Find $\mathbf{u}^o \in \mathcal{U}$ that satisfies the following variational inequality

$$g(\mathbf{u}^o)^T(\mathbf{u} - \mathbf{u}^o) \ge 0, \quad \forall \ \mathbf{u} \in \mathcal{U}.$$
 (8)

B. Extragradient-Based SCA

To solve problem (8), one may consider the basic projection algorithm (BPA) [34], which iteratively updates **u** using $g(\mathbf{u})$ and then projects it onto \mathcal{U} . However, the convergence of BPA requires operator $g(\mathbf{u})$ to be strongly monotone [34]. Following Definition 1, operator $\psi(\cdot)$ is said to be strongly monotone on \mathcal{Z} , if there exists a constant c > 0 such that $[\psi(\mathbf{z}) - \psi(\mathbf{z}')]^T (\mathbf{z} - \mathbf{z}')$ $\geq c ||\mathbf{z} - \mathbf{z}'||^2, \forall \mathbf{z}, \mathbf{z}' \in \mathcal{Z}$. Since $\mathcal{L}(\mathbf{x}, \boldsymbol{\eta})$ is linear with respect to (w.r.t.) $\boldsymbol{\eta}$, operator $g(\mathbf{u})$ is not strongly monotone on \mathcal{U} . Thus, we cannot apply BPA to our problem due to no convergence guarantee.

Instead of BPA, in this work, we adopt the extragradient method, a variant of BPA, proposed by Korpelevich in [33]. Compared with BPA, the extragradient method can guarantee convergence for a monotone operator, at the cost of an extra update-and-projection step at each iteration. For operator $g(\mathbf{u})$ in (7), since $\mathcal{L}(\mathbf{x}, \boldsymbol{\eta})$ is convex in $\mathbf{x} \in \mathbb{R}^{2K_{\text{tot}}}$ and $-\mathcal{L}(\mathbf{x}, \boldsymbol{\eta})$ is convex in $\boldsymbol{\eta} \in \mathbb{R}^{K_{\text{tot}}}_+$, $g(\mathbf{u})$ is monotone on \mathcal{U} by Definition 1. Thus, we develop a fast iterative algorithm to solve $\mathcal{P}^{\text{r}}_{1\text{SCA}}(\mathbf{y})$ by applying the extragradient algorithm in the problem dual domain.

The updating procedure of the extragradient algorithm for finding the saddle point is summarized as follows [33], [34]: At iteration n + 1,

$$\bar{\mathbf{x}}^{(n)} = \mathbf{x}^{(n)} - \alpha \nabla_{\mathbf{x}^{(n)}} \mathcal{L}(\mathbf{x}^{(n)}, \boldsymbol{\eta}^{(n)}), \tag{9}$$

$$\bar{\boldsymbol{\eta}}^{(n)} = \left[\boldsymbol{\eta}^{(n)} + \alpha \nabla_{\boldsymbol{\eta}^{(n)}} \mathcal{L}(\mathbf{x}^{(n)}, \boldsymbol{\eta}^{(n)})\right]^+, \qquad (10)$$

$$\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} - \alpha \nabla_{\bar{\mathbf{x}}^{(n)}} \mathcal{L}(\bar{\mathbf{x}}^{(n)}, \bar{\boldsymbol{\eta}}^{(n)}), \qquad (11)$$

$$\boldsymbol{\eta}^{(n+1)} = \left[\boldsymbol{\eta}^{(n)} + \alpha \nabla_{\bar{\boldsymbol{\eta}}^{(n)}} \mathcal{L}(\bar{\mathbf{x}}^{(n)}, \bar{\boldsymbol{\eta}}^{(n)})\right]^+$$
(12)

where $\bar{\mathbf{x}}^{(n)}$ and $\bar{\boldsymbol{\eta}}^{(n)}$ are the intermediate updates in the extra update-and-projection step in (9)(10), α is the step size, and notation $[\mathbf{z}]^+ \triangleq [[z_1]^+, \dots, [z_n]^+]^T$ with $[z_i]^+ \triangleq \max\{z_i, 0\}$, for $\mathbf{z} \in \mathbb{R}^n$.

The gradient $\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\eta})$ can be denoted as $\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\eta}) = [\nabla_{\mathbf{x}_1} \mathcal{L}(\mathbf{x}, \boldsymbol{\eta})^T, \cdots, \nabla_{\mathbf{x}_G} \mathcal{L}(\mathbf{x}, \boldsymbol{\eta})^T]^T$. From (6), we obtain

$$\nabla_{\mathbf{x}_{i}} \mathcal{L}(\mathbf{x}, \boldsymbol{\eta}) = 2\mathbf{A}_{i}^{T} \mathbf{A}_{i} \mathbf{x}_{i} + 2 \left(\sum_{j \neq i} \sum_{k=1}^{K_{j}} \gamma_{jk} \eta_{jk} \mathbf{F}_{ijk} \right) \mathbf{x}_{i} - 2 \left(\sum_{k=1}^{K_{i}} \eta_{ik} \mathbf{F}_{iik} \right)^{T} \mathbf{y}_{i}, \quad i \in \mathcal{G}.$$
(13)

Also from (6), we obtain the gradient $\nabla_{\eta} \mathcal{L}(\mathbf{x}, \eta)$ as

$$\nabla_{\boldsymbol{\eta}} \mathcal{L}(\mathbf{x}, \boldsymbol{\eta}) = [\boldsymbol{\phi}_1^T(\mathbf{x}), \cdots, \boldsymbol{\phi}_G^T(\mathbf{x})]^T.$$
(14)

Substituting the expressions in (13) and (14) into (9)–(12), we obtain the closed-form updates for $\mathbf{x}^{(n+1)}$ and $\boldsymbol{\eta}^{(n+1)}$.

For a monotone variational inequality problem with operator being L-Lipschitz continuous, the extragradient algorithm is guaranteed to converge to the optimal solution, provided that the step size satisfies $\alpha < 1/L$ [34]. Unfortunately, it is difficult to determine Lipschitz constant L for $g(\mathbf{u})$ in our problem. To overcome this difficulty, instead of a constant step size α , we adopt an adaptive strategy based on the prediction-correction procedure [35], [36] to adaptively set the step size $\alpha^{(n)}$ for each iteration.

Given a fixed step size α , the prediction-correction procedure contains two steps at iteration n + 1:

- 1) Prediction: Obtain $\bar{\mathbf{u}}^{(n)} \triangleq (\bar{\mathbf{x}}^{(n)}, \bar{\boldsymbol{\eta}}^{(n)})$ from (9) and (10) with fixed α . Compute $d_{\mathbf{u}}^{(n)} \triangleq \|\bar{\mathbf{u}}^{(n)} \mathbf{u}^{(n)}\|$ and $d_g^{(n)} \triangleq \|g(\bar{\mathbf{u}}^{(n)}) g(\mathbf{u}^{(n)})\|$. Compute step size $\hat{\alpha}^{(n)} = c \frac{d_{\mathbf{u}}^{(n)}}{d_g^{(n)}}$, where $c \in (0, 1)$ is a constant.
- 2) Correction: Set $\alpha^{(n)} = \min\{\alpha, \hat{\alpha}^{(n)}\}\$ for iteration n+1 to update $\mathbf{x}^{(n+1)}$ and $\boldsymbol{\eta}^{(n+1)}$ in (9)–(12).

We now show that the prediction-correction procedure guarantees the extragradient algorithm converges to the saddle point \mathbf{u}^{o} of problem (8). If we replace α by $\alpha^{(n)}$ in (9)–(12) of the extragradient algorithm, then for any step size sequence $\{\alpha^{(n)}\}$, the following holds [35], [36]

$$\|\mathbf{u}^{(n+1)} - \mathbf{u}^{o}\|^{2} \leq \|\mathbf{u}^{(n)} - \mathbf{u}^{o}\|^{2} - \|\bar{\mathbf{u}}^{(n)} - \mathbf{u}^{(n)}\|^{2} \left(1 - \left(\alpha^{(n)} \frac{d_{g}^{(n)}}{d_{\mathbf{u}}^{(n)}}\right)^{2}\right).$$
(15)

If we set $\alpha^{(n)}$ as in the correction step of the above predictioncorrection procedure, then $\|\mathbf{u}^{(n+1)} - \mathbf{u}^o\| < \|\mathbf{u}^{(n)} - \mathbf{u}^o\|$ and $\{\mathbf{u}^{(n)}\}$ move towards the saddle point \mathbf{u}^o of problem (8). It follows that the algorithm converges to the optimal solution to $\mathcal{P}_{1SCA}^r(\mathbf{y})$ in each SCA iteration.

Overall, the ESCA algorithm to solve \mathcal{P}_1 is summarized in Algorithm 1. The main computational complexity of ESCA lies in computing the gradients $\nabla_{\mathbf{x}}\mathcal{L}(\mathbf{x}, \boldsymbol{\eta})$ using (13) and $\nabla_{\boldsymbol{\eta}}\mathcal{L}(\mathbf{x}, \boldsymbol{\eta})$ in (14). At each extragradient iteration, the related matrix-vector computation for $\nabla_{\mathbf{x}}\mathcal{L}(\mathbf{x}, \boldsymbol{\eta})$ requires $\sum_{i=1}^{G} \left(\sum_{j=1}^{G} 4K_jK_i^2 + 32K_i^2 \right)$ flops, and that for $\nabla_{\boldsymbol{\eta}}\mathcal{L}(\mathbf{x}, \boldsymbol{\eta})$ requires $\sum_{i=1}^{G} \left(\sum_{j=1}^{G} (16K_j^2 + 8K_j)K_i - 8K_i^2 \right)$ flops. Note that ESCA consists of two layers of iterations: the outer-layer SCA and the inner-layer extragradient-based algorithm to solve each $\mathcal{P}_{1\text{SCA}}(\mathbf{v})$. The number of iterations at the two layers for convergence depend on the system parameters N and K_{tot} . From our experiments, for $N = 100 \sim 500$ antennas and $K_{\text{tot}} = 15 \sim 60$ users, it typically takes $20 \sim 200$ iterations for the inner-layer extragradient algorithm to converge at each SCA iteration, and the outer layer typically takes $2 \sim 60$ iterations to converge.

Remark 1. As mentioned earlier, BPA consists of only one update-and-projection step at each iteration (similar to (9)(10)), but requires the operator $g(\mathbf{u})$ to be strongly monotone for convergence. If $q(\mathbf{u})$ is only monotone, the updated point (the same as $\bar{\mathbf{u}}^{(n)}$ in (9)(10) in iteration n+1) may move away from, instead of closer to, the optimal point \mathbf{u}^{o} at each iteration. Thus, the updating procedure may not converge over iterations. In contrast, the extragradient algorithm adds an extra updateand-projection step as in (11)(12) at each iteration. This extra step can ensure convergence for a monotone operator $g(\mathbf{u})$. Specifically, it is shown in [34] that in iteration n + 1, after obtaining $\mathbf{\bar{u}}^{(n)}$ at the first update-and-projection step in (9)(10), a hyperplane $\mathcal{H}^{(n)} \triangleq \{\mathbf{u} | g(\mathbf{\bar{u}}^{(n)})^T (\mathbf{u} - \mathbf{\bar{u}}^{(n)}) = 0\}$ can be constructed at $\bar{\mathbf{u}}^{(n)}$ with normal vector $q(\bar{\mathbf{u}}^{(n)})$. For the monotone operator $q(\mathbf{u})$, it is proven that this hyperplane $\mathcal{H}^{(n)}$ separates the point $\mathbf{u}^{(n)}$ and the optimal point \mathbf{u}^{o} , where \mathbf{u}^{o} is in the halfspace in the direction $-g(\bar{\mathbf{u}}^{(n)})$. The second updateand-projection step in (11)(12) then utilizes $-q(\bar{\mathbf{u}}^{(n)})$ to update

Algorithm 1 ESCA Algorithm to Solve \mathcal{P}_1

1: Initialization: Set feasible initial point $\mathbf{y}^{(0)}$. Set α and c. Set l = 0. 2: repeat // solve $\mathcal{P}_{1\text{SCA}}^{r}(\mathbf{y}^{(l)})$ Initialization: $\mathbf{x}^{(0)} = \mathbf{y}^{(l)}, \ \boldsymbol{\eta}^{(0)} = \mathbf{0}, \ n = 0.$ 3: 4: 5: repeat Compute $\bar{\mathbf{x}}^{(n)}$ and $\bar{\boldsymbol{\eta}}^{(n)}$ by (9)(10) using α . 6: Compute $g(\bar{\mathbf{u}}^{(n)})$ by (13)(14). Compute $d_{\mathbf{u}}^{(n)}$ and $d_{g}^{(n)}$. 7: 8: Prediction: Compute $\hat{\alpha}^{(n)} = c d_{\mathbf{u}}^{(n)} / d_{g}^{(n)}$. 9: *Correction:* Update step size $\alpha^{(n)} = \min{\{\alpha, \hat{\alpha}^{(n)}\}}$. 10: if $\alpha^{(n)} = \alpha$ then 11: Update $x^{(n+1)}$ and $\eta^{(n+1)}$ by (11)(12). 12: else 13: Update $\mathbf{x}^{(n+1)}$ and $\boldsymbol{\eta}^{(n+1)}$ by (9)–(12) using $\alpha^{(n)}$. 14: end if 15: $n \leftarrow n+1.$ 16: until convergence 17: Set $\mathbf{y}^{(l+1)} = \mathbf{x}^{(n)}$. Set $l \leftarrow l+1$. 18: 19: until convergence

 $\mathbf{u}^{(n+1)}$. This ensures that $\mathbf{u}^{(n+1)}$ move towards the optimal point \mathbf{u}^{o} and thus is closer to \mathbf{u}^{o} than $\mathbf{u}^{(n)}$ is.

C. Initialization for ESCA

A challenge in using SCA to solve \mathcal{P}_1 is that it requires a feasible initial point satisfying the SINR constraint (4) as $\mathbf{v}^{(0)}$ for $\mathcal{P}_{1\text{SCA}}(\mathbf{v})$ (equivalently $\mathbf{y}^{(0)}$ for $\mathcal{P}_{1\text{SCA}}^r(\mathbf{y})$). It is necessary to generate a feasible initial point with a low computational complexity. Furthermore, a good initial point closer to the (locally) optimal point of \mathcal{P}_1 could accelerate the convergence. Below, we consider two initialization methods for ESCA.

1) Extragradient-based initialization method (EIM): Based on the extragradient-based algorithm in Section IV-B, we propose EIM as follows. EIM generates a feasible point by solving the following feasibility problem

$$\mathcal{P}_{1\text{fea}} : \text{Find } \{\mathbf{x}\}$$

s.t.
$$\frac{\mathbf{x}_{i}^{T}\mathbf{F}_{iik}\mathbf{x}_{i}}{\sum_{j\neq i}\mathbf{x}_{j}^{T}\mathbf{F}_{jik}\mathbf{x}_{j} + \sigma^{2}} \geq \gamma_{ik}, \ k \in \mathcal{K}_{i}, i \in \mathcal{G} \quad (16)$$

where \mathbf{x}_i and \mathbf{F}_{jik} are defined at the beginning of Section IV-A, and constraint (16) is an equivalent real representation of constraint (4) in \mathcal{P}_1 based on $|\mathbf{a}_j^H \mathbf{f}_{jik}|^2 = \mathbf{x}_j^T \mathbf{F}_{jik} \mathbf{x}_j$.

We solve $\mathcal{P}_{1\text{fea}}$ by applying the extragradient method with the adaptive step size proposed in Section IV-B. Specifically, the Lagrangian of $\mathcal{P}_{1\text{fea}}$ is given by $\widetilde{\mathcal{L}}(\mathbf{x}, \widetilde{\boldsymbol{\eta}}) = \sum_{i=1}^{G} \widetilde{\phi}_{i}^{T}(\mathbf{x}) \widetilde{\eta}_{i}$, where $\widetilde{\eta}_{ik}$ is the dual variable for constraint (16), $\widetilde{\boldsymbol{\eta}} \triangleq [\widetilde{\boldsymbol{\eta}}_{1}^{T}, \cdots, \widetilde{\boldsymbol{\eta}}_{G}^{T}]^{T}$ with $\widetilde{\boldsymbol{\eta}}_{i} \triangleq [\widetilde{\eta}_{i1}, \cdots, \widetilde{\eta}_{iK_{i}}]^{T}$, and $\widetilde{\phi}_{i}(\mathbf{x}) \triangleq [\widetilde{\phi}_{i1}(\mathbf{x}), \dots, \widetilde{\phi}_{iK_{i}}(\mathbf{x})]^{T}$ with $\widetilde{\phi}_{ik}(\mathbf{x}) \triangleq \gamma_{ik} \sum_{j \neq i} \mathbf{x}_{j}^{T} \mathbf{F}_{jik} \mathbf{x}_{j} - \mathbf{x}_{i}^{T} \mathbf{F}_{iik} \mathbf{x}_{i} + \gamma_{ik} \sigma^{2}$, for $k \in \mathcal{K}_{i}, i \in \mathcal{G}$. The gradient $\nabla_{\mathbf{x}_{i}} \widetilde{\mathcal{L}}(\mathbf{x}, \widetilde{\boldsymbol{\eta}})$, for $i \in \mathcal{G}$, is given by

$$\nabla_{\mathbf{x}_i} \widetilde{\mathcal{L}}(\mathbf{x}, \widetilde{\boldsymbol{\eta}}) = 2 \left(\sum_{j \neq i} \sum_{k=1}^{K_j} \gamma_{jk} \widetilde{\eta}_{jk} \mathbf{F}_{ijk} \right)^T \mathbf{x}_i - 2 \left(\sum_{k=1}^{K_i} \widetilde{\eta}_{ik} \mathbf{F}_{iik} \right)^T \mathbf{x}_i$$

Similar to (14), the gradient $\nabla_{\tilde{\eta}} \tilde{\mathcal{L}}(\mathbf{x}, \tilde{\eta})$ is given by $\nabla_{\tilde{\eta}} \tilde{\mathcal{L}}(\mathbf{x}, \tilde{\eta}) = [\tilde{\phi}_1^T(\mathbf{x}), \cdots, \tilde{\phi}_G^T(\mathbf{x})]^T$. The updating procedure of the extragradient method in (9)–(12) is then applied for solving $\mathcal{P}_{1\text{fea}}$, with η , $\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \eta)$, and $\nabla_{\eta} \mathcal{L}(\mathbf{x}, \eta)$ being replaced by $\tilde{\eta}$, $\nabla_{\mathbf{x}} \tilde{\mathcal{L}}(\mathbf{x}, \tilde{\eta})$, and $\nabla_{\tilde{\eta}} \tilde{\mathcal{L}}(\mathbf{x}, \tilde{\eta})$, respectively. Furthermore, the prediction-correction procedure in Section IV-B is used to set the adaptive step size at each iteration.

For simplicity, we randomly chose the initial point for EIM. Note that this point may not be feasible for SINR constraint (4) in \mathcal{P}_1 . Also, since $\mathcal{P}_{1\text{fea}}$ is a non-convex optimization problem, EIM may not be guaranteed to converge or the terminating point may not be feasible as required for $\mathbf{v}^{(0)}$ for $\mathcal{P}_{1\text{SCA}}(\mathbf{v})$. Thus, EIM is served as a heuristic algorithm. If the point produced by EIM is infeasible, we may consider different initial points for EIM until a feasible is obtained. Our extensive simulation experiments show that EIM converges and provides a feasible initial point with a very high probability.

2) SDR-based initialization method: Similar to [2], we can apply SDR along with Gaussian randomization to \mathcal{P}_1 to obtain an approximate solution as a feasible initial point $\mathbf{y}^{(0)}$ to $\mathcal{P}_{1SCA}^r(\mathbf{y})$ for ESCA. Recall that \mathcal{P}_1 is converted from the original problem \mathcal{P}_o and is of relatively smaller size (K_{tot} variables and constraints). SDR can provide a good initial point when the problem is of small to moderate size, leading to fast convergence for ESCA. However, as K_{tot} becomes large, the computational complexity of SDR will increase significantly, and the quality of the initial point it computes deteriorates. Thus, SDR for initialization is only suitable for small to moderate problems.

V. ADMM-BASED SCA ALGORITHM

In this section, we develop another computationally efficient algorithm based on ADMM [37], which is a different optimization technique from ESCA. ADMM has drawn growing popularity in recent years as a robust and fast numerical method for solving large-scale problems. We propose an ADMM-based algorithm to solve each SCA subproblem $\mathcal{P}_{1SCA}(\mathbf{v})$. Define the auxiliary variables $d_{jik} \triangleq \mathbf{a}_{j}^{H}\mathbf{f}_{jik}$, for $k \in \mathcal{K}_{i}, i, j \in \mathcal{G}$. Define d $\triangleq \left[\mathbf{d}_{11}^{H}, \cdots, \mathbf{d}_{GK_{G}}^{H}\right]^{H} \in \mathbb{C}^{GK_{tot}}$, with $\mathbf{d}_{ik} \triangleq \left[d_{1ik}, \cdots, d_{Gik}\right]^{T} \in \mathbb{C}^{G}$. Then, the problem $\mathcal{P}_{1SCA}(\mathbf{v})$ can be equivalently expressed as

$$\mathcal{P}_{1\text{ADMM}}(\mathbf{v}) : \min_{\mathbf{a},\mathbf{d}} \sum_{j=1}^{G} \|\mathbf{C}_{j}\mathbf{a}_{j}\|^{2}$$

s.t. $d_{jik} = \mathbf{a}_{j}^{H}\mathbf{f}_{jik}, \ k \in \mathcal{K}_{i}, i, j \in \mathcal{G}$ (17)
 $\gamma_{ik} \sum_{j \neq i} |d_{jik}|^{2} + |\mathbf{v}_{i}^{H}\mathbf{f}_{iik}|^{2} + \gamma_{ik}\sigma^{2} - 2\mathfrak{Re}\{d_{iik}\mathbf{f}_{iik}^{H}\mathbf{v}_{i}\} \leq 0,$
 $k \in \mathcal{K}_{i}, i \in \mathcal{G}.$ (18)

Denote the feasible set for the inequality constraint (18) by \mathcal{F} . Define the indicator function for the set \mathcal{F} as

$$I_{\mathcal{F}}(\mathbf{d}) \triangleq \begin{cases} 0 & \text{if } \mathbf{d} \in \mathcal{F}, \\ \infty & \text{otherwise.} \end{cases}$$

Then, $\mathcal{P}_{1ADMM}(\mathbf{v})$ is equivalent to the following problem

$$\mathcal{P}'_{\text{IADMM}}(\mathbf{v}) : \min_{\mathbf{a}, \mathbf{d}} \sum_{j=1}^{G} \|\mathbf{C}_{j}\mathbf{a}_{j}\|^{2} + I_{\mathcal{F}}(\mathbf{d})$$

s.t. $d_{jik} = \mathbf{a}_{j}^{H}\mathbf{f}_{jik}, \ k \in \mathcal{K}_{i}, i, j \in \mathcal{G}.$ (19)

By introducing the auxiliary vector **d** in constraint (18), we construct the equality-constrained problem $\mathcal{P}'_{1ADMM}(\mathbf{v})$, where the objective function contains two separate terms for **d** and **a** only. This allows us to apply ADMM to decompose the problem into separate subproblems [37]. The augmented Lagrangian of $\mathcal{P}'_{1ADMM}(\mathbf{v})$ is given by

$$\mathcal{L}_{\rho}(\mathbf{d}, \mathbf{a}, \mathbf{q}) = \sum_{j=1}^{G} \|\mathbf{C}_{j} \mathbf{a}_{j}\|^{2} + I_{\mathcal{F}}(\mathbf{d}) + \frac{\rho}{2} \sum_{j=1}^{G} \sum_{i=1}^{G} \sum_{k=1}^{K_{i}} |d_{jik} - \mathbf{a}_{j}^{H} \mathbf{f}_{jik} + q_{jik}|^{2} \quad (20)$$

where $\rho > 0$ is the penalty parameter, and $q_{jik} \in \mathbb{C}$ is the dual variable associated with constraint (19), and $\mathbf{q} \triangleq [\mathbf{q}_{11}^H, \cdots, \mathbf{q}_{GK_G}^H]^H$, with $\mathbf{q}_{ik} \triangleq [q_{1ik}, \cdots, q_{Gik}]^T$. We now decompose $\mathcal{L}_{\rho}(\mathbf{d}, \mathbf{a}, \mathbf{q})$ into two subproblems for \mathbf{d} and \mathbf{a} separately, and update $\{\mathbf{d}, \mathbf{a}, \mathbf{q}\}$ alternatively. Our proposed ADMM-based updating procedure for solving $\mathcal{P}'_{\text{IADMM}}(\mathbf{v})$ is given as follows:

$$\mathbf{d}^{(n+1)} = \arg\min_{\mathbf{d}} \mathcal{L}_{\rho}(\mathbf{d}, \mathbf{a}^{(n)}, \mathbf{q}^{(n)}), \tag{21}$$

$$\mathbf{a}^{(n+1)} = \operatorname*{arg\,min}_{\mathbf{a}} \mathcal{L}_{\rho}(\mathbf{d}^{(n+1)}, \mathbf{a}, \mathbf{q}^{(n)}), \tag{22}$$

$$q_{jik}^{(n+1)} = q_{jik}^{(n)} + \left(d_{jik}^{(n+1)} - \mathbf{a}_j^{(n+1)H} \mathbf{f}_{jik} \right)$$
(23)

where *n* is the iteration index. Since $\mathcal{P}'_{1\text{ADMM}}(\mathbf{v})$ is convex, the above ADMM procedure is guaranteed to converge to the optimal solution of $\mathcal{P}_{1\text{ADMM}}(\mathbf{v})$ [37].

The two main updating steps in (21) and (22) involve solving two optimization problems. In the following, we derive the closed-form solutions for the two optimization problems in (21) and (22), which leads to fast computation at each iteration.

A. Closed-Form d-Update

Given $\mathbf{a}^{(n)}$ and $\mathbf{q}^{(n)}$, from (20), the update of d in (21) is equivalent to solving the following problem

$$\mathcal{P}_{\mathsf{d}}(\mathbf{v}) : \min_{\mathbf{d}} \sum_{j=1}^{G} \sum_{i=1}^{G} \sum_{k=1}^{K_{i}} |d_{jik} - \mathbf{a}_{j}^{(n)H} \mathbf{f}_{jik} + q_{jik}^{(n)}|^{2} \quad \text{s.t. (18).}$$

Problem $\mathcal{P}_{d}(\mathbf{v})$ can be decomposed into K_{tot} convex subproblems, one for each user k in group i given by

$$\mathcal{P}_{dsub}(\mathbf{v}) : \min_{\mathbf{d}_{ik}} \sum_{j=1}^{G} |d_{jik} - e_{1,jik}^{(n)}|^2$$

s.t. $e_{2,ik} + \gamma_{ik} \sum_{j \neq i} |d_{jik}|^2 - 2\Re \mathfrak{e}\{d_{iik}e_{3,ik}\} \le 0$ (24)

where

$$e_{1,jik}^{(n)} \triangleq \mathbf{a}_{j}^{(n)H} \mathbf{f}_{jik} - q_{jik}^{(n)}, \ e_{2,ik} \triangleq |\mathbf{v}_{i}^{H} \mathbf{f}_{iik}|^{2} + \gamma_{ik} \sigma^{2}, \ e_{3,ik} \triangleq \mathbf{f}_{iik}^{H} \mathbf{v}_{i}.$$
(25)

Problem $\mathcal{P}_{dsub}(\mathbf{v})$ is a convex QCQP-1 problem, whose solution may be derived in closed-form. In particular, a problem of similar structure has been considered in [21], where the closedform solution is derived in [21, Appendix A]. We directly use this result and state the closed-form solution below. The optimal solution \mathbf{d}_{ik}^{o} for $\mathcal{P}_{dsub}(\mathbf{v})$ is given by

$$d_{jik}^{o} = \begin{cases} e_{1,iik}^{(n)} + \nu_{ik}^{o} e_{3,ik}^{*} & \text{if } j = i, \\ \frac{e_{1,jik}^{(n)}}{1 + \nu_{ik}^{o} \gamma_{ik}} & \text{otherwise} \end{cases}$$
(26)

where $\nu_{ik}^o = 0$ if $e_{2,ik} + \gamma_{ik} \sum_{j \neq i} |e_{1,jik}^{(n)}|^2 - 2\Re \mathfrak{e}\{e_{3,ik}e_{1,ik}^{(n)}\} \le 0$; otherwise, ν_{ik}^o is the unique real positive root of the following cubic equation of ν_{ik} :

$$e_{2,ik} + \gamma_{ik} \frac{\sum_{j \neq i} |e_{1,jik}^{(n)}|^2}{(1 + \nu_{ik}\gamma_{ik})^2} - 2\Re \mathfrak{e}\{e_{3,ik}e_{1,iik}^{(n)}\} - 2\nu_{ik}|e_{3,ik}|^2 = 0.$$
(27)

Since the roots of (27) are given by the cubic formula, ν_{ik}^{o} is obtained in closed-form. For the sake of completeness, the key steps leading to the above solution is provided in Appendix A.

B. Closed-Form a-Update

Given $d^{(n+1)}$ and $q^{(n)}$, the update of a in (22) is equivalent to solving the following problem

$$\mathcal{P}_{\mathbf{a}}(\mathbf{v}): \min_{\mathbf{a}} \sum_{j=1}^{G} \left(\|\mathbf{C}_{j}\mathbf{a}_{j}\|^{2} + \frac{\rho}{2} \sum_{i=1}^{G} \sum_{k=1}^{K_{i}} |d_{jik}^{(n+1)} - \mathbf{a}_{j}^{H} \mathbf{f}_{jik} + q_{jik}^{(n)}|^{2} \right)$$

Problem $\mathcal{P}_{a}(\mathbf{v})$ can be decomposed into G subproblems, one for each group j expressed as

$$\mathcal{P}_{\text{asub}}(\mathbf{v}):\min_{\mathbf{a}_{j}} \|\mathbf{C}_{j}\mathbf{a}_{j}\|^{2} + \frac{\rho}{2} \sum_{i=1}^{G} \sum_{k=1}^{K_{i}} |d_{jik}^{(n+1)} - \mathbf{a}_{j}^{H}\mathbf{f}_{jik} + q_{jik}^{(n)}|^{2}.$$

Problem $\mathcal{P}_{asub}(\mathbf{v})$ is an unconstrained convex optimization problem. Using the first-order optimality condition [18], we obtain the closed-form solution to $\mathcal{P}_{asub}(\mathbf{v})$ as

$$\mathbf{a}_{j}^{(n+1)} = \frac{\rho}{2} \left(\mathbf{C}_{j}^{H} \mathbf{C}_{j} + \frac{\rho}{2} \sum_{i=1}^{G} \sum_{k=1}^{K_{i}} \mathbf{f}_{jik} \mathbf{f}_{jik}^{H} \right)^{-1} \\ \cdot \sum_{i=1}^{G} \sum_{k=1}^{K_{i}} \left(d_{jik}^{(n+1)} + q_{jik}^{(n)} \right)^{*} \mathbf{f}_{jik}.$$
(28)

We summarize the ASCA algorithm in Algorithm 2. The main computational complexity of ASCA lies in computing $\{e_{1,jik}^{(n)}\}$ in (25) for $k \in \mathcal{K}_i, i, j \in \mathcal{G}$ and $\{\mathbf{a}_j^{(n+1)}\}$ in (28) for $j \in \mathcal{G}$ at each ADMM iteration. Note that computing $\mathbf{a}_j^{(n+1)}$ involves a matrix inversion with a complexity of $O(K_j^3)$. However, the matrix only depends on the channel vectors, and thus the matrix inversion only needs to be computed once at the beginning of ASCA. Thus for each ADMM iteration, only matrix-vector multiplications are involved. At each ADMM iteration, the related matrix-vector computation of $\{e_{1,jik}^{(n+1)}\}$ requires $12\left(\sum_{i=1}^{G}K_i\right)^2 + 2G\sum_{i=1}^{G}K_i$ flops, and that of $\{\mathbf{a}_j^{(n+1)}\}$ requires $6\left(\sum_{i=1}^{G}K_i\right)^2 + 2\sum_{i=1}^{G}(6K_i^2 + GK_i)$ flops. There are two layers of iterations in ASCA: the outer-layer SCA and the inner-layer ADMM to solve each $\mathcal{P}_{1SCA}(\mathbf{v})$. The convergence speed depends on the system parameters N and $K_{tot} = 15 \sim 60$ users, it typically takes $50 \sim 150$ iterations for the inner-layer ADMM-based algorithm to converge at each SCA iteration, and the outer layer typically takes $2 \sim 90$ iterations to converge.

Algorithm 2 ASCA Algorithm to Solve \mathcal{P}_1

1: Initialization: Set feasible initial point $\mathbf{v}^{(0)}$. Set ρ ; l = 0. 2: repeat // solve $\mathcal{P}_{1\text{ADMM}}(\mathbf{v}^{(l)})$ 3: Initialization: $\mathbf{a}^{(0)} = \mathbf{v}^{(l)}, \mathbf{d}^{(0)} = \mathbf{0}, \mathbf{q}^{(0)} = \mathbf{0}, n = 0.$ 4: repeat 5: Compute $\mathbf{d}^{(n+1)}$ by (21) with $\mathbf{a}^{(n)}$ and $\mathbf{q}^{(n)}$. 6: Compute $\mathbf{a}^{(n+1)}$ by (22) with $\mathbf{d}^{(n+1)}$ and $\mathbf{q}^{(n)}$. 7: Compute $q^{(n+1)}$ by (23) with $d^{(n+1)}$ and $a^{(n+1)}$. 8: $n \leftarrow n+1$. 9: until convergence 10: Set $\mathbf{v}^{(l+1)} = \mathbf{a}^{(n)}$. Set $l \leftarrow l+1$. 11: 12: **until** convergence

In Section VII, we will provide the simulation study to compare the convergence, performance, and the computational time of ESCA and ASCA.

C. Initialization for ASCA

As discussed earlier in Section IV-C, a feasible initial point is required by SCA to solve \mathcal{P}_1 . Below we propose an ADMMbased initialization method (AIM).

Using the ADMM-based algorithm above, we propose AIM for ASCA. The feasible point is computed by solving the following feasibility problem

$$\mathcal{P}'_{1 \text{fea}}$$
: Find $\{\mathbf{a}\}$ s.t. (4).

Using the auxiliary variables $d_{jik} = \mathbf{a}_j^H \mathbf{f}_{jik}$, for $k \in \mathcal{K}_i, i, j \in \mathcal{G}$, defined at the beginning of Section V, \mathcal{P}'_{1fea} is equivalently expressed as

$$\mathcal{P}'_{1\text{feaADMM}}: \text{ Find } \{\mathbf{a}, \mathbf{d}\}$$

s.t. $d_{jik} = \mathbf{a}_{j}^{H} \mathbf{f}_{jik}, \ k \in \mathcal{K}_{i}, i, j \in \mathcal{G}$
$$\frac{|d_{iik}|^{2}}{\sum_{j \neq i} |d_{jik}|^{2} + \sigma^{2}} \geq \gamma_{ik}, \ k \in \mathcal{K}_{i}, i \in \mathcal{G}.$$
(29)

Denote the feasible set for constraint (29) by $\widetilde{\mathcal{F}}$. The augmented Lagrangian of $\mathcal{P}'_{1\text{feaADMM}}$ is given by

$$\widetilde{\mathcal{L}}_{\tilde{\rho}}(\mathbf{d}, \mathbf{a}, \tilde{\mathbf{q}}) = I_{\widetilde{\mathcal{F}}}(\mathbf{d}) + \frac{\widetilde{\rho}}{2} \sum_{j=1}^{G} \sum_{i=1}^{G} \sum_{k=1}^{K_i} |d_{jik} - \mathbf{a}_j^H \mathbf{f}_{jik} + \widetilde{q}_{jik}|^2 \quad (30)$$

where indicator function I, penalty parameter $\tilde{\rho}$, dual variable \tilde{q}_{jik} are defined similarly as those in (20). Similar to the ADMM updating procedures in (21)–(23), based on (30), the AIM updating procedure for solving $\mathcal{P}'_{1\text{feaADMM}}$ is given as follows: At iteration n + 1,

$$\mathbf{d}^{(n+1)} = \arg\min_{\mathbf{d}\in\tilde{\mathcal{F}}} \sum_{j=1}^{G} \sum_{i=1}^{G} \sum_{k=1}^{K_i} |d_{jik} - \mathbf{a}_j^{(n)H} \mathbf{f}_{jik} + \tilde{q}_{jik}^{(n)}|^2, \quad (31)$$

$$\mathbf{a}^{(n+1)} = \arg\min_{\mathbf{a}} \sum_{j=1}^{G} \sum_{i=1}^{G} \sum_{k=1}^{K_i} |d_{jik}^{(n+1)} - \mathbf{a}_j^H \mathbf{f}_{jik} + \tilde{q}_{jik}^{(n)}|^2, \quad (32)$$

$$\tilde{q}_{jik}^{(n+1)} = \tilde{q}_{jik}^{(n)} + \left(d_{jik}^{(n+1)} - \mathbf{a}_j^{(n+1)H} \mathbf{f}_{jik} \right).$$
(33)

The updating steps in (31) and (32) can be derived in closedform, which are provided in Appendix B. The initial point for AIM is randomly chosen, which may not be feasible for SINR constraint (29) in $\mathcal{P}'_{1\text{feaADMM}}$. However, if AIM converges, the terminating point will satisfy SINR constraint (29), and the produced point a for $\mathbf{v}^{(0)}$ is feasible to $\mathcal{P}_{1\text{SCA}}(\mathbf{v})$. Note that since $\mathcal{P}'_{1\text{fea}}$ is a non-convex optimization problem, the ADMM procedure for AIM described above may not be guaranteed to converge. Thus, AIM is served as a heuristic algorithm. Our extensive simulation studies show that AIM converges to a feasible point with a very high probability.

Besides AIM, We can also use the SDR-based initialization method discussed in Section IV-C for initialization of ASCA, where a feasible initial point $\mathbf{v}^{(0)}$ for ASCA is obtained by applying SDR along with Gaussian randomization to solve \mathcal{P}_1 .

VI. SCALING SCHEME FOR THE MMF PROBLEM

In this section, with the proposed ESCA and ASCA for the QoS problem, we propose an efficient scheme to obtain a solution to the MMF problem S_o .

We transform the original MMF problem S_o into the following equivalent problem

$$\begin{aligned} \mathcal{S}'_{o}(\boldsymbol{\gamma}, P) &: \max_{\mathbf{w}, t} \quad t \\ \text{s.t. SINR}_{ik} \geq t \gamma_{ik}, \ k \in \mathcal{K}_{i}, i \in \mathcal{G} \\ &\sum_{i=1}^{G} \|\mathbf{w}_{i}\|^{2} \leq P \end{aligned}$$

where vector γ contains all SINR targets $\{\gamma_{ik}\}$ of all users in all groups. Furthermore, we parameterize the QoS problem \mathcal{P}_o as $\mathcal{P}_o(\gamma)$. It has been shown that $\mathcal{S}'_o(\gamma, P)$ and $\mathcal{P}_o(t\gamma)$ are the inverse problems to each other [5]. Specifically, denote the maximum objective value of $\mathcal{S}'_o(\gamma, P)$ by $t^o = \mathcal{S}'_o(\gamma, P)$. Let the minimum power objective value of $\mathcal{P}_o(\gamma')$ be $P = \mathcal{P}_o(\gamma')$ for some γ' . Then, we have the following inverse relation:

$$t^{o} = \mathcal{S}_{o}'(\boldsymbol{\gamma}, \mathcal{P}_{o}(t^{o}\boldsymbol{\gamma})), \quad P = \mathcal{P}_{o}(\mathcal{S}_{o}'(\boldsymbol{\gamma}, P)\boldsymbol{\gamma}).$$
(34)

Based on this relation, in the literature, the MMF problem $\mathcal{S}'_{o}(\boldsymbol{\gamma}, P)$ is typically solved via iteratively solving the QoS problem $\mathcal{P}_o(t\gamma)$ along with a bi-section search over t until the transmit power objective in $\mathcal{P}_{o}(t\gamma)$ is equal to P [2], [5], [6], [17], [21]. However, this approach is computationally expensive. As mentioned at the end of Section III-B, existing works use either SDR or SCA to compute a solution to $\mathcal{P}_o(t\gamma)$, where the second-order IPM is used to solve either relaxed or convexified approximate problem. As a result, iteratively solving $\mathcal{P}_o(t\gamma)$ incurs high computational complexity not suitable for largescale problems. Using the optimal beamforming structure \mathbf{w}_{i}^{o} in (3) can significantly reduce the required computation in the above approach, where $\mathcal{P}_o(t\gamma)$ can be converted to a much smaller weight optimization problem as in \mathcal{P}_1 . Nonetheless, resorting to the additional layer of iterations to solve the QoS problem is still computationally costly for the MMF problem as compared with the QoS problem itself, especially for large-scale problems [2].

To avoid the additional iterative procedure, we develop a closed-form scaling scheme for finding a solution to S_o directly. Specifically, we first obtain the solution to $\mathcal{P}_o(\gamma)$ by solving the smaller weight optimization problem \mathcal{P}_1 using either ESCA or ASCA proposed in Sections IV and V. Then, we scale this solution to $\mathcal{P}_o(\gamma)$ to obtain a solution to $S'_o(\gamma, P)$, such that the transmit power budget P is met. Parameterize S_o as $S_o(\gamma, P)$.

1: Solve $\mathcal{P}_{o}(\boldsymbol{\gamma})$ and attain solution $\mathbf{w}^{\mathsf{Q}}(\boldsymbol{\gamma})$. 2: Compute $P^{\mathsf{Q}}(\boldsymbol{\gamma}) = \sum_{i=1}^{G} \|\mathbf{w}_{i}^{\mathsf{Q}}(\boldsymbol{\gamma})\|^{2}$. 3: Compute $\mathbf{w}^{\mathsf{s}}(\boldsymbol{\gamma}, P) = \sqrt{\frac{P}{P^{\mathsf{Q}}(\boldsymbol{\gamma})}} \mathbf{w}^{\mathsf{Q}}(\boldsymbol{\gamma})$ as the solution to $\mathcal{S}_{o}(\boldsymbol{\gamma}, P)$.

This proposed scaling scheme is summarized in Algorithm 3. We show below that this scaling scheme provides a feasible solution to $S_o(\gamma, P)$, and we also bound its performance.

Proposition 1. Let $\mathbf{w}^{\mathsf{Q}}(\boldsymbol{\gamma})$ be a feasible beamforming solution to $\mathcal{P}_{o}(\boldsymbol{\gamma})$ produced by a given algorithm, with the achieved objective value denoted by $P^{\mathsf{Q}}(\boldsymbol{\gamma})$. Define $I_{ik}^{\mathsf{Q}}(\boldsymbol{\gamma}) \triangleq \sum_{j \neq i} |\mathbf{h}_{ik}^{H} \mathbf{w}_{j}^{\mathsf{Q}}(\boldsymbol{\gamma})|^{2}$. Then, the scaled beamforming vector $\mathbf{w}^{\mathsf{s}}(\boldsymbol{\gamma}, P) \triangleq \sqrt{\frac{P}{P^{\mathsf{Q}}(\boldsymbol{\gamma})}} \mathbf{w}^{\mathsf{Q}}(\boldsymbol{\gamma})$ is a feasible solution to $\mathcal{S}_{o}(\boldsymbol{\gamma}, P)$; the corresponding achieved objective value, denoted by $t^{\mathsf{s}}(\boldsymbol{\gamma}, P)$, satisfies

$$\frac{P}{P^{\mathsf{Q}}(\boldsymbol{\gamma})} \min_{i,k} \frac{I^{\mathsf{Q}}_{ik}(\boldsymbol{\gamma}) + \sigma^{2}}{\frac{P}{P^{\mathsf{Q}}(\boldsymbol{\gamma})}I^{\mathsf{Q}}_{ik}(\boldsymbol{\gamma}) + \sigma^{2}} \leq t^{\mathsf{s}}(\boldsymbol{\gamma}, P) \\
\leq \frac{P}{P^{o}(\boldsymbol{\gamma})} \max_{i,k} \frac{I^{\mathsf{Q}}_{ik}(\boldsymbol{\gamma}) + \sigma^{2}}{\frac{P}{P^{\mathsf{Q}}(\boldsymbol{\gamma})}I^{\mathsf{Q}}_{ik}(\boldsymbol{\gamma}) + \sigma^{2}} \quad (35)$$

where $P^{o}(\gamma)$ denotes the optimal objective value of $\mathcal{P}_{o}(\gamma)$.

Proof: See Appendix C.

Note that the tightness of the lower and upper bounds for $t^{s}(\gamma, P)$ in (35) depends on transmit power $P^{q}(\gamma)$, which is obtained by a given algorithm for $\mathcal{P}_{o}(\gamma)$ with solution $\mathbf{w}^{q}(\gamma)$. If the power budget P for the MMF problem $\mathcal{S}_{o}(\gamma, P)$ is more than the optimal power value for $\mathcal{P}_{o}(\gamma)$, *i.e.*, $P \geq P^{o}(\gamma)$, and the algorithm provides $P^{q}(\gamma) = P$, then $\mathbf{w}^{s}(\gamma, P) = \mathbf{w}^{q}(\gamma)$, and the bounds in (35) in this case are simplified to $1 \leq t^{s}(\gamma, P) \leq \frac{P}{P^{o}(\gamma)}$. Note that, often for the given algorithm, the solution $\mathbf{w}^{q}(\gamma)$ results in at least one SINR constraint being attained with equality, *i.e.*, SINR_{*ik*} = γ_{ik} , for some *i*, *k*, then we have $t^{s}(\gamma, P) = 1$. In a special case when $P^{q}(\gamma) = P = P^{o}(\gamma)$, we have $t^{s}(\gamma, P) = 1$, and also both the upper and lower bounds become 1. In this case, since $P = P^{o}(\gamma)$, from the inverse relation in (34), we have $t^{o} = 1$. Thus, $t^{s}(\gamma, P) = t^{o} = 1$, and the bounds in (35) are tight.

Comparing Algorithm 3 with the iterative bi-section search method using (34), we see that our proposed closed-form scaling scheme avoids iteratively solving $\mathcal{P}_2(t\gamma)$ along with a bi-section search over t, and thereby, enjoys a significant reduction of computational complexity. Moreover, we can directly apply the fast algorithm ESCA proposed in Section IV or ASCA proposed in Section V along with the optimal structure in (3) to provide a solution to $\mathcal{P}_o(\gamma)$ with this scaling scheme. This approach leads to two fast first-order algorithms for solving $\mathcal{S}_o(\gamma, P)$ in large-scale systems.

Remark 2. We point out that our scaling scheme for the multigroup multicast beamforming MMF problem is different from a similar scaling scheme proposed for the MMF problem in a single-group scenario in [22]. Specifically, the scheme in [22] first applies ZF beamforming for each group to eliminate inter-group interference. Once the interference is removed, for the equivalent single-group multicast beamforming problem,

[22] scales the beamforming solution of the single-group QoS problem to obtain a feasible solution to the MMF problem. In our scheme, we directly handle the original multi-group MMF problem containing inter-group interference. Our scheme scales the solution of the multi-group QoS problem $\mathcal{P}_{o}(\gamma)$ to solve the original MMF problem $S_o(\gamma, P)$. Note that the scheme in [22] has the lower and upper bounds for the objective value t (similar to t in $S'_o(\gamma, P)$) as $\left[\frac{P}{P^O(\gamma)}, \frac{P}{P^O(\gamma)}\right]$ with no interference present. In contrast, for our scheme, the lower and upper bounds in (35) contain additional terms w.r.t. $I_{ik}^{Q}(\gamma)$. Note that $I_{ik}^{Q}(\gamma)$ represents the inter-group interference to user k in group i by solution $w^{Q}(\gamma)$. Following this, both terms $I_{ik}^{Q}(\gamma) + \sigma^{2}$ and $\frac{P}{P^{Q}(\gamma)}I_{ik}^{Q}(\gamma) + \sigma^{2}$ in (35) are the interference plus noise term for user k in group i, where the latter is based on the scaled beamforming solution $\mathbf{w}^{s}(\boldsymbol{\gamma}, P)$. These additional terms in the lower and upper bounds in (35) represent the minimum and maximum inter-group interference ratios, respectively. In the special case of single group G = 1, the bounds in (35) reduces to $\left[\frac{P}{P^{Q}(\gamma)}, \frac{P}{P^{o}(\gamma)}\right]$ in [22]. Thus, the bounds in (35) can be viewed as a generalization of the bounds in [22] from singlegroup to multi-group settings with inter-group interference.

VII. SIMULATION RESULTS

We consider a default setup for downlink multi-group multicast beamforming, where G = 3 groups, $K_i = K$ users/group, $\forall i \in \mathcal{G}$, and the same SINR target for all users as $\gamma_{ik} = \gamma = 10$ dB, $\forall k, i$. The user channels are generated independently with an identical distribution as $\mathbf{h}_{ik} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I})$. The performance plots are obtained by averaging the results over 100 channel realizations per user.

For QoS problem \mathcal{P}_o , we consider the proposed two fast algorithms: ESCA in Algorithm 1 and ASCA in Algorithm 2. For ESCA, we set the step size $\alpha = 0.1$ and the constant c = 0.8. For ASCA, we set the penalty parameter $\rho =$ 0.2.¹ We consider different initialization methods discussed in Sections IV-C and V-C for ESCA and ASCA, respectively. Therefore, we refer to our algorithms as follows: 1) EIM-ESCA: ESCA with EIM initialization; 2) SDR-ECSA: ESCA with SDR initialization; 3) AIM-ASCA: ASCA with AIM initialization; 4) SDR-ASCA: ASCA with SDR initialization. Note that each algorithm solves the weight optimization problem \mathcal{P}_1 using the optimal beamforming structure in (3).² Besides the above four algorithms for solving \mathcal{P}_1 , we also consider the following methods for comparison:

- Lower Bound for \mathcal{P}_o : Obtained by solving the relaxed version of \mathcal{P}_o via SDR.
- SDR-CSCA [2]³: Solve \mathcal{P}_1 via SCA using $\mathcal{P}_{1SCA}(\mathbf{v})$,

¹We have studied different values of ρ and found $\rho = 0.2$ generally provides overall good performance and convergence speed for ASCA.

²Note that in (3), the exact expression of $\mathbf{R}(\boldsymbol{\lambda}^{o})$ is used and its computation is discussed below (3). One may utilize the asymptotic expression of $\mathbf{R}(\boldsymbol{\lambda}^{o})$ obtained in [2] to simplify the computation of $\mathbf{R}(\boldsymbol{\lambda}^{o})$. Since the fixed-point iterative method [2] for computing $\mathbf{R}(\boldsymbol{\lambda}^{o})$ is computationally inexpensive, using the asymptotic expression of $\mathbf{R}(\boldsymbol{\lambda}^{o})$ only brings a minor reduction of computation cost and thus is not considered in the simulation.

³Note that we only consider SDR-CSCA as the benchmark for comparison against our proposed algorithms. This is because SDR-CSCA is the stateof-the-art method with a near-optimal performance and substantially lower computational complexity than other existing algorithms in the literature. The comparison of SDR-CSCA and other existing algorithms in both performance and computational complexity has already been provided in [2], and adding other algorithms here will not provide additional insight or observation than what has been shown in [2].



Fig. 1. Convergence behavior for \mathcal{P}_o : Normalized power objective P_t/σ^2 over the inner-layer iterations at the first outer-layer SCA iteration (N = 500, K = 10).



Fig. 2. Convergence behavior for \mathcal{P}_o : Relative difference over the inner-layer iterations at the first outer-layer SCA iteration (N = 500, K = 10).

where each $\mathcal{P}_{1SCA}(\mathbf{v})$ is solved by the standard convex solver CVX, which uses IPM. The SDR method is used to generate an initial point.

For MMF problem S_o , we consider ESCA and ASCA with the proposed scaling scheme in Algorithm 3. They are denoted as ESCA-Scaling and ASCA-Scaling, respectively. For comparison, we also consider the following methods:

- Upper Bound for S_o : Obtained by solving the relaxed version of \mathcal{P}_o using SDR along with the bi-section search over t.
- ESCA-Bisection: Solve S_o via iteratively solving \mathcal{P}_o along with bi-section search over t. For solving \mathcal{P}_o , SDR-ESCA is used, where the optimal beamforming structure in (3) with the asymptotic expression of $\mathbf{R}(\lambda^o)$ obtained in [2] is applied.
- ASCA-Bisection: Similar to ESCA-Bisection, except that SDR-ASCA is used to solve \mathcal{P}_o .
- CSCA-Bisection [2]: Similar to ESCA-Bisection, except that SDR-CSCA is used to solve \mathcal{P}_o .
- SDR-Bisection [2]: Solve S_o via iteratively solving \mathcal{P}_o along with bi-section search over t. For solving \mathcal{P}_o , SDR along with the Gaussian randomization method is used, where the optimal beamforming structure in (3) with the



Fig. 3. Convergence behavior for \mathcal{P}_o : Normalized power objective P_t/σ^2 over the outer-layer SCA iterations (N = 500, K = 10).

asymptotic expression of $\mathbf{R}(\boldsymbol{\lambda}^{o})$ is applied.

A. Convergence Analysis

In this subsection, we study the convergence behavior of the two proposed fast algorithms (ESCA and ASCA) for the QoS problem \mathcal{P}_o . Note that each algorithm consists of two layers of iterations: Outer-layer SCA and the inner-layer iterative algorithm to solve each $\mathcal{P}_{1SCA}(\mathbf{v})$. Let $P_t = \sum_{i=1}^G \|\mathbf{w}_i\|^2$ denote the total transmit power. Fig. 1 shows the trajectory of normalized transmit power $P_{\rm t}/\sigma^2$ over the inner-layer iterations at the first outer-layer SCA iteration. We set N = 500 and K = 10. We observe that SDR-ESCA and SDR-ASCA converge faster and result in lower $P_{\rm t}/\sigma^2$ than EIM-ESCA and AIM-ASCA. This shows that the SDR initialization method provides a better initial point than the other initialization methods. Between SDR-ESCA and SDR-ASCA, we observe a similar convergence rate. Comparing EIM-ESCA and AIM-ASCA, we see that AIM-ASCA converges faster than EIM-ESCA. Next, we consider the convergence behavior of the relative difference of the optimization variable for each algorithm. Specifically, we consider the normalized difference $\frac{\|\mathbf{x}^{(n+1)}-\mathbf{x}^{(n)}\|}{\|\mathbf{x}^{(n+1)}\|}$ in two consecutive inner-layer iterations of ESCA in Algorithm 1, and similarly, $\frac{\|\mathbf{a}^{(n+1)}-\mathbf{a}^{(n)}\|}{\|\mathbf{a}^{(n+1)}\|}$ for ASCA in Algorithm 2. Fig. 2 shows these relative differences for different algorithms at the first outerlayer SCA iteration for N = 500 and K = 10. Again, we observe that the SDR initialization method results in a faster convergence than the other initialization methods. Both SDR-ESCA and SDR-ASCA reach a relative difference of 10^{-3} within 50 iterations, with SDR-ASCA converging slightly faster than SDR-ESCA. Comparing EIM-ESCA and AIM-ASCA, we see that AIM-ASCA provides a faster convergence than EIM-ESCA, reaching a relative difference of 10^{-3} within 50 iterations.

We now study the convergence behavior of different algorithms over the outer-layer SCA iterations. Fig. 3 shows the convergence behavior of our proposed algorithms at the outer-layer SCA iterations, for N = 500 and K = 10. We set the inner-layer convergence threshold for the proposed algorithms such that they converge to nearly the same value. For EIM-ESCA and SDR-ESCA, we set the inner-layer convergence threshold to be $\frac{\|\mathbf{x}^{(n+1)}-\mathbf{x}^{(n)}\|}{\|\mathbf{x}^{(n+1)}\|} \leq 10^{-3}$. Note that although Fig. 1



Fig. 4. QoS: Normalized transmit power $P_{\rm t}/\sigma^2$ vs. N~(G=3,~K=10). TABLE I

QOS: AVERAGE COMPUTATION	Time over N (sec.)	(G = 3, K = 10)	

N	I	100	200	300	400	500
EIM-ESCA	I	2.694	3.032	2.708	2.589	2.711
AIM-ASCA	I	1.150	2.415	3.195	4.092	4.339
SDR-ESCA	I	0.327	0.237	0.233	0.224	0.243
SDR-ASCA	I	0.076	0.051	0.061	0.071	0.088
SDR-CSCA [2]		7.179	6.126	6.400	5.963	6.500
EIM (Init. method)	I	0.037	0.044	0.045	0.050	0.057
AIM (Init. method)	I	0.0068	0.0058	0.0059	0.0063	0.0073
SDR (Init. method)	I	1.050	1.038	1.104	1.103	1.137

shows that at the first outer-layer SCA iteration, EIM-ESCA converges to a higher value of $P_{\rm t}/\sigma^2$ than that of SDR-ESCA over the inner-layer iterations, Fig. 3 shows that EIM-ESCA eventually converges to nearly the same value of $P_{\rm t}/\sigma^2$ as that of SDR-ESCA over the outer-layer SCA iterations. For ASCA, we set the inner-layer convergence threshold $\frac{\|\mathbf{a}^{(n+1)}-\mathbf{a}^{(n)}\|}{\|\mathbf{a}^{(n+1)}\|}$ $\leq 10^{-3}$ for SDR-ASCA and $\frac{\|\mathbf{a}^{(n+1)}-\mathbf{a}^{(n)}\|}{\|\mathbf{a}^{(n+1)}\|} \leq 0.2 \times 10^{-3}$ for AIM-ASCA. Our experiments show that a tighter innerlayer threshold for AIM-ASCA than that of SDR-ASCA is needed to converge to the same value of $P_{\rm t}/\sigma^2$ as the rest of algorithms at the outer-layer iterations. As we see in Fig. 3, all algorithms converge to the same value of the normalized transmit power $P_{\rm t}/\sigma^2$ within 35 SCA iterations. Again, for different initialization methods, since the SDR method provides a better initial point than EIM and AIM methods, it leads to a faster convergence for both ESCA and ASCA, where only less than 5 SCA iterations are required to reach convergence. When the same SDR initialization is used, ESCA, ASCA, and CSCA have a similar convergence behavior.

B. Performance Comparison for the QoS Problem

We now compare the performance of different algorithms for the QoS problem. We set SINR target $\gamma = 10$ dB. Also, for all algorithms in comparison, we set the convergence threshold for the outer-layer SCA iterations as $\frac{\|\mathbf{v}^{(l+1)}-\mathbf{v}^{(l)}\|}{\|\mathbf{v}^{(l+1)}\|} \leq 10^{-3}$. Fig. 4 shows the normalized transmit power P_t/σ^2 vs. N by different algorithms, for G = 3 and K = 10. We see that our proposed algorithms have a similar performance to that of SDR-CSCA in



Fig. 5. QoS: Normalized transmit power P_t/σ^2 vs. K (G = 3, N = 100).

TABLE II QOS: AVERAGE COMPUTATION TIME OVER K (SEC.) (G = 3, N = 100).

K	I	5	7	10	15	20	35
EIM-ESCA	I	1.098	1.689	2.737	5.498	8.685	14.33
AIM-ASCA	I	0.519	0.725	1.234	2.043	3.146	5.426
SDR-ESCA	I	0.056	0.112	0.373	0.856	1.666	4.883
SDR-ASCA	I	0.012	0.024	0.087	0.216	0.433	2.594
SDR-CSCA [2]	ļ	1.747	3.502	7.592	13.79	21.52	57.85
EIM (Init. method)	I	0.019	0.025	0.038	0.063	0.090	0.450
AIM (Init. method)	1	0.0048	0.0055	0.0067	0.012	0.018	0.077
SDR (Init. method)	I	0.545	0.727	1.094	1.852	3.489	22.79
SDR-CSCA [2]EIM (Init. method)AIM (Init. method)SDR (Init. method)		1.747 0.019 0.0048 0.545	3.502 0.025 0.0055 0.727	7.592 0.038 0.0067 1.094	13.790.0630.0121.852	21.52 0.090 0.018 3.489	5 0 0 2

[2], and all algorithms nearly attain the lower bound for \mathcal{P}_o .⁴ This demonstrates that our proposed fast algorithms achieve a near-optimal performance. The computational advantages of ESCA and ASCA are shown in Table I, where we list the average computation time by each algorithm used for the plots in Fig. 4. The first five rows show the computation times of different algorithms, excluding the initialization. We observe that, by using the optimal structure in (3), the computation times of all algorithms remain roughly unchanged as N increases. Under the same SDR initialization method, the computation times of SDR-ESCA and SDR-ASCA are only about 4% and 1% of that of SDR-CSCA, respectively, and SDR-ASCA has a smaller computation time than SDR-ESCA. The computation times of EIM-ESCA and AIM-ASCA are both about 40% of that of SDR-CSCA. The computation time of AIM-ASCA increases with N more noticeably than other algorithms. As a result, AIM-ASCA is initially faster than EIM-ESCA for $N \leq 200$, but its computation time increases and becomes slower than EIM-ESCA for N > 300. The reason is due to the quality of initialization as will be explained below. The last three rows in Table I show the computation times of different initialization methods. We see that the EIM is a fast initialization method, with its computation time about 4% of that of SDR. The AIM is the fastest one among all three initialization methods, with its computation time about 15%and 0.5% of that of EIM and SDR, respectively. However, the quality of AIM initialization deteriorates as N increases,

⁴Note that the lower bound is only shown until N = 300 in Fig. 4 due to the high computational complexity involved in generating the lower bound, which increases fast with N and becomes impractical for N beyond 300. Similarly, the upper bound shown in Fig. 6 is provided until N = 200.

unlike other initialization methods, leading to more iterations for convergence and a longer computation time. This is evidenced by the computation time of AIM-ASCA, which increases more noticeably as N increases.

Fig. 5 shows P_t/σ^2 vs. K by different algorithms, for G = 3and N = 100. Again, the performance of SDR-ESCA and SDR-ASCA are nearly identical to that of SDR-CSCA and nearly attain the lower bound. The performance gaps of EIM-ESCA and AIM-ASCA to the lower bound are more noticeable as Kbecomes large, although they are still within 0.6 dB. Similar to Table I, the corresponding average computation times of these algorithms are shown in Table II. We see that both ESCA and ASCA are considerably faster than CSCA to compute the solution. In particular, unlike CSCA, their computation times only mildly increase with K. For the initialization methods, EIM and AIM are faster and much more scalable over Kthan SDR. The complexity of SDR increases noticeably over K, and thus, it is not a suitable initialization method for very large systems. Overall, in terms of the total computation time (including initialization and the algorithm itself), for N = 100and $K \leq 20$, SDR-ASCA is the fastest one among all algorithms. AIM-ASCA offers comparable computation time as SDR-ASCA. For N = 100 and $K \ge 35$, AIM-ASCA and EIM-ESCA offers faster computation time than the rest using SDR for initialization.

Comparison Summary: Based on the above simulation analysis, we have the following comparison summary of the two proposed algorithms:

- Among the initialization methods (EIM, AIM, and SDR), SDR provides the best initialization point, which leads to faster computation for both ESCA and ASCA (*i.e.*, SDR-ESCA and SDR-ASCA). However, the computational complexity of SDR is high. As K becomes large (*e.g.*, K > 20), SDR-based initialization becomes computationally expensive and is not recommended. EIM and AIM are both very low-complexity methods over N and K. AIM provides faster initialization than EIM. However, the quality of the initial point that AIM provides deteriorates over N, leading to a noticeable increase of computation time by AIM-ASCA over N. In contrast, the computation time of EIM-ESCA remains roughly unchanged over N.
- Both ESCA and ASCA provide closed-form updates in each iteration. The computational complexity of ESCA and ASCA grow over K. From our study, we conclude that:
 - For the system with users per group K ≤ 15, SDR-ASCA is generally the fastest among all algorithms. When N ≤ 100, AIM-ASCA is similar to SDR-ASCA. However, the complexity of SDR-ASCA increases over K, and that of AIM-ASCA increases over N. As both K and N grow, SDR-ASCA and AIM-ASCA have higher computation time and are no longer preferred.
 - For a large-scale system where both N and K are large (e.g., $N \ge 500$, $K \ge 35$), EIM-ESCA is expected to provide the fastest computation time among these algorithms and thus is preferred.

C. Performance Comparison for the MMF Problem

We now present the performance of our proposed Algorithm 3 for the MMF problem S_o . We set the maximum transmit



Fig. 6. MMF: Minimum SINR vs. N (G = 3, K = 10).

TABLE III MMF: Average Computation Time over N (sec.) (G = 3, K = 10).

N	50	100	200	300	400
ESCA-Scaling	0.482	0.327	0.237	0.233	0.224
ASCA-Scaling	0.099	0.076	0.051	0.061	0.071
ESCA-Bisection	16.43	21.11	33.02	46.18	44.63
ASCA-Bisection	12.35	12.11	13.06	14.09	17.94
CSCA-Bisection [2]	99.04	87.33	81.62	84.10	99.86
SDR-Bisection [2]	11.19	15.46	19.85	15.82	22.12

power budget against noise variance as $P/\sigma^2 = 10$ dB. Fig. 6 plots the average minimum SINR vs. N, and Table III shows the corresponding computation time by these algorithms, for G = 3 and K = 10. Both ESCA-Bisection and ASCA-Bisection provide near-identical performance to that of CSCA-Bisection, and they are nearly optimal as compared with the upper bound, but with much lower computation times. The proposed simple ESCA-Scaling and ASCA-Scaling algorithms for the MMF problem nearly attain the upper bound for N <100. Their performance gap to the upper bound increases as N increases, indicating the accuracy of the scaling degrades as N becomes large. At N = 400, the gap is about 1 dB. Nonetheless, the ESCA-Scaling and ASCA-Scaling are several orders of magnitude faster than ESCA-Bisection and ASCA-Bisection. SDR-Bisection has worse performance than all the rest algorithms. This is because, for the QoS problem \mathcal{P}_1 at each bi-section iteration, the SDR approximation deteriorates when the number of constraints (*GK*) for \mathcal{P}_1 is large. Finally, Fig. 7 shows the average minimum SINR vs. K by different algorithms, for G = 3 and N = 100, with the corresponding computation time shown in Table IV. Except for the SDR-Bisection, all the proposed algorithms nearly attain the upper bound and thus are nearly optimal. In particular, ESCA-Scaling and ASCA-Scaling maintain their near-optimal performance as K increases. The computational advantage of ESCA-Scaling and ASCA-Scaling is clearly seen in Table IV, where both algorithms are substantially faster in computing a solution than the other algorithms.

VIII. CONCLUSION

In this work, exploiting the optimal multicast beamforming structure, we proposed two fast computational algorithms,



Fig. 7. MMF: Minimum SINR vs. K (G = 3, N = 100).

TABLE IV MMF: Average Computation Time over K (sec.) (G = 3, N = 100).

K	I	5	7	10	15
ESCA-Scaling	I	0.056	0.112	0.373	0.856
ASCA-Scaling	I	0.012	0.024	0.087	0.216
ESCA-Bisection	I	7.297	12.62	21.26	39.59
ASCA-Bisection	I	6.192	8.628	12.13	25.18
CSCA-Bisection [2]	I	24.91	48.47	88.40	178.5
SDR-Bisection [2]	I	6.300	10.10	15.53	25.75

ESCA and ASCA, for multi-group multicast beamforming design. For the QoS problem solved by the SCA method, these two algorithms provide efficient computational methods for solving the convex subproblems of SCA. At each SCA iteration, ESCA implements a dual saddle point reformulation along with the extragradient method to solve the convex subproblem; ASCA constructs an ADMM procedure in a form that decomposes the convex subproblem into multiple smaller subproblems for parallel computing with closed-form updates. To provide effective initial feasible points for SCA to facilitate fast convergence, we proposed three initialization methods based on the extragradient method, ADMM, and SDR. For the MMF problem, we further proposed a simple closed-form scaling approach based on the solution to the QoS problem, avoiding the high computational complexity involved in iteratively solving the QoS problems, while offering bounded performance guarantee. Our simulation studies showed that the proposed ESCA and ASCA provide near-optimal performance with substantially lower computational complexity than the state-of-the-art algorithms for largescale systems.

Finally, note that in this work, we assumed single-antenna users in our system model for multicast beamforming design. For the case of multi-antenna users, with a given linear receiver processing technique implemented by the receiver, each MIMO channel can be converted into an equivalent MISO channel. Therefore, our proposed algorithms can be relatively straightforward to be applied to the case of multi-antenna users with their receiver processing techniques given.

APPENDIX A Derivation of \mathbf{d}_{ik}^{o} in (26)

The Lagrangian of $\mathcal{P}_{dsub}(\mathbf{v})$ is given by

$$\mathcal{L}(\mathbf{d}_{ik},\nu_{ik}) = \sum_{j=1}^{G} |d_{jik} - e_{1,jik}^{(n)}|^2 + \nu_{ik}e_{2,ik} + \nu_{ik}\gamma_{ik}\sum_{j\neq i} |d_{jik}|^2 - 2\nu_{ik}\mathfrak{Re}\{d_{iik}e_{3,ik}\}$$

where $\nu_{ik} \geq 0$ is the Lagrange multiplier associated with constraint (24). Since $\mathcal{P}_{dsub}(\mathbf{v})$ is convex, the optimal \mathbf{d}_{ik}^{o} and ν_{ik}^{o} satisfy the KKT conditions [18]. Setting the gradient $\nabla_{\mathbf{d}_{ik}} \mathcal{L}(\mathbf{d}_{ik}, \nu_{ik}^{o}) = 0$, we obtain \mathbf{d}_{ik}^{o} in (26). Substituting the expression of d_{jik}^{o} in (26) into the constraint of $\mathcal{P}_{dsub}(\mathbf{v})$ yields

$$f(\nu_{ik}^{o}) \triangleq e_{2,ik} + \gamma_{ik} \sum_{j \neq i} \frac{|e_{1,jik}^{(n)}|^2}{(1 + \nu_{ik}^{o}\gamma_{ik})^2} - 2\mathfrak{Re}\{e_{3,ik}e_{1,iik}^{(n)}\} - 2\nu_{ik}^{o}|e_{3,ik}|^2 \le 0.$$

It is shown in [21, Appendix A] that the function $f(\nu_{ik}^o)$ is strictly decreasing for $\nu_{ik}^o \ge 0$. By the complementary slackness condition, we have $\nu_{ik}^o f(\nu_{ik}^o) = 0$. If $f(0) \le 0$, then $f(\nu_{ik}^o) < 0$ for any $\nu_{ik}^o \ge 0$, and we have $\nu_{ik}^o = 0$. Otherwise, $\nu_{ik}^o = 0$ is not an feasible dual solution for $\mathcal{P}_{dsub}(\mathbf{v})$; thus, $f(\nu_{ik}^o) = 0$, and ν_{ik}^o is the unique real positive root of the cubic equation (27), whose roots can be obtained by the cubic formula [21].

APPENDIX B CLOSED-FORM UPDATING STEPS FOR AIM

1) Closed-Form d-update in (31): Similar to that in Section V-A, given $\mathbf{a}^{(n)}$ and $\tilde{\mathbf{q}}^{(n)}$, the optimization problem in (31) can be decomposed into K_{tot} subproblems, one for each user k in group i given by

$$\mathcal{P}'_{dsub} : \min_{\mathbf{d}_{ik}} \sum_{j=1}^{G} |d_{jik} - \tilde{e}^{(n)}_{1,jik}|^2$$

s.t. $\gamma_{ik} \sum_{j \neq i} |d_{jik}|^2 + \gamma_{ik} \sigma^2 - |d_{iik}|^2 \le 0.$ (36)

where $\tilde{e}_{1,jik}^{(n)} \triangleq \mathbf{a}_{j}^{(n)H} \mathbf{f}_{jik} - \tilde{q}_{jik}^{(n)}$, for $k \in \mathcal{K}_i, i, j \in \mathcal{G}$. Problem \mathcal{P}'_{dsub} is a non-convex QCQP-1 problem similar to $\mathcal{P}_{dsub}(\mathbf{v})$. Since it satisfies Slater's condition, the strong duality holds [18, Appendix B], and the optimal \mathbf{d}_{ik}^o satisfies the KKT conditions. Let $\mu_{ik}^o \ge 0$ be the optimal Lagrange multiplier associated with constraint (36). For \mathcal{P}'_{dsub} being feasible, we have $0 \le \mu_{ik}^o \le 1$ [18, Appendix B]. Thus, using the KKT conditions and with a procedure similar to that in Appendix A, we have the closed-form optimal solution \mathbf{d}_{ik} to \mathcal{P}'_{dsub} given by

$$d^{o}_{jik} = \begin{cases} \frac{\tilde{e}^{(n)}_{1,ik}}{1-\mu^{o}_{ik}} & \text{if } j=i, \\ \frac{\tilde{e}^{(n)}_{1,jik}}{1+\mu^{o}_{ik}\gamma_{ik}} & \text{otherwise} \end{cases}$$

where $\mu_{ik}^o = 0$ if $\gamma_{ik} \sum_{j \neq i} |\tilde{e}_{1,jik}^{(n)}|^2 + \gamma_{ik}\sigma^2 - |\tilde{e}_{1,iik}^{(n)}|^2 \leq 0$; otherwise the following quartic equation is guaranteed to have a unique root in (0, 1), and μ_{ik}^o is this unique root:

$$\gamma_{ik} \frac{\sum_{j \neq i} |\tilde{e}_{1,jik}^{(n)}|^2}{(1 + \mu_{ik}^o \gamma_{ik})^2} + \gamma_{ik} \sigma^2 - \frac{|\tilde{e}_{1,iik}^{(n)}|^2}{(1 - \mu_{ik}^o)^2} = 0.$$

2) Closed-Form a-update in (32): Given $d^{(n+1)}$ and $\tilde{q}^{(n)}$, the optimization problem in (32) can be decomposed into G subproblems, one for each group j given by

$$\mathcal{P}'_{\text{asub}} : \min_{\mathbf{a}} \sum_{i=1}^{G} \sum_{k=1}^{K_i} |d_{jik}^{(n+1)} - \mathbf{a}_j^H \mathbf{f}_{jik} + \tilde{q}_{jik}^{(n)}|^2,$$

which is an unconstrained quadratic convex optimization problem. The closed-form solution to \mathcal{P}'_{asub} is expressed as

$$\mathbf{a}_{j}^{(n+1)} = \left(\sum_{i=1}^{G}\sum_{k=1}^{K_{i}}\mathbf{f}_{jik}\mathbf{f}_{jik}^{H}\right)^{-1}\sum_{i=1}^{G}\sum_{k=1}^{K_{i}}\left(d_{jik}^{(n+1)} + \tilde{q}_{jik}^{(n)}\right)^{*}\mathbf{f}_{jik}.$$

APPENDIX C PROOF OF PROPOSITION 1

Proof: It is straightforward to check that the scaled beamforming vector $\mathbf{w}^{s}(\boldsymbol{\gamma}, P) = \sqrt{\frac{P}{P^{\mathbb{Q}}(\boldsymbol{\gamma})}} \mathbf{w}^{\mathbb{Q}}(\boldsymbol{\gamma})$ satisfies constraint (2) and therefore is a feasible solution to $\mathcal{S}_{o}(\boldsymbol{\gamma}, P)$. The achieved objective value $t^{s}(\boldsymbol{\gamma}, P)$ corresponding to $\mathbf{w}^{s}(\boldsymbol{\gamma}, P)$ is expressed in terms of $\mathbf{w}^{\mathbb{Q}}(\boldsymbol{\gamma})$ as

$$t^{s}(\boldsymbol{\gamma}, P) = \min_{i,k} \frac{1}{\gamma_{ik}} \frac{\frac{P}{P^{Q}(\boldsymbol{\gamma})} |\mathbf{h}_{ik}^{H} \mathbf{w}_{i}^{Q}(\boldsymbol{\gamma})|^{2}}{\frac{P}{P^{Q}(\boldsymbol{\gamma})} I_{ik}^{Q}(\boldsymbol{\gamma}) + \sigma^{2}}$$
$$= \min_{i,k} \frac{\frac{P}{P^{Q}(\boldsymbol{\gamma})} |\mathbf{h}_{ik}^{H} \mathbf{w}_{i}^{Q}(\boldsymbol{\gamma})|^{2}}{\gamma_{ik} I_{ik}^{Q}(\boldsymbol{\gamma}) + \gamma_{ik} \sigma^{2}} \frac{I_{ik}^{Q}(\boldsymbol{\gamma}) + \sigma^{2}}{\frac{P}{P^{Q}(\boldsymbol{\gamma})} I_{ik}^{Q}(\boldsymbol{\gamma}) + \sigma^{2}}.$$
 (37)

Define $t^{\mathsf{Q}}(\boldsymbol{\gamma}) \triangleq \min_{i,k} \frac{|\mathbf{h}_{ik}^{H} \mathbf{w}_{i}^{\mathsf{Q}}(\boldsymbol{\gamma})|^{2}}{\gamma_{ik} I_{ik}^{\mathsf{Q}} (\boldsymbol{\gamma}) + \gamma_{ik} \sigma^{2}}$. Since $\mathbf{w}^{\mathsf{Q}}(\boldsymbol{\gamma})$ satisfies constraint (1) in $\mathcal{P}_{o}(\boldsymbol{\gamma})$ as a feasible solution, we have $\frac{|\mathbf{h}_{ik}^{H} \mathbf{w}_{i}^{\mathsf{Q}}(\boldsymbol{\gamma})|^{2}}{I_{ik}^{\mathsf{Q}}(\boldsymbol{\gamma}) + \sigma^{2}} \geq \gamma_{ik}$ for $k \in \mathcal{K}_{i}, i \in \mathcal{G}$. It follows that $t^{\mathsf{Q}}(\boldsymbol{\gamma}) \geq 1$. Based on this, from (37), we have

$$t^{s}(\boldsymbol{\gamma}, P) \geq \frac{P}{P^{q}(\boldsymbol{\gamma})} t^{q}(\boldsymbol{\gamma}) \min_{i,k} \frac{I_{ik}^{q}(\boldsymbol{\gamma}) + \sigma^{2}}{\frac{P}{P^{q}(\boldsymbol{\gamma})} I_{ik}^{q}(\boldsymbol{\gamma}) + \sigma^{2}}$$
$$\geq \frac{P}{P^{q}(\boldsymbol{\gamma})} \min_{i,k} \frac{I_{ik}^{q}(\boldsymbol{\gamma}) + \sigma^{2}}{\frac{P}{P^{q}(\boldsymbol{\gamma})} I_{ik}^{q}(\boldsymbol{\gamma}) + \sigma^{2}}.$$
(38)

From (37), we can also obtain an upper bound on $t^{s}(\gamma, P)$ as

$$t^{s}(\boldsymbol{\gamma}, P) \leq \frac{P}{P^{q}(\boldsymbol{\gamma})} t^{q}(\boldsymbol{\gamma}) \frac{I^{q}_{i'k'}(\boldsymbol{\gamma}) + \sigma^{2}}{\frac{P}{P^{q}(\boldsymbol{\gamma})} I^{q}_{i'k'}(\boldsymbol{\gamma}) + \sigma^{2}}$$
$$\leq \frac{P}{P^{q}(\boldsymbol{\gamma})} t^{q}(\boldsymbol{\gamma}) \max_{i,k} \frac{I^{q}_{ik}(\boldsymbol{\gamma}) + \sigma^{2}}{\frac{P}{P^{q}(\boldsymbol{\gamma})} I^{q}_{ik}(\boldsymbol{\gamma}) + \sigma^{2}}$$
(39)

where $\{i', k'\} = \underset{i,k}{\operatorname{arg\,min}} \frac{|\mathbf{h}_{ik}^{H} \mathbf{w}_{i}^{Q}(\gamma)|^{2}}{\gamma_{ik} I_{ik}^{Q}(\gamma) + \gamma_{ik} \sigma^{2}}$. Note that the scaled

beamforming vector $\frac{\mathbf{w}^{\mathbb{Q}}(\gamma)}{\sqrt{t^{\mathbb{Q}}(\gamma)}}$ results in the transmit power $\frac{P^{\mathbb{Q}}(\gamma)}{t^{\mathbb{Q}}(\gamma)}$. Then, with $\frac{\mathbf{w}^{\mathbb{Q}}(\gamma)}{\sqrt{t^{\mathbb{Q}}(\gamma)}}$, the achieved weighted SINR, for any $k \in$

$$\mathcal{K}_i, i \in \mathcal{G}, \text{ is given by}$$

$$\frac{1}{\gamma_{ik}} \frac{\frac{1}{t^{\mathbb{Q}}(\gamma)} |\mathbf{h}_{ik}^{H} \mathbf{w}_{i}^{\mathbb{Q}}(\gamma)|^{2}}{\frac{1}{t^{\mathbb{Q}}(\gamma)} I_{ik}^{\mathbb{Q}}(\gamma) + \sigma^{2}} \stackrel{(a)}{\geq} \frac{1}{\gamma_{ik} t^{\mathbb{Q}}(\gamma)} \frac{|\mathbf{h}_{ik}^{H} \mathbf{w}_{i}^{\mathbb{Q}}(\gamma)|^{2}}{I_{ik}^{\mathbb{Q}}(\gamma) + \sigma^{2}} \stackrel{(b)}{\geq} 1 \quad (40)$$

where (a) is due to $t^{\varrho}(\gamma) \geq 1$, and (b) is from the definition of $t^{\varrho}(\gamma)$. Thus, the beamforming vector $\frac{\mathbf{w}^{\varrho}(\gamma)}{\sqrt{t^{\varrho}(\gamma)}}$ is feasible to $\mathcal{P}_{o}(\gamma)$ and $\frac{P^{\varrho}(\gamma)}{t^{\varrho}(\gamma)} \geq P^{o}(\gamma)$, where $P^{o}(\gamma)$ is the minimum transmit power in $\mathcal{P}_o(\gamma)$. Applying $\frac{P^{\mathbb{Q}}(\gamma)}{t^{\mathbb{Q}}(\gamma)} \geq P^o(\gamma)$ to the RHS of (39) yields

$$t^{s}(\boldsymbol{\gamma}, P) \leq \frac{P}{P^{o}(\boldsymbol{\gamma})} \max_{i,k} \frac{I_{ik}^{q}(\boldsymbol{\gamma}) + \sigma^{2}}{\frac{P}{P^{Q}(\boldsymbol{\gamma})}I_{ik}^{q}(\boldsymbol{\gamma}) + \sigma^{2}}.$$
 (41)

Combining (38) and (41), we have (35).

REFERENCES

- C. Zhang, M. Dong, and B. Liang, "First-order fast algorithm for structurally optimal multi-group multicast beamforming in large-scale systems," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, Jun. 2021, pp. 4790–4794.
- [2] M. Dong and Q. Wang, "Multi-group multicast beamforming: Optimal structure and efficient algorithms," *IEEE Trans. Signal Process.*, vol. 68, pp. 3738–3753, May 2020.
- [3] 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.
- [4] 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.
- [5] 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.
- [6] 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.
- [7] M. Jordan, X. Gong, and G. Ascheid, "Multicell multicast beamforming with delayed SNR feedback," in 2009 IEEE Global Telecommun. Conf., Nov. 2009, pp. 1–6.
- [8] 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.
- [9] N. Bornhorst, M. Pesavento, and A. B. Gershman, "Distributed beamforming for multi-group multicasting relay networks," *IEEE Trans. Signal Process.*, vol. 60, no. 1, pp. 221–232, Jan. 2012.
- Process., vol. 60, no. 1, pp. 221–232, Jan. 2012.
 [10] M. Dong and B. Liang, "Multicast relay beamforming through dual approach," in *Proc. IEEE Int. Workshop Comput. Advances Multi-Sensor Adaptive Process.*, Dec. 2013, pp. 492–495.
- [11] K. T. Phan, S. A. Vorobyov, N. D. Sidiropoulos, and C. Tellambura, "Spectrum sharing in wireless networks via QoS-aware secondary multicast beamforming," *IEEE Trans. Signal Process.*, vol. 57, no. 6, pp. 2323–2335, Jun. 2009.
- [12] Y. Huang, Q. Li, W.-K. Ma, and S. Zhang, "Robust multicast beamforming for spectrum sharing-based cognitive radios," *IEEE Trans. Signal Process.*, vol. 60, no. 1, pp. 527–533, Jan. 2012.
- [13] Z.-Q. Luo, W.-K. Ma, A. M.-C. So, Y. Ye, and S. Zhang, "Semidefinite relaxation of quadratic optimization problems," *IEEE Signal Process. Mag.*, vol. 27, no. 3, pp. 20–34, May 2010.
- [14] B. R. Marks and G. P. Wright, "A general inner approximation algorithm for nonconvex mathematical programs," *Oper. Res.*, vol. 26, no. 4, pp. 681–683, Aug. 1978.
- [15] L. 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, Jul. 2015.
- [17] D. Christopoulos, S. Chatzinotas, and B. Ottersten, "Multicast multigroup beamforming for per-antenna power constrained large-scale arrays," in *Proc. IEEE Int. Workshop Signal Process. Advances Wireless Commun.*, Jun. 2015, pp. 271–275.
- [18] S. Boyd and L. Vandenberghe, *Convex Optimization*, Cambridge, U.K.: Cambridge Univ. Press, 2004.
- [19] A. Konar and N. D. Sidiropoulos, "Fast approximation algorithms for a class of non-convex QCQP problems using first-order methods," *IEEE Trans. Signal Process.*, vol. 65, no. 13, pp. 3494–3509, Jul. 2017.
- [20] M. S. Ibrahim, A. Konar, and N. D. Sidiropoulos, "Fast algorithms for joint multicast beamforming and antenna selection in massive MIMO," *IEEE Trans. Signal Process.*, vol. 68, pp. 1897–1909, Mar. 2020.
- [21] 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.
- [22] 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.

- [23] J. Yu and M. Dong, "Low-complexity weighted MRT multicast beamforming in massive MIMO cellular networks," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, Apr. 2018, pp. 3849–3853.
- [24] J. Yu and M. Dong, "Distributed low-complexity multi-cell coordinated multicast beamforming with large-scale antennas," in *Proc. IEEE Int. Workshop Signal Process. Advances Wireless Commun.*, Jun. 2018, pp. 1–5.
- [25] 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.
- [26] M. Sadeghi and C. Yuen, "Multi-cell multi-group massive MIMO multicasting: An asymptotic analysis," in *Proc. IEEE Global Commun. Conf.*, Dec. 2015, pp. 1–6.
- [27] O. Tervo, L. Tran, H. Pennanen, S. Chatzinotas, B. Ottersten, and M. Juntti, "Energy-efficient multicell multigroup multicasting with joint beamforming and antenna selection," *IEEE Trans. Signal Process.*, vol. 66, no. 18, pp. 4904–4919, Sept. 2018.
- [28] M. Sadeghi, E. Björnson, E. G. Larsson, C. Yuen, and T. Marzetta, "Joint unicast and multi-group multicast transmission in massive MIMO systems," *IEEE Trans. Wireless Commun.*, vol. 17, no. 10, pp. 6375–6388, Oct. 2018.
- [29] S. Mohammadi, M. Dong, and S. ShahbazPanahi, "Fast algorithm for joint unicast and multicast beamforming in large-scale systems," in *Proc. IEEE Workshop on Signal Processing Advances in Wireless Commun.(SPAWC)*, Sept. 2021, pp. 91–95.

- [30] H. Joudeh and B. Clerckx, "Rate-splitting for max-min fair multigroup multicast beamforming in overloaded systems," *IEEE Trans. Wireless Commun.*, vol. 16, no. 11, pp. 7276–7289, Nov. 2017.
- [31] O. Tervo, L.-N. Tran, S. Chatzinotas, B. Ottersten, and M. Juntti, "Multigroup multicast beamforming and antenna selection with rate-splitting in multicell systems," in *Proc. IEEE Int. Workshop Signal Process. Advances Wireless Commun.*, Jun. 2018, pp. 1–5.
- [32] H. Chen, D. Mi, B. Clerckx, Z. Chu, J. Shi, and P. Xiao, "Joint power and subcarrier allocation optimization for multigroup multicast systems with rate splitting," *IEEE Trans. Veh. Technol.*, vol. 69, no. 2, pp. 2306–2310, Feb. 2020.
- [33] G. M. Korpelevich, "The extragradient method for finding saddle points and other problems," *Matecon*, vol. 13, pp. 35–49, 1977.
- [34] F. Facchinei and J.-S. Pang, *Finite-Dimensional Variational Inequalities* and Complementarity Problems, New York, NY: Springer-Verlag, 2003.
- [35] E. N. Khobotov, "Modification of the extra-gradient method for solving variational inequalities and certain optimization problems," USSR Comput. Math. Math. Phys., vol. 27, no. 5, pp. 120–127, 1987.
- [36] P. Marcotte, "Application of Khobotov's algorithm to varational inequalities and network equilibrium problems," *Inf. Syst. Oper. Res.*, vol. 29, no. 4, pp. 258–270, Nov. 1991.
- [37] 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, 2011.