# Max-Min Fair Precoder Design and Power Allocation for MU-MIMO NOMA

Ahmet Zahid Yalcin, Mustafa Kagan Cetin and Melda Yuksel *IEEE, Senior Member*

*Abstract*—In this paper, a downlink multiple input multiple output (MIMO) non-orthogonal multiple access (NOMA) wireless communication system is considered. In NOMA systems, the base station has unicast data for all users, and multiple users in a group share the same resources. The objective is to design transmit precoders and power allocation coefficients jointly that provide max-min fairness (MMF) among the strongest users in each group, while maintaining minimum target rates for all the other users. The problem is solved via two main iterative approaches. The first method is based on semi-definite relaxation (SDR) and successive convex approximation (SCA), and the second method is based on the equivalency between achievable rate and minimum mean square error (MMSE) expressions. For the latter approach, Karush-Kuhn-Tucker (KKT) optimality conditions are derived and the expressions satisfied by the optimal receivers, MMSE weights and the optimal precoders are obtained. Proposed algorithms are compared with rate-splitting (RS), orthogonal multiple access (OMA) and multi-user linear precoding (MULP) schemes in terms of MMF rates, energy efficiency and complexity. It is shown that while RS has the best MMF rates and energy efficiency, the MMSE approach based on KKT optimality conditions has the least complexity. Moreover, the SDR/SCA approach offers an excellent tradeoff. It offers high MMF rates, low complexity and superior energy efficiency.

*Index Terms*—Max-min fairness, mean square error, MIMO, NOMA, precoder design, quality-of-service, rate splitting, successive convex approximation.

## I. INTRODUCTION

The demand for data traffic is steadily increasing and wireless networks of the next decade have to meet the high data rate requirements for many different applications [1]. To handle this high data rate, non-orthogonal multiple access (NOMA) is considered as a breakthrough technique, which enables simultaneous multiple access in the power domain for 5G wireless networks [2]. Specifically, downlink NOMA is an application of broadcast channels [3] and it relies on superposition coding (SPC) at the transmitter to transfer multiple data streams in the same resource block, and successive interference cancellation (SIC) at the receiver to cancel co-channel interference. NOMA has the potential to deliver higher system throughput [4], [5] and higher ergodic sum capacity [6], and to achieve better outage performance [7] compared to the existing orthogonal multiple access (OMA) techniques. In practical power domain NOMA schemes, more power is allocated to users with poor channel conditions to guarantee

their required minimum rates [8]. This way NOMA presents an advantage in providing higher spectral efficiency and fairness.

### A. Related Work

In NOMA systems, each user can have a dedicated precoding vector, or a cluster of users can share the same precoding vector. The former has the advantage of custom precoding for each user, but suffers from the rank constraints in the downlink multiuser MIMO broadcast channel [9]. The latter is not limited by rank, but messages are not individually precoded, so channel gain vectors and precoders are mismatched.

Assuming the transmit signals of each user are coded by a dedicated precoding vector, sum rate maximization, total power minimization and max-min fairness for NOMA systems are studied under different constraints and with different methods in the literature. The paper [9] solves the sum rate maximization problem by approximating the problem with a minorization-maximization algorithm. The paper [10] presents a precoding design for maximizing the sum rate of all users under decoding order and quality-of-service (QoS) constraints. Similarly, to maximize sum rate, [11] studies the channel state information based singular value decomposition precoding scheme. Total power minimization with QoS requirements and total power minimization under target interference level constraints are respectively investigated in [12] and [13]. In addition, a max-min fair (MMF) precoder design problem for a multiple antenna base station is also studied in [12]. Power allocation (PA) problems for achieving MMF in NOMA systems with single antenna transmitters are studied in [14] and [15].

As mentioned above, in NOMA, a single precoder vector can be shared by a cluster of users. For this case, weighted sum rate optimization under a total power constraint when two users exist in each cluster is studied in [16]. For clustered downlink NOMA systems, a sub-optimal user clustering algorithm is proposed and the optimal power allocation policy that maximizes the weighted sum rate is derived in [17], [18]. Joint power allocation and precoder design to maximize the strong users' sum rate subject to QoS constraints on weak users' rates is solved via successive convex approximation (SCA) and semi-definite relaxation (SDR) in [19]. The same problem is generalized to the multi-cell networks in [20]. Finally, minimizing total transmission power for downlink clustered NOMA is studied in [21] and [22].

### B. Motivation and Contributions

In this work, we study downlink MIMO clustered NOMA system from a fairness standpoint and we investigate joint

A. Z. Yalcin, M. K. Cetin and M. Yuksel are with TOBB University of Economics and Technology, Ankara 06560, Turkey (email: {azyalcin, mustafakagancetin, yuksel}@etu.edu.tr).

precoder design and PA problem that provide MMF among the strongest users in each cluster and ensure the minimum rate requirements for all the other users. To the best of our knowledge, there is no joint MMF precoder design and PA optimization for a clustered downlink NOMA system. Our contributions are listed below:

1) Firstly, we define a joint precoder design and PA problem to attain max-min fairness among the best users in each cluster, while guaranteeing target data rates for the rest of the users. Due to the non-convexity of the defined problem, we apply Taylor series expansion, SDR, to simplify the original problem. Next, we propose a suboptimal iterative SCA based algorithm.

2) Secondly, we use the equivalency between weighted mean square error (WMMSE) and achievable rate expressions, and restate the original problem as an equivalent MMF WMMSE problem. To do that, we apply the achievable rate-WMMSE relationship. Due to the non-convexity of the main problem, we split it into two different problems: i) to design optimal precoders for given power allocation coefficients (PAC), and ii) to obtain optimal PAC for given precoders. We derive a sub-optimal PA scheme while designing the optimal precoders. Employing the CVX toolbox to obtain the precoders, we then propose a suboptimal iterative WMMSE based algorithm, which updates transmit precoders, receivers, weights and PAC sequentially.

3) Thirdly, employing the Karush-Kuhn-Tucker (KKT) optimality conditions, we find the expressions the optimal receivers, MMSE weights and the optimal precoders have to satisfy. Utilizing these expressions, we propose a low-complexity iterative algorithm to evaluate precoders and receivers. We use the exponential penalty method to evaluate the Lagrange multipliers. We find that this approach significantly decreases complexity, while ensuring a similar MMF rate performance as the CVX solution in the second item.

4) To the best of our knowledge, there is no work in the literature, which studies both SDR/SCA and WMMSE based approaches for the same optimization problem. We discover that SDR/SCA performs better than the latter as it solves a tighter approximation.

5) We compare the proposed schemes with and without power allocation to observe that power optimization does not significantly increase complexity and its advantages in terms of MMF rates are justified.

6) We also compare our results with rate-splitting (RS) [23], [24], OMA and multi-user linear precoding (MULP) schemes in terms of MMF rates, complexity and energy efficiency. Our results reveal that the SDR/SCA based scheme offers an excellent tradeoff in all three aspects.

Next, we explain the system model and define the optimization problem in Section II. We propose the SDR/SCA and WMMSE based precoder designs respectively in Sections III and IV. We present the numerical results in Section V. Finally we provide conclusions and future work in Section VI.

## II. System Model and Problem Definitions

In this paper, we investigate a downlink multiuser MIMO system. The base station has $M$ transmit antennas and communicates with $K$ clusters. There are $L$ single antenna users in each cluster[1] and each user belongs to only one cluster.

The base station aims to send the data $s_{k,l}$ to the $l$-th user in the $k$-th cluster, for all $l \in \{1, \ldots, L\}$, and $k \in \{1, \ldots, K\}$. All $s_{k,l}$ are independent and $\mathbb{E}\{s_{k,l}s_{k,l}^*\} = \alpha_{k,l}$. Here $\alpha_{k,l}$ is the ratio of power allocated to the data stream $s_{k,l}$. The PAC vector is defined as $\mathbf{A} = [\alpha_{1,1}, \ldots, \alpha_{1,L}, \ldots, \alpha_{K,1}, \ldots, \alpha_{K,L}]$. Moreover, $\sum_{l=1}^{L} \alpha_{k,l} = 1$. To send all the messages, the base station superposes all the messages in a cluster as $s_k = \sum_{l=1}^{L} s_{k,l}$ and forms $\mathbf{s} = [s_1, \ldots, s_K]^T \in \mathbb{C}^{K \times 1}$. When $\mathbf{p}_k \in \mathbb{C}^{M \times 1}$ indicates the precoder vector for the $k$-th cluster, the base station transforms $\mathbf{s}$ with the precoder matrix $\mathbf{P} = [\mathbf{p}_1, \ldots, \mathbf{p}_K] \in \mathbb{C}^{M \times K}$. Then, the base station transmits $\mathbf{x} \in \mathbb{C}^{M \times 1}$, which is equal to

$$\mathbf{x} = \mathbf{P}\mathbf{s} = \sum_{k=1}^{K} \mathbf{p}_k s_k = \sum_{k=1}^{K} \sum_{l=1}^{L} \mathbf{p}_k s_{k,l}. \tag{1}$$

The base station has an average total power constraint $E_{tx}$, which is written as

$$\mathbb{E}\{\mathbf{x}^H\mathbf{x}\} = \text{Tr}(\mathbf{P}\mathbf{P}^H) \leq E_{tx}. \tag{2}$$

Then, the received signal at the $l$-th user in the $k$-th cluster becomes

$$y_{k,l} = \mathbf{h}_{k,l}^H \mathbf{p}_k \sum_{l=1}^{L} s_{k,l} + \mathbf{h}_{k,l}^H \sum_{i=1,i\neq k}^{K} \mathbf{p}_i s_i + n_{k,l}. \tag{3}$$

Here, $\mathbf{h}_{k,l} \in \mathbb{C}^{M \times 1}$ is the effective channel gain vector of the $l$-th user in the $k$-th cluster. The effective channel gain is defined as $\mathbf{h}_{k,l} = \tilde{\mathbf{h}}_{k,l}/\sqrt{d_{k,l}^\rho}$, where $d_{k,l}$ is the distance between the $l$-th user in the $k$-th cluster and the base station, and $\rho$ is the path loss exponent. The entries in $\tilde{\mathbf{h}}_{k,l}$ are independent and identically distributed (i.i.d.) and complex valued random variables. Moreover, the effective channel gain magnitudes are ordered as $|h_{k,L}| > |h_{k,L-1}| > \ldots > |h_{k,1}|$. It means that the user with the smallest effective channel gain magnitude is the first user in a cluster and the $L$-th user has the largest channel gain magnitude. The noise component $n_{k,l}$ is a circularly symmetric complex Gaussian random variable with zero mean and unit variance, and $n_{k,l}$ are i.i.d. for all $k$ and $l$. The base station is informed about all effective channel gains $\mathbf{h}_{k,l}$, while the receivers know only their own $\mathbf{h}_{k,l}$.

### A. Achievable Data Rates

For this NOMA system we investigate, the messages for different clusters will be treated as noise, while SIC will be carried out within a cluster to limit intra-cluster interference. Due to SIC, in the $k$-th cluster, the $l$-th user's message is decoded at the $i$-th user, for which $l \leq i$. In other words, decoding is ordered and starts from the first user's message.

---

[1]In fact, the results can easily be extended to cover for unequal number of users in each group. However, to keep the notation simple we adhere to a fixed number of users in each cluster.

The first user in the cluster decodes its own message only, and the $L$-th user decodes all users' messages within the cluster. To simplify the notation, we define the sets $\mathcal{K} \triangleq \{1, ..., K\}$, $\mathcal{L} \triangleq \{1, ..., L\}$, $\bar{\mathcal{L}} \triangleq \{1, ..., L-1\}$ and $\mathcal{I} \triangleq \{l, ..., L\}$. Then, the signal to interference ratio (SINR) for decoding the $l$-th user's message at the $i$-th user in the $k$-th cluster, $i \in \mathcal{I}, l \in \mathcal{L}, k \in \mathcal{K}$ can be written as

$$\gamma_{k,i \to l} = \alpha_{k,l} |\mathbf{h}_{k,i}^H \mathbf{p}_k|^2 r_{k,i \to l}^{-1}. \tag{4}$$

In the above equation, $r_{k,i \to l}$ is the effective noise variance and is defined as

$$r_{k,i \to l} = \sum_{j=l+1}^{L} \alpha_{k,j} |\mathbf{h}_{k,i}^H \mathbf{p}_k|^2 + \sum_{t=1, t \neq k}^{K} |\mathbf{h}_{k,i}^H \mathbf{p}_t|^2 + 1. \tag{5}$$

Then, in the $k$-th cluster, the $i$-th user's achievable rate[2] for decoding the $l$-th user's message is

$$R_{k,i \to l} = \log\left(1 + \gamma_{k,i \to l}\right). \tag{6}$$

Overall, the achievable rate for the $l$-th user's message in the $k$-th cluster is defined as the minimum of all $R_{k,i \to l}$, and is denoted as

$$R_{k,l} = \min_{i, i \in \mathcal{I}} R_{k,i \to l}, \forall l \in \bar{\mathcal{L}}. \tag{7}$$

Note that, due to this definition, $R_{k,L} = R_{k,L \to L}$.

### B. Max-Min Fair Problem Definition

In this subsection, we define the MMF rate optimization problem, which aims to find the optimal precoder matrix $\mathbf{P}$ and optimal PAC vector $\mathbf{A}$, such that the minimum of the strongest users' rates is maximized subject to a total power constraint and a minimum rate constraint for the rest of the users. Then, the optimization problem is stated as

$$\max_{\mathbf{P}, \mathbf{A}} \quad \min_{k \in \mathcal{K}} R_{k,L} \tag{8a}$$

$$\text{s.t.} \quad R_{k,l}^{th} \leq R_{k,i \to l}, \forall k \in \mathcal{K}, \forall l \in \bar{\mathcal{L}}, \forall i \in \mathcal{I}, \tag{8b}$$

$$\sum_{l=1}^{L} \alpha_{k,l} = 1, \alpha_{k,l} \geq 0, \ \forall k \in \mathcal{K}, \forall l \in \mathcal{L}, \tag{8c}$$

$$\text{Tr}(\mathbf{P}\mathbf{P}^H) \leq E_{tx}, \tag{8d}$$

where $R_{k,l}^{th} \geq 0$ is the threshold data rate that has to be provided to the $l$-th user in the $k$-th cluster $\forall k \in \mathcal{K}, \forall l \in \bar{\mathcal{L}}$. Note that, due to SIC, the $l$-th user's message in the $k$-th cluster has to be decoded by all $i$, $\forall i \in \mathcal{I}$, resulting in the inequality in (8b). The equality in (8c) indicates that the superposed data $s_k$ for the $k$-th cluster has normalized power. In addition, (8d) is the total power constraint at the base station.

To solve this problem, we need to restate (8), as the minimum operation in the objective function is not a convex

[2]In all the derivations, all rate expressions are expressed in nats/channel use. In Section V, without loss of generality, simulation results are presented in bits/channel use.

function. Thus, we add an auxiliary variable $R_g$ and convert (8) to a new constrained optimization problem as

$$\max_{\mathbf{P}, \mathbf{A}, R_g} \quad R_g \tag{9a}$$

$$\text{s.t.} \quad R_g \leq R_{k,L}, \forall k \in \mathcal{K}, \tag{9b}$$

$$(8b), (8c), (8d). \tag{9c}$$

## III. Successive Convex Approximation Solution

The problem defined in (9) is still a non-convex optimization problem. In this section, we further modify the optimization problem in (9) to obtain an equivalent semi-definite programming problem.

To achieve this objective, we introduce and optimize the auxiliary optimization matrix $\mathbf{Q}_k = \mathbf{p}_k \mathbf{p}_k^H$. Note that, $\mathbf{Q}_k \in \mathbb{C}^{M \times M}$ is a rank-one positive semi-definite matrix. Then, we can rewrite our optimization problem as

$$\max_{\substack{\mathbf{Q}, \mathbf{A}, \\ R_g}} \quad R_g \tag{10a}$$

$$\text{s.t.} \quad R_g \leq \log\left(1 + \tilde{\gamma}_{k,L}\right), \forall k \in \mathcal{K}, \tag{10b}$$

$$R_{k,l}^{th} \leq \log\left(1 + \tilde{\gamma}_{k,i \to l}\right), \forall k \in \mathcal{K}, \forall l \in \bar{\mathcal{L}}, \forall i \in \mathcal{I}, \tag{10c}$$

$$\sum_{l=1}^{L} \alpha_{k,l} = 1, \alpha_{k,l} \geq 0, \ \forall k \in \mathcal{K}, \forall l \in \mathcal{L}, \tag{10d}$$

$$\mathbf{Q}_k \succeq 0, \forall k \in \mathcal{K}, \tag{10e}$$

$$\text{rank}\left(\mathbf{Q}_k\right) \leq 1, \forall k \in \mathcal{K}, \tag{10f}$$

$$\sum_{k=1}^{K} \text{Tr}\left(\mathbf{Q}_k\right) \leq E_{tx}, \tag{10g}$$

where

$$\tilde{\gamma}_{k,L} = \frac{\alpha_{k,L} \mathbf{h}_{k,L}^H \mathbf{Q}_k \mathbf{h}_{k,L}}{\sum_{t=1, t \neq k}^{K} \mathbf{h}_{k,L}^H \mathbf{Q}_t \mathbf{h}_{k,L} + 1}, \tag{11}$$

$$\tilde{\gamma}_{k,i \to l} = \frac{\alpha_{k,l} \mathbf{h}_{k,i}^H \mathbf{Q}_k \mathbf{h}_{k,i}}{\sum_{j=l+1}^{L} \alpha_{k,j} \mathbf{h}_{k,i}^H \mathbf{Q}_k \mathbf{h}_{k,i} + \sum_{t=1, t \neq k}^{K} \mathbf{h}_{k,i}^H \mathbf{Q}_t \mathbf{h}_{k,i} + 1}. \tag{12}$$

Convex optimization solvers are not efficient when operating with logarithmic functions. To eliminate the logarithms in (10b) and (10c), we define a new auxiliary variable $\delta$ and new constants $\zeta_{k,l}$, $\forall k \in \mathcal{K}, \forall l \in \bar{\mathcal{L}}$ as

$$\delta = e^{R_g} - 1$$
$$\zeta_{k,l} = e^{R_{k,l}^{th}} - 1.$$

Then, we can reformulate (10) as

$$\max_{\mathbf{Q}, \mathbf{A}, \delta} \quad \delta \tag{13a}$$

$$\text{s.t.} \quad \delta \leq \frac{\alpha_{k,L} \phi_{k,L}}{\omega_{k,L}}, \ \forall k \in \mathcal{K}, \tag{13b}$$

$$\zeta_{k,l} \leq \frac{\mu_{k,l} \phi_{k,i}}{\omega_{k,i}} \ \forall k \in \mathcal{K}, \forall l \in \bar{\mathcal{L}}, \forall i \in \mathcal{I}, \tag{13c}$$

$$(10d), (10e), (10f), (10g), \tag{13d}$$

where

$$\phi_{k,l} = \mathbf{h}_{k,l}^H \mathbf{Q}_k \mathbf{h}_{k,l} \tag{14}$$

$$\omega_{k,i} = 1 + \sum_{t=1, t \neq k}^{K} \mathbf{h}_{k,i}^H \mathbf{Q}_t \mathbf{h}_{k,i} \tag{15}$$

$$\mu_{k,l} = \alpha_{k,l} - \zeta_{k,l} \sum_{j=l+1}^{L} \alpha_{k,j}. \tag{16}$$

The constraints (13b) and (13c) are not convex, since $\alpha_{k,L}\phi_{k,L}$ and $\mu_{k,l}\phi_{k,i}$ are both bilinear functions. To change (13b) and (13c) into convex constraints, we need to apply the Schur complement [25]. Introducing new auxiliary variables $\tau_{k,i,l}$, we can replace (13b) and (13c) with the following 4 new constraints

$$\begin{bmatrix} \alpha_{k,L} & \tau_{k,L,L} \\ \tau_{k,L,L} & \phi_{k,L} \end{bmatrix} \succeq \mathbf{0}, \ \forall k \in \mathcal{K} \tag{17}$$

$$\begin{bmatrix} \mu_{k,l} & \tau_{k,i,l} \\ \tau_{k,i,l} & \phi_{k,i} \end{bmatrix} \succeq \mathbf{0}, \ \forall k \in \mathcal{K}, \forall l \in \bar{\mathcal{L}}, \forall i \in \mathcal{I} \tag{18}$$

and

$$\delta \leq \frac{\tau_{k,L,L}^2}{\omega_{k,L}}, \ \forall k \in \mathcal{K} \tag{19}$$

$$\zeta_{k,l} \leq \frac{\tau_{k,i,l}^2}{\omega_{k,i}}, \ \forall k \in \mathcal{K}, \ \forall l \in \bar{\mathcal{L}}, \ \forall i \in \mathcal{I}. \tag{20}$$

The right-hand side of (19) is convex in both $\tau_{k,L,L}$ and $\omega_{k,L}$, and the right-hand side of (20) is convex in both $\tau_{k,i,l}$ and $\omega_{k,i}$. In other words, right hand sides of both (19) and (20) are difference-of-convex functions [26]. Therefore, we can apply the first-order Taylor expansions [27] to obtain a tight lower bound on these two functions. For given fixed points $(\tilde{\tau}_{k,L,L}, \tilde{\omega}_{k,L})$ $\forall k \in \mathcal{K}$ and $(\tilde{\tau}_{k,i,l}, \tilde{\omega}_{k,i}), \forall k \in \mathcal{K}, \forall l \in \bar{\mathcal{L}}, \forall i \in \mathcal{I}$, we write

$$\delta \leq \frac{2\tilde{\tau}_{k,L,L}}{\tilde{\omega}_{k,L}} \tau_{k,L,L} - \frac{\tilde{\tau}_{k,L,L}^2}{\tilde{\omega}_{k,L}^2} \omega_{k,L} \leq \frac{\tau_{k,L,L}^2}{\omega_{k,L}} \tag{21}$$

$$\zeta_{k,l} \leq \frac{2\tilde{\tau}_{k,i,l}}{\tilde{\omega}_{k,i}} \tau_{k,i,l} - \frac{\tilde{\tau}_{k,i,l}^2}{\tilde{\omega}_{k,i}^2} \omega_{k,i} \leq \frac{\tau_{k,i,l}^2}{\omega_{k,i}} \tag{22}$$

where $\tilde{\tau}_{k,i,l} \geq 0$ and $\tilde{\omega}_{k,i} \geq 1$.

Finally, we relax the equality in (10d) as an inequality, omit the constraint in (10f) and transform the optimization defined in (10) as

---

**Algorithm 1** SDR/SCA Based MMF Algorithm with PA

1: **Input:** $\mathbf{A}^{(0)}$, $E_{tx}$, $\Upsilon$, $R_{k,l}^{th}$, $n_{max}$
2: **Initialize:** $\delta^{(0)} = 0$, $\tilde{\omega}_{k,i}^{(0)}, \tilde{\tau}_{k,i,l}^{(0)}$, and $n = 0$;
3: **iterate** $\forall j, l, k$;
4:     $n = n + 1$
5:     Update $\left\{ \mathbf{Q}_k^{(n)}, \mathbf{A}^{(n)}, \delta^{(n)}, \tau_{k,i,l}^{(n)} \right\}$ by solving (23) for given $\tilde{\tau}_{k,i,l}^{(n-1)}$ and $\tilde{\omega}_{k,i}^{(n-1)}$
6:     Update $\tilde{\tau}_{k,i,l}^{(n)}$ and $\tilde{\omega}_{k,i}^{(n)}$ using (25)
7:     **If** $(\delta^{(n)} - \delta^{(n-1)})/\delta^{(n-1)} < \Upsilon$ **or** $n = n_{max}$ **then**
8:       Terminate
9:     **else then**
10:       Go to Step 3

---

$$\max_{\mathbf{Q}, \mathbf{A}, \delta, \tau} \quad \delta \tag{23a}$$

$$\text{s.t.} \quad \delta \leq \frac{2\tilde{\tau}_{k,L,L}}{\tilde{\omega}_{k,L}} \tau_{k,L,L} - \frac{\tilde{\tau}_{k,L,L}^2}{\tilde{\omega}_{k,L}^2} \omega_{k,L}, \forall k \in \mathcal{K}, \tag{23b}$$

$$\zeta_{k,l} \leq \frac{2\tilde{\tau}_{k,i,l}}{\tilde{\omega}_{k,i}} \tau_{k,i,l} - \frac{\tilde{\tau}_{k,i,l}^2}{\tilde{\omega}_{k,i}^2} \omega_{k,i},$$
$$\forall k \in \mathcal{K}, \forall l \in \bar{\mathcal{L}}, \forall i \in \mathcal{I}, \tag{23c}$$

$$\begin{bmatrix} \alpha_{k,L} & \tau_{k,L,L} \\ \tau_{k,L,L} & \phi_{k,L} \end{bmatrix} \succeq \mathbf{0}, \forall k \in \mathcal{K}, \tag{23d}$$

$$\begin{bmatrix} \mu_{k,l} & \tau_{k,i,l} \\ \tau_{k,i,l} & \phi_{k,i} \end{bmatrix} \succeq \mathbf{0}, \forall k \in \mathcal{K}, \forall l \in \bar{\mathcal{L}}, \forall i \in \mathcal{I}, \tag{23e}$$

$$\sum_{l=1}^{L} \alpha_{k,l} = 1, \alpha_{k,l} \geq 0, \ \forall k \in \mathcal{K}, \forall l \in \mathcal{L}, \tag{23f}$$

$$\mathbf{Q}_k \succeq 0, \forall k \in \mathcal{K}, \tag{23g}$$

$$\sum_{k=1}^{K} \text{Tr}(\mathbf{Q}_k) \leq E_{tx}. \tag{23h}$$

The problem in (23) is a constrained convex optimization problem when $\tilde{\tau}_{k,i,l}$ and $\tilde{\omega}_{k,i}$ are given. In [20], the authors prove that omitting the rank constraint (10f) in (23) does not alter the problem. They discuss that the solution is always rank one. However, their proof assumes that the principle eigenvalue of the positive semi-definite matrix in [20, eqn. (33)] is always unique. This may not be the case and there can be more than one principle eigenvalue. However, adding independent rank one matrices in [20, eqn. (33)] results in full rank matrices with very high probability and this does not pose a significant issue.

The optimization problem (23) is an approximation to the original problem in (8). To solve (8), we use Algorithm 1. The algorithm solves (23) when $\tilde{\tau}_{k,i,l}$ and $\tilde{\omega}_{k,i}$ are given, and updates these values in each iteration. While solving (23), we employ the CVX optimization toolbox [28].

In Algorithm 1, we can initialize $\tilde{\tau}_{k,i,l}^{(0)}$ and $\tilde{\omega}_{k,i}^{(0)}$ arbitrarily, as long as $\tilde{\tau}_{k,i,l}^{(0)} \geq 0$ and $\tilde{\omega}_{k,i}^{(0)} \geq 1$. However, one can initialize $\tilde{\tau}_{k,i,l}^{(0)}$ and $\tilde{\omega}_{k,i}^{(0)}$ more efficiently. To do so, we create a random rank one positive semi-definite matrix for each $\mathbf{Q}_k$ and a uniform vector $\mathbf{A}$, calculate $\phi_{k,l}^{(0)}$ and $\omega_{k,i}^{(0)}$ and compute $\tilde{\tau}_{k,i,l}^{(0)}$

and $\tilde{\omega}_{k,i}^{(0)}$ as

$$\tilde{\tau}_{k,i,l}^{(0)} = \begin{cases} \sqrt{\alpha_{k,L}\phi_{k,L}^{(0)}}, & i = l = L \\ \sqrt{\mu_{k,l}\phi_{k,i}^{(0)}}, & l \leq i, l < L \end{cases} \quad \text{and} \quad \tilde{\omega}_{k,i}^{(0)} = \omega_{k,i}^{(0)} \tag{24}$$

In (24), the $\tilde{\tau}_{k,i,l}^{(0)}$ values satisfy (23d) and (23e) with equality. Using the randomly generated matrices for $\mathbf{Q}_k$, we calculate $\tilde{\omega}_{k,i}^{(0)}$ using (15). After initialization, in Algorithm 1, we update $\tilde{\tau}_{k,i,l}$ and $\tilde{\omega}_{k,i}$ in each iteration as

$$\tilde{\tau}_{k,i,l}^{(n)} = \tau_{k,i,l}^{(n-1)} \quad \text{and} \quad \tilde{\omega}_{k,i}^{(n)} = \omega_{k,i}^{(n-1)}. \tag{25}$$

This way, the bounds in (23b) and (23c) become tighter in each iteration. The algorithm convergence can be proved in a similar manner as in [20].

## IV. WMMSE BASED SOLUTIONS

In this section, we provide an alternative solution to the MMF problem defined in (8) using the MMSE approach. In this approach, we utilize the relation between mutual information and MMSE [29], [30]. We can state $R_{k,i \to l}$ in terms of error variances, assuming MMSE receivers are employed at the receivers.

We remind that the effective channel gain magnitudes are ordered as $|h_{k,L}| > |h_{k,L-1}| > \ldots |h_k, 1|$ as described in Section II. Therefore, the $l$-th user in the $k$-th cluster decodes messages in order starting from the first user's message, and decodes its own message in the last step. SIC is employed in each step. In other words, to estimate the $l$-th user's message, the $i$-th user ($i \in \mathcal{I}$) in the $k$-th cluster employs the SIC receiver $V_{k,i \to l}$ on its equivalent received signal $\hat{y}_{k,i}$ where

$$\hat{y}_{k,i} = y_{k,i} - \mathbf{h}_{k,i}^H \mathbf{p}_k \sum_{j=1}^{l-1} s_{k,j}. \tag{26}$$

The $i$-th user's estimate $\hat{s}_{k,i \to l}$ about $s_{k,l}$ becomes

$$\hat{s}_{k,i \to l} = V_{k,i \to l} \hat{y}_{k,i}. \tag{27}$$

Then, the MSE of the $i$-th user's estimate of the $l$-th user's message in the $k$-th cluster can be written as

$$\varepsilon_{k,i \to l} = \mathbb{E}\{||\hat{s}_{k,i \to l} - s_{k,l}||^2\}$$
$$= |V_{k,i \to l}|^2 T_{k,i \to l} + \alpha_{k,l} - 2\Re\{\alpha_{k,l} V_{k,i \to l} \mathbf{h}_{k,i}^H \mathbf{p}_k\}, \tag{28}$$

where $T_{k,i \to l} = |_{k,i}\mathbf{p}_k|^2 \alpha_{k,l} + r_{k,i \to l}$. Given above, the optimal MMSE receiver is

$$V_{k,i \to l}^{mmse} = \arg \min_{V_{k,i \to l}} \varepsilon_{k,i \to l} = \alpha_{k,l} \mathbf{p}_k^H \mathbf{h}_{k,i} T_{k,i \to l}^{-1}. \tag{29}$$

When this MMSE receiver in (29) is employed, the resulting error variance expression in (28) becomes

$$\varepsilon_{k,i \to l}^{mmse} = \left(\frac{1}{\alpha_{k,l}} + |\mathbf{h}_{k,i}^H \mathbf{p}_k|^2 r_{k,i \to l}^{-1}\right)^{-1}. \tag{30}$$

As the message for the $l$-th user has to be decoded by all users $i$ for which $i \geq l$ in the $k$-th cluster, we define $\varepsilon_{k,l}^{mmse}$ as

$$\varepsilon_{k,l}^{mmse} = \max_{i, i \in \{l, \ldots, L\}} \varepsilon_{k,i \to l}^{mmse}. \tag{31}$$

Note that, by simply comparing the rate and MMSE expressions in (6) and (30) we observe that

$$R_{k,i \to l} = \log\left[\alpha_{k,l} \varepsilon_{k,i \to l}^{mmse^{-1}}\right]. \tag{32}$$

### A. Equivalent MMF WMMSE Problem

To convert (9) into an equivalent WMMSE problem, we use the above relation between rate and MMSE. We define the augmented weighted MSE [30] as

$$\xi_{k,i \to l} = b_{k,i \to l} \varepsilon_{k,i \to l} - \log(\alpha_{k,l} b_{k,i \to l}), \tag{33}$$

where $b_{k,i \to l} > 0$ is the weight for MSE. We also define the minimum of the augmented WMSEs as

$$\xi_{k,i \to l}^{mmse} \triangleq \arg \min_{\{b_{k,i \to l}, V_{k,i \to l}\}} \xi_{k,i \to l}, \tag{34}$$

$$= b_{k,i \to l}^{mmse} \varepsilon_{k,i \to l}^{mmse} - \log(\alpha_{k,l} b_{k,i \to l}^{mmse}). \tag{35}$$

It is seen that the augmented WMSE $\xi_{k,i \to l}$ is convex in the receiver $V_{k,i \to l}$. Solving for the first order optimality conditions in (33), we find the optimum receiver in (35) as $V_{k,i \to l}^{\star} = V_{k,i \to l}^{mmse}$ and the optimum weights as

$$b_{k,i \to l}^{\star} = b_{k,i \to l}^{mmse} = \frac{1}{\varepsilon_{k,i \to l}^{mmse}}, \tag{36}$$

where the MMSE receiver $V_{k,i \to l}^{mmse}$ is given in (29) and the MMSE error variance $\varepsilon_{k,i \to l}^{mmse}$ is given in (30).

One can obtain the relation between rate expressions and augmented WMSEs by checking the first order optimality conditions [30] to find that

$$\xi_{k,i \to l}^{mmse} = 1 - R_{k,i \to l}. \tag{37}$$

Utilizing the equality in (37), the optimization problem in (9) can be written as:

$$\max_{\substack{\mathbf{P}, \mathbf{A}, \\ \mathbf{R}, R_g}} \quad R_g \tag{38a}$$

$$\text{s.t.} \quad R_g \leq R_{k,L}, \forall k \in \mathcal{K}, \tag{38b}$$

$$R_{k,l}^{th} \leq R_{k,l}, \forall k \in \mathcal{K}, \forall l \in \bar{\mathcal{L}}, \tag{38c}$$

$$R_{k,l} \leq 1 - \xi_{k,i \to l}^{mmse}, \forall k \in \mathcal{K}, \forall l \in \mathcal{L}, \forall i \in \mathcal{I}, \tag{38d}$$

$$(8c), (8d), \tag{38e}$$

where $\mathbf{R} = [R_{1,1}, \ldots, R_{1,L}, \ldots, R_{K,1}, \ldots, R_{K,L}]$ is a new auxiliary variable vector.

The optimization problem in (38) assumes that the optimal MMSE receiver defined in (29) is employed at all users, and finds the optimal precoders at the transmitter. Below, we first define a generalized problem which allows for arbitrary receivers $V_{k,i \to l}$ that attain $\varepsilon_{k,i \to l}$ in (28).

$$\max_{\substack{\mathbf{P}, \mathbf{A}, \mathbf{R}, \\ R_g, \mathbf{V}}} \quad R_g \tag{39a}$$

$$\text{s.t.} \quad R_g \leq R_{k,L}, \forall k \in \mathcal{K}, \tag{39b}$$

$$R_{k,l}^{th} \leq R_{k,l}, \forall k \in \mathcal{K}, \forall l \in \bar{\mathcal{L}}, \tag{39c}$$

$$R_{k,l} \leq 1 - \xi_{k,i \to l}, \forall k \in \mathcal{K}, \forall l \in \mathcal{L}, \forall i \in \mathcal{I}, \tag{39d}$$

$$(8c), (8d). \tag{39e}$$

The problem defined in (39) is hard to solve and there are no closed form expressions for the optimal precoders and PAC. Instead, in this section we propose an iterative precoder design algorithm. To do that, we need to split this problem into two different problems. In the first part of (39), we investigate the optimal precoders for a given set of PAC. Then, we update PAC using the updated precoders. For the first part we assume $\mathbf{A}$ is given and solve

$$\max_{\substack{\mathbf{P},\mathbf{R}, \\ R_g,\mathbf{V}}} \quad R_g \tag{40a}$$

$$\text{s.t.} \quad R_{k,l} \leq 1 - \xi_{k,i\to l}, \forall k \in \mathcal{K}, \forall l \in \mathcal{L}, \forall i \in \mathcal{I}, \tag{40b}$$

$$(39b), (39c), (8c), (8d). \tag{40c}$$

where $\mathbf{V} = [V_{1,1\to 1}, \ldots, V_{K,L\to L}]$ and $\mathbf{b} = [b_{1,1\to 1}, \ldots, b_{K,L\to L}]$ consist of all receivers and weights respectively. The optimization problem in (40) is convex if either the precoder matrix $\mathbf{P}$ or the receiver matrix $\mathbf{V}$ is given. Thus, an iterative algorithm can solve (40) sub-optimally, starting from an initial precoder $\mathbf{P}^{init}$. In each iteration, the algorithm can update $\mathbf{V}$ and $\mathbf{P}$ using a standard convex program solver such as CVX [28].

### B. Power Allocation for WMMSE

In this subsection, we discuss the optimal PAC selection for the second part of (39) for given precoders and receivers, and weights $b_{k,i\to l}$, which are already calculated using (36).

Note that, PAC in each cluster are not related with the coefficients in other clusters, as inter-cluster power allocation is already a part of the precoder optimization in (40). Therefore, we can consider PAC optimization as intra-cluster power allocation and write a simplified problem for each cluster $k$ as

$$\max_{\mathbf{A}} \quad R_{k,L} \tag{41a}$$

$$\text{s.t.} \quad R_{k,l}^{th} \leq R_{k,i\to l}, \forall l \in \bar{\mathcal{L}}, \forall i \in \mathcal{I}, \tag{41b}$$

$$\sum_{l=1}^{L} \alpha_{k,l} = 1, \alpha_{k,l} \geq 0, \ \forall l \in \mathcal{L}. \tag{41c}$$

Although this problem is non-convex, we can make it affine using (32). Then, (41) becomes

$$\max_{\mathbf{A}} \quad \log\left(\alpha_{k,L} b_{k,L\to L}\right) \tag{42a}$$

$$\text{s.t.} \quad \Psi_{k,l} \leq \log\left(\alpha_{k,l} b_{k,i\to l}\right), \forall l \in \bar{\mathcal{L}}, \forall i \in \mathcal{I}, \tag{42b}$$

$$\sum_{l=1}^{L} \alpha_{k,l} = 1, \alpha_{k,l} \geq 0, \ \forall l \in \mathcal{L}. \tag{42c}$$

Here, $\Psi_{k,l} = R_{k,l}^{th}$ [3]. The objective function in (42a) is monotonically increasing in $\alpha_{k,L}$. As

$$\alpha_{k,L} = 1 - \sum_{l=1}^{L-1} \alpha_{k,l}, \tag{43}$$

---

[3]In the next subsection, we will alter this definition.

---

**Algorithm 2** WMMSE1: MMF Algorithm with PA

1: **Init:** $\mathbf{A}^{(0)}$, $E_{tx}$, $\Upsilon$, $\mathbf{P}^{(0)}$, $R_{k,l}^{th}$, $R_g^{(0)} = 0$, $n_{max}$, $n = 0$;
2: **iterate** $\forall j, l, k$;
3: $\quad n = n + 1$
4: $\quad$ Compute $V_{k,i\to l}^{(n)}$ using (29)
5: $\quad$ Compute $\varepsilon_{k,i\to l}^{(n)}$ using (28)
6: $\quad$ Compute $b_{k,i\to l}^{(n)}$ using (36)
7: $\quad$ Update $\left\{\mathbf{P}_k^{(n)}, R_g^{(n)}\right\}$ by solving (40) for given $V_{k,i\to l}^{(n)}$ and $b_{k,i\to l}^{(n)}$
8: $\quad$ Update $\mathbf{A}^{(n)}$ using (46)
9: $\quad$ **If** $(R_g^{(n)} - R_g^{(n-1)})/R_g^{(n-1)} < \Upsilon$ **or** $n = n_{max}$ **then**
10: $\quad\quad$ Terminate
11: **else then**
12: $\quad$ Go to Step 2

---

we can restate (42) as

$$\min_{\mathbf{A}} \quad \sum_{l=1}^{L-1} \alpha_{k,l} \tag{44a}$$

$$\text{s.t.} \quad \frac{e^{\Psi_{k,l}}}{b_{k,i\to l}} \leq \alpha_{k,l}, \forall l \in \bar{\mathcal{L}}, \forall i \in \mathcal{I}, \tag{44b}$$

$$\alpha_{k,l} \geq 0, \ \forall l \in \mathcal{L}. \tag{44c}$$

Then, we can obtain the optimal PAC as

$$\alpha_{k,l}^{opt} = \max_{i\in\mathcal{I}} \frac{e^{\Psi_{k,l}}}{b_{k,i\to l}}, \forall l \in \bar{\mathcal{L}}. \tag{45}$$

and $\alpha_{k,L}^{opt}$ can be obtained using (43). The optimal $\alpha_{k,l}^{opt}$ always satisfies (44c) $\forall l \in \bar{\mathcal{L}}$. However, this may not be true for the $L$-th user in each cluster $k$. Therefore, we update $\alpha_{k,l}^{(n)}$ as

$$\alpha_{k,l}^{(n)} = \begin{cases} \alpha_{k,l}^{(n-1)}, & \alpha_{k,L}^{opt} \leq 0 \\ \alpha_{k,l}^{opt}, & \text{otherwise} \end{cases}, \quad \forall l \in \mathcal{L}. \tag{46}$$

### C. A Low Complexity WMMSE Solution

Algorithm 2 resorts to convex solvers in Step 7 to solve (40) for given receivers and MMSE weights. Although, this results in the optimal solution in Step 7, it significantly increases computational complexity. In this subsection, we propose a low complexity solution. As the objective function and the constraints in (40) are all continuously differentiable, we can make use of the KKT conditions to reduce the search space and thus to decrease complexity.

When $\theta_k, \Gamma_{k,l}, \eta_{k,i\to l}$ and $\beta$ denote Lagrange multipliers, the Lagrangian objective function of (40) is written as

$$h(\mathbf{P}, \mathbf{R}, R_g, \mathbf{V}) = -R_g + \sum_{k=1}^{K} \theta_k(R_g - R_{k,L})$$

$$+ \sum_{k=1}^{K} \sum_{l=1}^{L} \sum_{i=l}^{L} \eta_{k,i\to l}(R_{k,l} - 1 + \xi_{k,\to l})$$

$$+ \sum_{k=1}^{K} \sum_{l=1}^{L-1} \Gamma_{k,l}(R_{k,l}^{th} - R_{k,l}) + \beta(\text{Tr}(\mathbf{PP}^H) - E_{tx}). \tag{47}$$

The optimal precoders and the receivers have to satify the KKT conditions for (47), and are given in the following theorem.

*Theorem 1:* For the optimization problem defined in (40), the following receivers $V_{k,i\to l}$, the Lagrange multiplier $\beta$, and the transmit precoder vectors $\mathbf{p}_k$ satisfy the KKT conditions.

$$V_{k,i\to l} = \alpha_{k,l}\mathbf{p}_k^H \mathbf{h}_{k,i} T_{k,i\to l}^{-1}, \tag{48}$$

$$\beta = \frac{1}{E_{tx}}\left[\sum_{k=1}^K \sum_{l=1}^L \sum_{i=l}^L \eta_{k,i\to l} b_{k,i\to l}|V_{k,i\to l}|^2\right], \tag{49}$$

$$\mathbf{p}_k = \left[\beta\mathbf{I} + \sum_{l=1}^L \sum_{i=l}^L \sum_{j=l}^L \alpha_{k,j}\eta_{k,i\to l}b_{k,i\to l}\mathbf{h}_{k,i}|V_{k,i\to l}|^2\mathbf{h}_{k,i}^H\right.$$
$$\left. + \sum_{t=1,t\neq k}^L \sum_{l=1}^L \sum_{i=l}^L \eta_{t,i\to l}b_{t,i\to l}\mathbf{h}_{t,i}|V_{t,i\to l}|^2\mathbf{h}_{t,i}^H\right]^{-1}$$
$$\times \left[\sum_{l=1}^L \sum_{i=l}^L \eta_{k,i\to l}b_{k,i\to l}\alpha_{k,l}\mathbf{h}_{k,i}V_{k,i\to l}^*\right]. \tag{50}$$

*Proof:* The proof is provided in Appendix A. ∎

*Remark 1:* The receiver $V_{k,i\to l}$ in (48) is exactly equal to the MMSE receiver $V_{k,i\to l}^{mmse}$ given in (29).

*Remark 2:* When the optimal MMSE receiver $V_{k,i\to l}^{mmse}$ and the weights $b_{k,i\to l}^{mmse}$ in (36) are substituted in $\xi_{k,i\to l}$ of (33), then $\xi_{k,i\to l}$ becomes equal to $\xi_{k,i\to l}^{mmse}$.

Utilizing Theorem 1, we propose solving for the receivers (48), the Lagrange multiplier (49) and the precoders (50) in an iterative fashion in Algorithm 3. However, calculating the Lagrange multipliers set $\{\theta_k, \Gamma_{k,l}, \eta_{k,i\to l}, \forall i, l, k\}$ is not trivial. In [31], an exponential penalty method is suggested to solve min-max type problems. According to the exponential penalty method, in each iteration of the algorithm, we update $\{\theta_k, \Gamma_{k,l}, \eta_{k,i\to l}\}$ as

$$\theta_k = \frac{\exp\{\nu(R_g - R_{k,L})\}}{\sum_{k=1}^K \exp\{\nu(R_g - R_{k,L})\}}, \forall k \in \mathcal{K}, \tag{51}$$

$$\Gamma_{k,l} = \exp\{\nu(R_{k,l}^{th} - R_{k,l})\}, \forall k \in \mathcal{K}, \forall l \in \bar{\mathcal{L}}, \tag{52}$$

$$\eta_{k,i\to l} = \Gamma_{k,l}\frac{\exp\{\nu(R_{k,l} - R_{k,i\to l})\}}{\sum_{i=l}^L \exp\{\nu(R_{k,l} - R_{k,i\to l})\}},$$
$$\eta_{k,L\to L} = \theta_k, \forall k \in \mathcal{K}, \forall l \in \bar{\mathcal{L}}, \forall i \in \mathcal{I}. \tag{53}$$

In the above equations, $\nu$ is a constant and as long as $\nu \geq (\log KL)/\epsilon$, the solution is $\epsilon$-optimal. Note that, this choice satisfies the KKT conditions on $\{\theta_k, \Gamma_{k,l}, \eta_{k,i\to l}\}$ since $\sum_{k=1}^K \theta_k = 1$, $\sum_{i=l}^L \eta_{k,i\to l} = \Gamma_{k,l}, \forall l \in \bar{\mathcal{L}}$ and $\eta_{k,L\to L} = \theta_k$, $\theta_k \geq 0, \Gamma_{k,l} \geq 0, \eta_{k,i\to l} \geq 0$.

In each iteration, Algorithm 3 increases the objective function, since there is a total power constraint. Thus, the proposed WMMSE algorithm converges to an upper limit. This limit is within an $\epsilon$ neighborhood of a local optimum, as the algorithm utilizes the equations found via the KKT conditions, and the exponential penalty method is employed. Following similar steps as in [30, Section IV-A] and [32], one can prove convergence in full detail.

*1) Power Allocation for Low Complexity WMMSE:* Algorithm 2 always returns a solution at Step 7, as the CVX approach returns the final result for (40) for given receivers and weights. On the other hand, the low-complexity WMMSE

---

**Algorithm 3** WMMSE2: Low Complexity MMF Algorithm with PA

1: **Init:** $\epsilon$, $\mathbf{A}^{(0)}$, $E_{tx}$, $\Upsilon$, $\Delta$, $\mathbf{P}^{(0)}$, $R_{k,l}^{th}$, $R_g^{(0)} = 0$,
$\nu = \log(KL)/\epsilon$, $n_{max}$, $n = 0$;
2: **iterate** $\forall j, l, k$;
3:     $n = n + 1$
4:     Compute $V_{k,i\to l}^{(n)}$ using (48)
5:     Compute $\varepsilon_{k,i\to l}^{(n)}$ using (28)
6:     Compute $b_{k,i\to l}^{(n)}$ using (36)
7:     Compute $\Gamma_{k,l}^{(n)}$ using (52)
8:     Compute $\theta_k^{(n)}$ using (51)
9:     Compute $\eta_{k,i\to l}^{(n)}$ using (53)
10:    Compute $\beta^{(n)}$ using (49)
11:    Compute $\mathbf{P}^{(n)}$ using (50)
12:    Scale $\mathbf{P}^{(n)}$ such that $\text{Tr}(\mathbf{P}^{(n)}\mathbf{P}^{(n)H}) = E_{tx}$
13:    Update $\mathbf{A}^{(n)}$ using (46)
14:    **If** $(R_g^{(n)} - R_g^{(n-1)})/R_g^{(n-1)} < \Upsilon$ **or** $n = n_{max}$ **then**
15:      **If** (38c) satisfied **then**
16:       Terminate
17:      **else then**
18:       $\nu = \nu + \Delta$, Go to Step 2
19:    **else then**
20:      Go to step 2

---

approach may not be feasible in each iteration, as it only provides a step in the favorable direction in each iteration. Therefore, in Algorithm 3 at Step 12, the updated precoder $\mathbf{P}^{(n)}$ may not satisfy the threshold rate constraints in (38c), and the algorithm may not find a feasible PAC at Step 13. One approach would be to skip power optimization, immediately update $\nu$ and proceed with the next iteration. However, we choose to find the best PAC that satisfies the current achievable rates. Thus, we update $\Psi_{k,l}$ in (42b) in each iteration as

$$\Psi_{k,l}^{(n)} = \min\left(R_{k,l}^{th}, R_{k,l}^{(n)}\right). \tag{54}$$

## V. NUMERICAL RESULTS

In this section, we present numerical results to evaluate the performance of the proposed transmission strategies given in Algorithms 1, 2 and 3. We compare these algorithms with OMA, MULP and RS in terms of MMF rates, energy efficiency and computational complexity. All three algorithms we propose carry out power optimization. We also compare them with their fixed power allocation versions.

### A. Orthogonal Multiple Access, Multiuser Linear Precoding and Rate Splitting

Before presenting any simulation results, in this subsection, we first describe the schemes used as benchmarks: OMA, MULP and RS.

*1) OMA:* In OMA, the transmission time is divided into $L$ equal slots. The base station communicates with the $l$-th strongest users in each cluster in each time slot-$l$. The input data vector for time slot $l$ is denoted as $\mathbf{s}^{l,OMA} =$

$[s_{1,l}, \ldots, s_{K,l}]^T \in \mathbb{C}^{K \times 1}$. We assume all $s_{k,l}$ are independent and $\mathbb{E}\{s_{k,l}s_{k,l}^*\} = 1$. The input data vector $\mathbf{s}^{l,OMA}$ is linearly processed by a precoder matrix $\mathbf{P}^{l,OMA} = [\mathbf{p}_1^{l,OMA}, \ldots, \mathbf{p}_K^{l,OMA}] \in \mathbb{C}^{M \times K}$, where the precoding vector $\mathbf{p}_k^{l,OMA} \in \mathbb{C}^{M \times 1}$ is dedicated to the $k$-th user in time slot-$l$. The overall transmit data vector $\mathbf{x}^{l,OMA} \in \mathbb{C}^{M \times 1}$ at the base station can be written as $\mathbf{x}^{l,OMA} = \mathbf{P}^{l,OMA}\mathbf{s}^{l,OMA}$. Then, the SINR at user-$k$ in time slot-$l$ is given by

$$\gamma_k^{l,OMA} = \frac{\left|\mathbf{h}_{k,l}^H \mathbf{p}_k^{l,OMA}\right|^2}{\sum_{i=1,i\neq k}^{K} \left|\mathbf{h}_{k,l}^H \mathbf{p}_i^{l,OMA}\right|^2 + 1} \tag{55}$$

and the corresponding rate expression is calculated as $R_{k,l}^{OMA} = \frac{1}{L}\log(1 + \gamma_k^{l,OMA})$.

Given these assumptions, the MMF OMA problem is equivalent to providing fairness in the last time slot, while satisfying the threshold rate constraints in earlier time slots. We can formulize the MMF OMA optimization problem as

$$\max_{\mathbf{P}^{l,OMA},\forall l \in \mathcal{L}} \quad \min_{k \in \mathcal{K}} R_{k,L}^{OMA} \tag{56a}$$

$$\text{s.t.} \quad R_{k,l}^{th} \leq R_{k,l}^{OMA}, \forall k \in \mathcal{K}, \forall l \in \bar{\mathcal{L}} \tag{56b}$$

$$\text{Tr}(\mathbf{P}^{l,OMA}\mathbf{P}^{l,OMA^H}) \leq E_{tx}, \forall l \in \mathcal{L}. \tag{56c}$$

Then, the MMF OMA rate $R^{OMA}$ can be calculated as

$$R^{OMA} = \min_{k \in \mathcal{K}} R_{k,L}^{OMA}$$

using the optimal precoders $\mathbf{P}^{l,OMA^*}$, $\forall l \in \mathcal{L}$ that solve (56). Note that precoders $\mathbf{P}^{l,OMA^*}$, $\forall l \in \bar{\mathcal{L}}$ are required to satisfy the rate constraints in (56b), whereas $\mathbf{P}^{L,OMA^*}$ provides fairness among the strongest users in each cluster.

*2) MULP:* In MULP precoding, the base station transmits data to all $KL$ users simultaneously. The input data vector is denoted as $\mathbf{s}^{MULP} = [s_{1,1}, \ldots, s_{1,L}, \ldots, s_{K,1}, \ldots, s_{K,L}]^T \in \mathbb{C}^{KL \times 1}$. We assume all $s_{k,l}$ are independent and $\mathbb{E}\{s_{k,l}s_{k,l}^*\} = 1$. The input data vector $\mathbf{s}^{MULP}$ is linearly processed by a precoder matrix $\mathbf{P}^{MULP} = [\mathbf{p}_{1,1}^{MULP}, \ldots, \mathbf{p}_{1,L}^{MULP}, \ldots, \mathbf{p}_{K,1}^{MULP}, \ldots, \mathbf{p}_{K,L}^{MULP}] \in \mathbb{C}^{M \times KL}$, where the precoding vector $\mathbf{p}_{k,l}^{MULP} \in \mathbb{C}^{M \times 1}$ is dedicated to the $l$-th user in the $k$-th cluster. Then, the overall transmit data vector $\mathbf{x}^{MULP} \in \mathbb{C}^{M \times 1}$ at the base station can be written as $\mathbf{x}^{MULP} = \mathbf{P}^{MULP}\mathbf{s}^{MULP}$. The SINR at user-$l$ in the $k$-th cluster is given by

$$\gamma_{k,l}^{MULP} = \frac{|\mathbf{h}_{k,l}^H \mathbf{p}_{k,l}^{MULP}|^2}{\sum_{\substack{j=1 \\ j\neq l}}^{L} |\mathbf{h}_{k,l}^H \mathbf{p}_{k,j}^{MULP}|^2 + \sum_{\substack{i=1 \\ i\neq k}}^{K} \sum_{l=1}^{L} |\mathbf{h}_{i,l}^H \mathbf{p}_{i,l}^{MULP}|^2 + 1}, \tag{57}$$

and the corresponding rate expression is calculated as $R_{k,l}^{MULP} = \log(1 + \gamma_{k,l}^{MULP})$.

For a fair comparison, we assume that fairness among the strongest users is needed while satisfying the threshold

rate constraints on other users. The MMF MULP problem is written as

$$\max_{\mathbf{P}^{MULP}} \quad \min_{k \in \mathcal{K}} R_{k,L}^{MULP} \tag{58a}$$

$$\text{s.t.} \quad R_{k,l}^{th} \leq R_{k,l}^{MULP}, \forall k \in \mathcal{K}, \forall l \in \bar{\mathcal{L}} \tag{58b}$$

$$\text{Tr}(\mathbf{P}^{MULP}\mathbf{P}^{MULP^H}) \leq E_{tx}. \tag{58c}$$

Then, the MMF MULP rate $R^{MULP}$ can be calculated as

$$R^{MULP} = \min_{k \in \mathcal{K}} R_{k,L}^{MULP}$$

using the optimal precoder $\mathbf{P}^{MULP^*}$ that solve (58).

*3) 1-Layer RS:* In 1-Layer RS, we use the same signal model proposed in [23]. In this strategy, the message stream of the $l$-th user in the $k$-th cluster is split into common and private parts. The common part is at rate $C_{k,l}^{RS}$ and the private part is at rate $R_{k,l}^{RS}$. The common parts are collectively encoded as a common message $s_c$ at rate $R_c^{RS} = \sum_{k \in \mathcal{K}} \sum_{l \in \mathcal{L}} C_{k,l}^{RS}$. The private messages are encoded as $s_{k,l,p}$. To send all the messages, the base station encodes the input data vector $\mathbf{s}^{RS} = [s_c, s_{1,1,p}, \ldots, s_{1,L,p}, \ldots, s_{K,1,p} \ldots, s_{K,L,p}]^T \in \mathbb{C}^{(KL+1)\times 1}$ by a precoder matrix $\mathbf{P}^{RS} = [\mathbf{p}_c^{RS}, \mathbf{p}_{k,1}^{RS}, \ldots, \mathbf{p}_{k,L}^{RS}, \ldots, \mathbf{p}_{K,l}^{RS}, \ldots, \mathbf{p}_{K,L}^{RS}]$. Here, $\mathbf{p}_c^{RS}$ and $\mathbf{p}_{k,l}^{RS} \in \mathbb{C}^{M \times 1}$ respectively indicate the precoder vectors for the common data $s_c$ and the private data $s_{k,l,p}$. The base station transmits $\mathbf{x}^{RS} \in \mathbb{C}^{M \times 1}$, which is equal to $\mathbf{x}^{RS} = \mathbf{P}^{RS}\mathbf{s}^{RS}$ Then, the SINR at the $l$-th user in the $k$-th cluster for common and private data messages respectively become

$$\gamma_{k,l,c}^{RS} = \frac{|\mathbf{h}_{k,l}^H \mathbf{p}_c^{RS}|^2}{\sum_{i=1}^{K}\sum_{l=1}^{L}|\mathbf{h}_{k,l}^H \mathbf{p}_{k,l}^{RS}|^2 + 1}, \tag{59}$$

$$\gamma_{k,l,p}^{RS} = \frac{|\mathbf{h}_{k,l}^H \mathbf{p}_{k,l}^{RS}|^2}{\sum_{\substack{j=1 \\ j\neq l}}^{L}|\mathbf{h}_{k,l}^H \mathbf{p}_{k,j}^{RS}|^2 + \sum_{\substack{i=1 \\ i\neq k}}^{K}\sum_{l=1}^{L}|\mathbf{h}_{i,l}^H \mathbf{p}_{i,l}^{RS}|^2 + 1}, \tag{60}$$

and the corresponding rate expressions are calculated as $R_{k,l,c}^{RS} = \log(1 + \gamma_{k,l,c}^{RS})$ and $R_{k,l,p}^{RS} = \log(1 + \gamma_{k,l,p}^{RS})$. As the common rate has to be decoded by all users, we define $R_c^{RS} = \min_{k \in \mathcal{K}} \min_{l \in \mathcal{L}} R_{k,l,c}^{RS}$. Then, the MMF RS problem can be stated as

$$\max_{\mathbf{P}^{RS}} \quad \min_{k \in \mathcal{K}} \left(C_{k,L}^{RS} + R_{k,L}^{RS}\right) \tag{61a}$$

$$\text{s.t.} \quad R_{k,l}^{th} \leq \left(C_{k,l}^{RS} + R_{k,l}^{RS}\right), \forall k \in \mathcal{K}, \forall l \in \bar{\mathcal{L}} \tag{61b}$$

$$\sum_{k \in \mathcal{K}} \sum_{l \in \mathcal{L}} C_{k,l}^{RS} \leq R_c^{RS}, \tag{61c}$$

$$0 \leq C_{k,l}^{RS}, \forall k \in \mathcal{K}, \forall l \in \mathcal{L} \tag{61d}$$

$$\text{Tr}(\mathbf{P}^{RS}\mathbf{P}^{RS^H}) \leq E_{tx}. \tag{61e}$$

As a result, the MMF RS rate $R^{RS}$ becomes

$$R^{RS} = \min_{k, \in \mathcal{K}} \left(C_{k,L}^{RS} + R_{k,L}^{RS}\right)$$

employing the optimal precoder that solves (61).

To solve all optimization problems stated for OMA, MULP and RS, we first find their equivalent weighted MMSE problems and solve them in an iterative fashion as done in Algorithm 2 in Section IV.
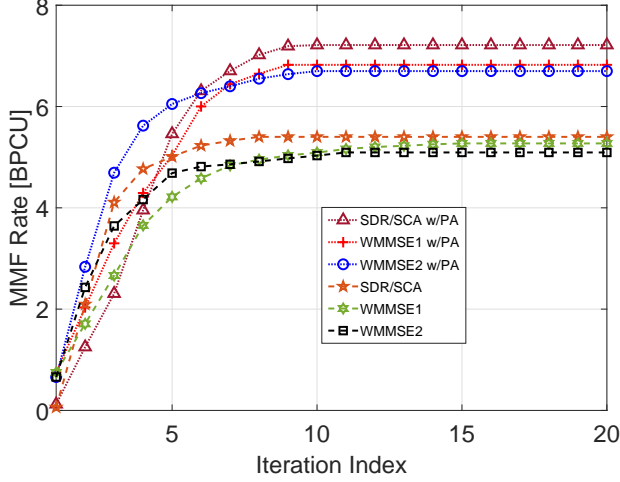
Fig. 1. MMF rate convergence performance of proposed algorithms for $M = K = 3, L = 2, R_{k,l}^{th} = 0.1$ bpcu. Transmit SNR is set to 15 dB.
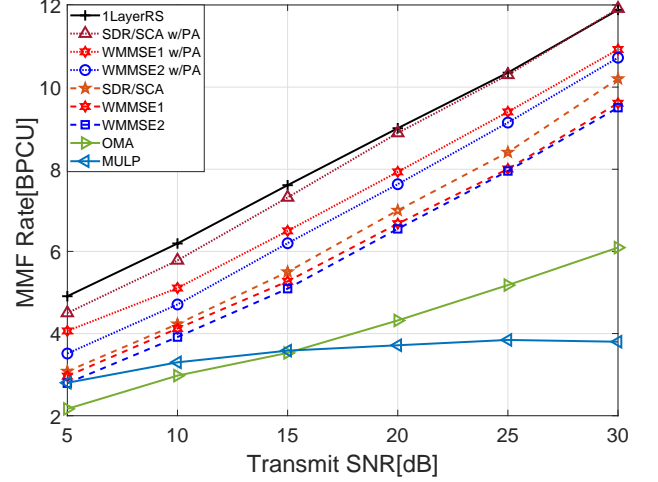


Fig. 2. MMF rates for different precoder schemes for $M = K = 3, L = 2, R_{k,l}^{th} = 0.1$ bpcu.



Fig. 3. MMF rates for different precoder schemes for $M = K = 4, L = 2, R_{k,l}^{th} = 0.1$ bpcu.

### B. Assumptions

In the simulations, the entries in $\tilde{\mathbf{h}}_{k,l}$ are assumed to be i.i.d. circularly symmetric complex Gaussian random variables with zero mean and unit variance. The path loss exponent is set to $\rho = 4$. The users are uniformly distributed in a circular region of radius 1. These users are clustered according to the scheme proposed in [17, Algorithm 1, Figure 3]. In this clustering scheme, the aim is to put users, which have highly different effective channel gain magnitudes in the same cluster. For example, for $L = 2$, the base station puts the user with the highest effective channel gain magnitude in the same cluster with the worst user among all users. The second best and and the second worst users are grouped as a second cluster. The remaining clusters are formed in a similar fashion. Note that for all the NOMA schemes, the base station has to inform the users about their order and the other users in their own cluster so that users within a cluster can perform SIC.

For the fixed power allocation versions of Algorithms 1, 2 and 3, we assume the power allocation scheme suggested in [17, Table 1], which assigns more power to weak users and less power to strong users. This idea is in line with power domain NOMA and widely used in the literature [16], [18]. This fixed power allocation vector is also used as the initial value of $\mathbf{A}^{(0)}$ in Algorithms 1, 2 and 3.

For Algorithms 1, 2 and 3, the presented results are averaged over $10^2$ channel realizations. The maximum number of iterations $n_{max}$ is limited to 100 and $\Upsilon$ are set to $10^{-3}$. The transmit signal to noise ratio (SNR) is defined as $E_{tx}/\sigma^2$. Here $\sigma^2$ is the noise variance and set to 1. For Algorithm 3, $\epsilon$ and $\Delta$ are set to $10^{-3}$ and 3 respectively. The parameter $\Delta$ is used to tune the algorithm to satisfy the rate constraint $R_{k,l}^{th}$. Finally, if a particular algorithm is infeasible, we set its MMF rate to zero to make a fair comparison among all algorithms under consideration [20].

In the following simulation results, we consider algorithm convergence, MMF rate and energy efficiency results for

different settings. Rates are expressed in terms of bits per channel use (bpcu).

### C. Simulation Results

Fig. 1 shows the convergence behavior of the proposed schemes given by Algorithms 1, 2 and 3 with and without PA for $M = 3, K = 3, L = 2, R_{k,l}^{th} = 0.1 \ \forall k \in \mathcal{K}, \forall l \in \bar{\mathcal{L}}$, when the total transmit power is set to 15 dB. The initial precoder matrix, $\mathbf{P}^{(0)}$ in Algorithms 1, 2 and 3 is assumed to be the identity matrix, scaled to satisfy the power constraint. The figure confirms that the proposed algorithms converge fast.

Figs. 2 and 3 compare MMF rates for the proposed algorithms with 1-layer RS, OMA and MULP schemes for $R_{k,l}^{th} = 0.1$ bits $\forall k \in \mathcal{K}, \forall l \in \bar{\mathcal{L}}$ for $M = 3, K = 3, L = 2$, and for $M = 4, K = 4, L = 2$ respectively. We observe that 1-layer RS has the best performance in terms of MMF

| | RUN TIME IN SECONDS | |
|---|---|---|
| ALGORITHM | $M = K = 3, L = 2$ | $M = K = 4, L = 2$ |
| WMMSE2 | 6.38 | 13.42 |
| WMMSE2 w/PA | 38.35 | 51.52 |
| SDR/SCA | 442.13 | 591.21 |
| SDR/SCA w/PA | 491.23 | 677.29 |
| WMMSE1 | 2134.45 | 3743.80 |
| WMMSE1 w/PA | 2752.55 | 4167.50 |
| 1-layer RS | 60457.80 | 75994.30 |

TABLE I
COMPLEXITY OF ALGORITHMS



Fig. 4. MMF rates for different precoder schemes for $M = 6, K = 3, L = 2$, $R_{k,l}^{th} = 0.1$ bpcu.

rates. It can effectively mitigate interference by adjusting the common message rate. Algorithm 1 (SDR/SCA w/PA) has similar performance with 1-layer RS and as SNR increases the gap between the two algorithms diminish. Both Algorithms 2 (WMMSE1 w/PA) and 3 (WMMSE2 w/PA) perform worse than Algorithm 1. This is because semi-definite programming with successive convex approximation is an effective approximation. In each iteration, the constraints in (23b) and (23c) become tighter and (23) approaches the original optimization problem in (8). As expected, Algorithms 2 and 3 have similar results. All algorithms are several dB better than their fixed power allocation versions (SDR/SCA, WMMSE1, WMMSE2). MULP is very inefficient in interference management, and displays very poor performance. The MMF rate for MULP converges for high SNR.

From Figs. 2 and 3 we also observe that all Algorithms 1, 2 and 3 (with or without power allocation) and the RS scheme present full degrees of freedom (DoF); i.e. 1. DoF is calculated as the MMF rate (in bpcu) over $\log_2$ SNR [33]. While OMA can accommodate all users in each time slot, it suffers from time division and its DoF is limited with 0.5. Although a detailed DoF analysis is out of the scope of this paper, we conjecture that the DoF for MULP for the overloaded settings in Figs. 2 and 3 is 0. This is because, the MMF rate calculation for MULP is similar to the MMF rate calculation for the designated beamforming scheme in the multigroup multicasting scenario examined in [33]. For the latter, the DoF is proved to be 0 either for $M = 3$ and there are 3 groups with 2 users each or for $M = 4$ and there are 4 groups with 2 users each.

Figs. 2 and 3 should be interpreted together with the complexity results given in Table I. Table I shows the complexity of all the algorithms under consideration. We observe that Algorithm 3 has the least complexity either with or without power optimization. For Algorithms 1, 2 and 3, power optimization does not change algorithm complexity and run time values are on the same order. Although 1-layer RS has the highest MMF rates in Figs. 2 and 3, it also has the highest complexity. The run time for 1-layer RS is 3-4 orders of magnitude larger than the run time for Algorithm 3, which is based on the closed form expressions of Theorem 1. Algorithm 2 has 2 orders of magnitude larger complexity than Algorithm 3 either with or without power optimization. As they achieve similar MMF rates, we conclude that Algorithm 3 is more advantageous than Algorithm 2. In conclusion, 1-layer RS has the best MMF rate performance, Algorithm 3 has the least complexity, and
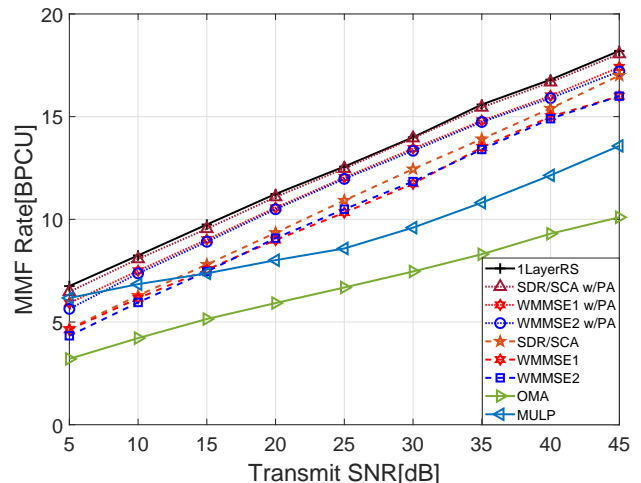
Algorithm 1 provides a good tradeoff between complexity and MMF rates. It performs almost the same as 1-layer RS in MMF rates, and its complexity is only an order of magnitude larger than that of Algorithm 3.

Note that, while solving (61) with the MMSE approach, one could apply the KKT optimality conditions and the ordinary penalty method, instead of calling for CVX. However, the common and private rate expressions for rate splitting are complex and numerous, and finding the expressions the optimal precoders, receivers, weights and Lagrange multipliers as in Theorem 1 is complicated, keeping the complexity for 1-layer RS high.

Figs. 2 and 3 are for overloaded systems. Fig. 4 shows how the MMF rates change, when $M$ is at least as large as $KL$. In the figure $M = 6$, $K = 3$ and $L = 2$ and $R_{k,l}^{th} = 0.1$ bpcu $\forall k \in \mathcal{K}$ and $\forall l \in \bar{\mathcal{L}}$. For this setting, the system is not overloaded, intense interference mitigation is not necessary, and benefits of rate splitting is less. Thus, SDR/SCA and WMMSE based schemes with or without power optimization are closer to 1-layer RS. RS and all the algorithms have full DoF equal to 1. OMA, by definition, still suffers from time division and its DoF is limited with 0.5. As the number of base station antennas is sufficient to serve all the users simultaneously, MULP also presents full DoF. This result is expected because the DoF for the designated beamforming scheme in [33] is proved to be 1, when there is a single user in each group and the number of base station antennas is equal to the number of groups. However, MULP does not achieve this performance easily, its DoF result does not converge until 30 dB or higher. MULP is quite inefficient in interference mitigation and the additional threshold rate constraints for the weakest users in each group makes the MULP problem in (58) harder to solve especially for low to medium SNR.

Fig. 5 shows that MMF rates decrease, when the threshold rates $R_{k,l}^{th}$, which is assumed to be the same $\forall k \in \mathcal{K}$ and $\forall l \in \bar{\mathcal{L}}$, increase from 0.1 to 0.4 bpcu for $M = 3$, $K = 3$ and $L = 2$. The transmit SNR is set to 15 dB. Note that,
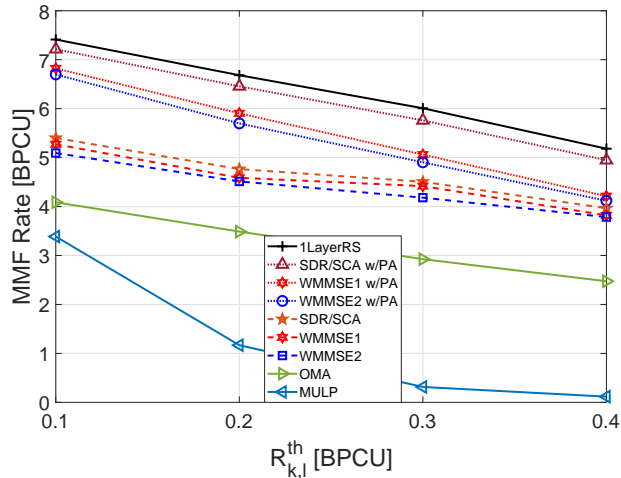
Fig. 5. MMF rates vs. $R_{k,l}^{th}$ for different precoder schemes for $M = K = 3, L = 2$. Transmit SNR is set to 15 dB.
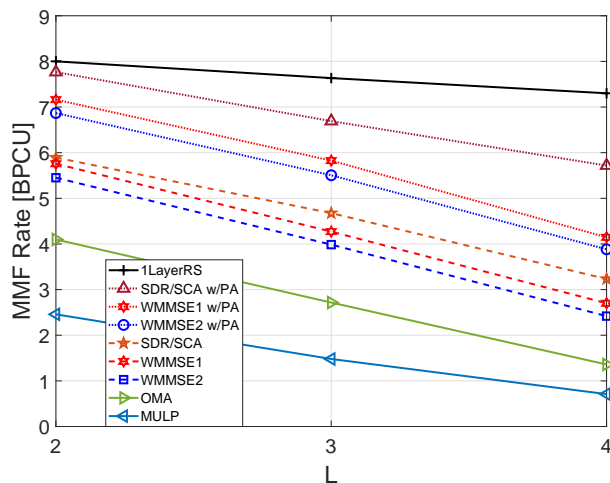


Fig. 7. Energy efficiency for different precoder schemes for $M = K = 3, L = 2, R_{k,l}^{th} = 0.1$ bpcu.



Fig. 6. MMF rates for different precoder schemes for $M = 2, K = 2, L = \{2, 3, 4\}, R_{k,l}^{th} = 0.1$ bpcu.

one could expect MMF rate for OMA to be constant with increasing SNR. Time is divided into slots and all the strongest users are served in the last time slot, seemingly unaffected from all the other users. However, unless all the threshold rate constraints are satisfied, OMA is infeasible and MMF rate for OMA is zero. Therefore, MMF rate for OMA also decreases with increasing $R_{k,l}^{th}$. MULP rates decrease much faster than other schemes as the feasible set quickly shrinks with increasing $R_{k,l}^{th}$.

Fig. 6 presents the effect of increasing number of users in each cluster for $M = 2, K = 2, L = \{2, 3, 4\}, R_{k,l}^{th} = 0.1$ bpcu $\forall k \in \mathcal{K}$ and $\forall l \in \bar{\mathcal{L}}$. The decrease in MMF rates for 1-layer RS is much slower than all the other schemes as it provides excellent interference mitigation.

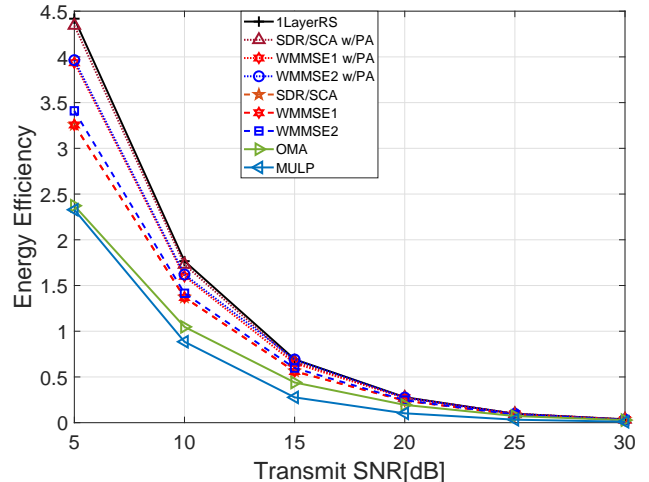Finally, in Fig. 7, we compare all the schemes in terms of

energy efficiency. Energy efficiency is defined as

$$EE = \frac{\sum_{k \in \mathcal{K}} \sum_{l \in \mathcal{L}} R_{k,l}}{\text{Tr}\left(\mathbf{PP}^H\right)} \quad (62)$$

for all precoding schemes. We observe that Algorithm 1 has the same energy efficiency as 1-layer RS. The results show that the gap between the algorithms are smaller. Together with the results in Figs. 2 and 3, and Table I, we conclude that SDR/SCA with power allocation is an excellent scheme with high MMF rates, low complexity and high energy efficiency.

## VI. CONCLUSION

We consider a joint precoder and power allocation design problem in downlink MIMO-NOMA to achieve max-min fairness among the strongest users in each cluster, while satisfying threshold rate constraints for all the other users. We propose 3 algorithms: (i) SDR/SCA, (ii) WMMSE1 and (iii) WMMSE2. The first algorithm is based on semi-definite relaxation and successive convex approximation, and the latter two are based on the relation between rate and minimum mean square error. WMMSE2 incorporates further simplifications in WMMSE1 based on the KKT optimality conditions and the ordinary penalty method. We compare our results with RS, OMA and MULP schemes. The results reveal that SDR/SCA scheme offers high MMF rates and superior energy efficiency at very low complexity. Future work includes designing precoders for imperfect channel state information and for finite block length channel coding.

## APPENDIX A

In this appendix, we prove Theorem 1. Taking the derivative of the objective function $h$ in (47) with respect to $V_{k,i \rightarrow l}$, then

equating it to zero, we obtain

$$\alpha_{k,l}\mathbf{p}_k^H\mathbf{h}_{k,i} = \sum_{j=l}^{L}\alpha_{k,j}\mathbf{h}_{k,i}^H\mathbf{p}_k\mathbf{p}_k^H\mathbf{h}_{k,i}V_{k,i\to l}$$
$$+ \sum_{t=1,t\neq k}^{K}\mathbf{h}_{k,i}^H\mathbf{p}_t\mathbf{p}_t^H\mathbf{h}_{k,i}V_{k,i\to l} + V_{k,i\to l}. \tag{63}$$

Then, when $\eta_{k,i\to l} > 0$,

$$V_{k,i\to l} = \alpha_{k,l}\mathbf{p}_k^H\mathbf{h}_{k,i}T_{k,i\to l}^{-1}. \tag{64}$$

Secondly, taking the gradient of (47) with respect to $\mathbf{p}_k^H$, and equating it to zero, we have the following equation

$$\sum_{l=1}^{L}\sum_{i=l}^{L}\eta_{k,i\to l}b_{k,i\to l}\mathbf{h}_{k,i}V_{k,i\to l}^*\alpha_{k,l}$$
$$= \sum_{l=1}^{L}\sum_{i=l}^{L}\sum_{j=l}^{L}\alpha_{k,j}\eta_{k,i\to l}b_{k,i\to l}\mathbf{h}_{k,i}V_{k,i\to l}^*V_{k,i\to l}\mathbf{h}_{k,i}^H\mathbf{p}_k$$
$$+ \sum_{\substack{t=1\\t\neq k}}^{K}\sum_{l=1}^{L}\sum_{i=l}^{L}\eta_{t,i\to l}b_{t,i\to l}\mathbf{h}_{t,i}|V_{t,i\to l}|^2\mathbf{h}_{t,i}^H\mathbf{p}_k + \beta\mathbf{p}_k. \tag{65}$$

Then,

$$\mathbf{p}_k = \left[\beta\mathbf{I} + \sum_{l=1}^{L}\sum_{i=l}^{L}\sum_{j=l}^{L}\alpha_{k,j}\eta_{k,i\to l}b_{k,i\to l}\mathbf{h}_{k,i}|V_{k,i\to l}|^2\mathbf{h}_{k,i}^H\right.$$
$$\left. + \sum_{t=1,t\neq k}^{K}\sum_{l=1}^{L}\sum_{i=l}^{L}\eta_{t,i\to l}b_{t,i\to l}\mathbf{h}_{t,i}|V_{t,i\to l}|^2\mathbf{h}_{t,i}^H\right]^{-1}$$
$$\times \left[\sum_{l=1}^{L}\sum_{i=l}^{L}\alpha_{k,l}\eta_{k,i\to l}b_{k,i\to l}\mathbf{h}_{k,i}V_{k,i\to l}^*\right]. \tag{66}$$

To calculate $\beta$, we post-multiply both sides of (63) by $V_{k,i\to l}^*\eta_{k,i\to l}b_{k,i\to l}$ and perform $\sum_{k=1}^{K}\sum_{l=1}^{L}\sum_{i=l}^{L}$ on both sides. We also pre-multiply (65) with $\mathbf{p}_k^H$ and sum over $k, k=\{1,2,\ldots,K\}$. After calculating the trace of these two resulting equations, we observe that the left sides of both equations are equal. Then, the right sides are also equal to each other. As we assume that the power constraint in (2) is satisfied with equality we can find that

$$\beta = \frac{1}{E_{tx}}\left[\sum_{k=1}^{K}\sum_{l=1}^{L}\sum_{i=l}^{L}\eta_{k,i\to l}b_{k,i\to l}|V_{k,i\to l}|^2\right]. \tag{67}$$

## REFERENCES

[1] J. G. Andrews, S. Buzzi, W. Choi, S. V. Hanly, A. Lozano, A. C. K. Soong, and J. C. Zhang, "What will 5G be?" *IEEE Journal on Selected Areas in Communications*, vol. 32, no. 6, pp. 1065–1082, 2014.

[2] L. Dai, B. Wang, Y. Yuan, S. Han, C. I, and Z. Wang, "Non-orthogonal multiple access for 5G: solutions, challenges, opportunities, and future research trends," *IEEE Communications Magazine*, vol. 53, no. 9, pp. 74–81, 2015.

[3] T. Cover, "Broadcast channels," *IEEE Transactions on Information Theory*, vol. 18, no. 1, pp. 2–14, 1972.

[4] Y. Saito, Y. Kishiyama, A. Benjebbour, T. Nakamura, A. Li, and K. Higuchi, "Non-orthogonal multiple access (NOMA) for cellular future radio access," in *2013 IEEE 77th Vehicular Technology Conference (VTC Spring)*, 2013, pp. 1–5.

[5] A. Benjebbour, A. Li, Y. Kishiyama, H. Jiang, and T. Nakamura, "System-level performance of downlink NOMA combined with SU-MIMO for future LTE enhancements," in *2014 IEEE Globecom Workshops (GC Wkshps)*, 2014, pp. 706–710.

[6] M. Zeng, A. Yadav, O. A. Dobre, G. I. Tsiropoulos, and H. V. Poor, "Capacity comparison between MIMO-NOMA and MIMO-OMA with multiple users in a cluster," *IEEE Journal on Selected Areas in Communications*, vol. 35, no. 10, pp. 2413–2424, 2017.

[7] Z. Ding, Z. Yang, P. Fan, and H. V. Poor, "On the performance of non-orthogonal multiple access in 5G systems with randomly deployed users," *IEEE Signal Processing Letters*, vol. 21, no. 12, pp. 1501–1505, 2014.

[8] Z. Yang, Z. Ding, P. Fan, and N. Al-Dhahir, "A general power allocation scheme to guarantee quality of service in downlink and uplink NOMA systems," *IEEE Transactions on Wireless Communications*, vol. 15, no. 11, pp. 7244–7257, 2016.

[9] M. F. Hanif, Z. Ding, T. Ratnarajah, and G. K. Karagiannidis, "A minorization-maximization method for optimizing sum rate in the downlink of non-orthogonal multiple access systems," *IEEE Transactions on Signal Processing*, vol. 64, no. 1, pp. 76–88, 2016.

[10] F. Zhu, Z. Lu, J. Zhu, J. Wang, and Y. Huang, "Beamforming design for downlink non-orthogonal multiple access systems," *IEEE Access*, vol. 6, pp. 10956–10965, 2018.

[11] C. Chen, W. Cai, X. Cheng, L. Yang, and Y. Jin, "Low complexity beamforming and user selection schemes for 5G MIMO-NOMA systems," *IEEE Journal on Selected Areas in Communications*, vol. 35, no. 12, pp. 2708–2722, 2017.

[12] F. Alavi, K. Cumanan, Z. Ding, and A. G. Burr, "Beamforming techniques for nonorthogonal multiple access in 5G cellular networks," *IEEE Transactions on Vehicular Technology*, vol. 67, no. 10, pp. 9474–9487, 2018.

[13] Z. Chen, Z. Ding, P. Xu, and X. Dai, "Optimal precoding for a QoS optimization problem in two-user MISO-NOMA downlink," *IEEE Communications Letters*, vol. 20, no. 6, pp. 1263–1266, 2016.

[14] S. Timotheou and I. Krikidis, "Fairness for non-orthogonal multiple access in 5G systems," *IEEE Signal Processing Letters*, vol. 22, no. 10, pp. 1647–1651, 2015.

[15] J. Choi, "Power allocation for max-sum rate and max-min rate proportional fairness in NOMA," *IEEE Communications Letters*, vol. 20, no. 10, pp. 2055–2058, 2016.

[16] Xiaofang Sun, D. Duran-Herrmann, Zhangdui Zhong, and Yaoqing Yang, "Non-orthogonal multiple access with weighted sum-rate optimization for downlink broadcast channel," in *MILCOM 2015 - 2015 IEEE Military Communications Conference*, 2015, pp. 1176–1181.

[17] M. S. Ali, H. Tabassum, and E. Hossain, "Dynamic user clustering and power allocation for uplink and downlink non-orthogonal multiple access (NOMA) systems," *IEEE Access*, vol. 4, pp. 6325–6343, 2016.

[18] S. Ali, E. Hossain, and D. I. Kim, "Non-orthogonal multiple access (NOMA) for downlink multiuser MIMO systems: User clustering, beamforming, and power allocation," *IEEE Access*, vol. 5, pp. 565–577, 2017.

[19] X. Sun, C. Shen, Y. Xu, S. M. Al-Basit, Z. Ding, N. Yang, and Z. Zhong, "Joint beamforming and power allocation design in downlink non-orthogonal multiple access systems," in *2016 IEEE Globecom Workshops (GC Wkshps)*, Dec. 2016, pp. 1–6.

[20] X. Sun, N. Yang, S. Yan, Z. Ding, D. W. K. Ng, C. Shen, and Z. Zhong, "Joint beamforming and power allocation in downlink NOMA multiuser MIMO networks," *IEEE Transactions on Wireless Communications*, vol. 17, no. 8, pp. 5367–5381, 2018.

[21] J. Choi, "Minimum power multicast beamforming with superposition coding for multiresolution broadcast and application to NOMA systems," *IEEE Transactions on Communications*, vol. 63, no. 3, pp. 791–800, 2015.

[22] Z. Liu, L. Lei, N. Zhang, G. Kang, and S. Chatzinotas, "Joint beamforming and power optimization with iterative user clustering for MISO-NOMA systems," *IEEE Access*, vol. 5, pp. 6872–6884, 2017.

[23] Y. Mao, B. Clerckx, and V. O. K. Li, "Rate-splitting multiple access for downlink communication systems: bridging, generalizing, and outperforming SDMA and NOMA," *EURASIP Journal on Wireless Communications and Networking*, 2018.

[24] ——, "Rate-splitting for multi-antenna non-orthogonal unicast and multicast transmission: Spectral and energy efficiency analysis," *IEEE Transactions on Communications*, vol. 67, no. 12, pp. 8754–8770, 2019.

[25] F. Zhang, *The Schur Complement and its Applications*. Springer, Boston, MA, 2005.

[26] A. Khabbazibasmenj, F. Roemer, S. A. Vorobyov, and M. Haardt, "Sum-rate maximization in two-way AF MIMO relaying: Polynomial time

solutions to a class of DC programming problems," *IEEE Transactions on Signal Processing*, vol. 60, no. 10, pp. 5478–5493, 2012.

[27] S. Boyd, L. Xiao, A. Mutapic, and J. Mattingley, "Sequential convex programming notes for EE364b Stanford University," http://www.stanford.edu/class/EE364b/, 2007.

[28] M. Grant and S. Boyd, "CVX: Matlab software for disciplined convex programming, version 2.2," http://cvxr.com/cvx, Jan. 2020.

[29] Dongning Guo, S. Shamai, and S. Verdu, "Mutual information and minimum mean-square error in Gaussian channels," *IEEE Transactions on Information Theory*, vol. 51, no. 4, pp. 1261–1282, 2005.

[30] S. S. Christensen, R. Agarwal, E. De Carvalho, and J. M. Cioffi, "Weighted sum-rate maximization using weighted MMSE for MIMO-BC beamforming design," *IEEE Transactions on Wireless Communications*, vol. 7, no. 12, pp. 4792–4799, 2008.

[31] X. Li and S. Pan, "Solving the finite min-max problem via an exponential penalty method," *Vychislitelnye Tekhnologii*, vol. 8, pp. 3–15, 01 2003.

[32] J. Kaleva, A. Tlli, and M. Juntti, "Decentralized sum rate maximization with QoS constraints for interfering broadcast channel via successive convex approximation," *IEEE Transactions on Signal Processing*, vol. 64, no. 11, pp. 2788–2802, 2016.

[33] H. Joudeh and B. Clerckx, "Rate-splitting for max-min fair multigroup multicast beamforming in overloaded systems," *IEEE Transactions on Wireless Communications*, vol. 16, no. 11, pp. 7276–7289, 2017.