# Rate-Splitting Multiple Access for Downlink MIMO: A Generalized Power Iteration Approach

Jeonghun Park, Jinseok Choi, Namyoon Lee, Wonjae Shin, and H. Vincent Poor

arXiv:2108.06844v2 [eess.SP] 2 Jun 2022

## Abstract

Rate-splitting multiple access (RSMA) is a general multiple access scheme for downlink multi-antenna systems embracing both classical spatial division multiple access and more recent non-orthogonal multiple access. Finding a linear precoding strategy that maximizes the sum spectral efficiency of RSMA is a challenging yet significant problem. In this paper, we put forth a novel precoder design framework that jointly finds the linear precoders for the common and private messages for RSMA. Our approach is first to approximate the non-smooth minimum function part in the sum spectral efficiency of RSMA using a LogSumExp technique. Then, we reformulate the sum spectral efficiency maximization problem as a form of the log-sum of Rayleigh quotients to convert it into a tractable form. By interpreting the first-order optimality condition of the reformulated problem as an eigenvector-dependent nonlinear eigenvalue problem, we reveal that the leading eigenvector of the derived optimality condition is a local optimal solution. To find the leading eigenvector, we propose an algorithm inspired by a power iteration. Simulation results show that the proposed RSMA transmission strategy provides significant improvement in the sum spectral efficiency compared to the state-of-the-art RSMA transmission methods.

## Index Terms

Rate-splitting multiple access (RSMA), multi-user MIMO, imperfect channel state information (CSI), sum spectral efficiency maximization, generalized power iteration.

J. Park is with the School of Electronics Engineering, Kyungpook National University, South Korea (e-mail: jeonghun.park@knu.ac.kr). J. Choi is with Department of Electrical Engineering, Ulsan National Institute of Science and Technology, South Korea (e-mail: jinseokchoi@unist.ac.kr). N. Lee is with Department of Electrical Engineering, POSTECH, South Korea (e-mail: nylee@postech.ac.kr). W. Shin is with Department of Electrical and Computer Engineering, Ajou University, South Korea (email: wjshin@ajou.ac.kr). H. V. Poor is with Department of Electrical and Computer Engineering, Princeton University, Princeton, NJ, USA. (email: poor@princeton.edu)

# I. INTRODUCTION

Multi-user multiple-input multiple-output (MU-MIMO) downlink transmissions can provide extensive gains in spectral efficiency by serving multiple users with a shared time-frequency resource [2]–[4]. Assuming perfect channel state information at the transmitter (CSIT), a transmitter is able to send information symbols along with multiple linear precoding vectors to different users simultaneously by mitigating inter-user interference. In practice, however, the theoretical gains of downlink MU-MIMO transmissions can greatly vanish due to the inaccuracies of the CSIT. For example, considering the frequency division duplex (FDD) systems, the downlink channel has to be estimated at the receiver first and sent back to the transmitter via a finite-rate feedback link [5], [6], wherein quantization error on CSIT is inevitable. For this reason, in order to attain the de facto MU-MIMO spectral efficiency gains, it is crucial to design a downlink MU-MIMO transmission strategy that achieves high spectral efficiency under imperfect CSIT.

Rate-splitting multiple access (RSMA) is a robust downlink multiple access technique, especially when a transmitter has inaccurate knowledge for downlink CSI. Unlike the conventional spatial division multiple access (SDMA), in RSMA [7]–[10], the transmitter harnesses the rate-splitting strategy that breaks user messages into common and private parts in order to dynamically manage interference caused by imperfect CSIT. The transmitter constructs a common message by jointly encoding the common parts of the users' split messages. The rate for this common message is carefully controlled so that all the users can decode it. The transmitter also encodes the private parts of the users' messages to generate private information symbols. Then, the transmitter sends the common and private information symbols along with linear precoding vectors in a non-orthogonal manner. Each user decodes and eliminates the common message by performing successive interference cancellation (SIC) while treating the residual interference as noise. It then decodes the desired private message. Thanks to the rate-splitting encoding and SIC decoding, RSMA has been shown to outperform dirty paper coding (DPC) when imperfect CSIT is given [11].

To clearly understand the gains of RSMA over SDMA, it is instructive to consider a simple case of a two-user multi-antenna broadcast channel with imperfect CSIT. From an information-theoretic viewpoint, when applying linear precoding with imperfect CSIT in a two-user multi-

antenna broadcast channel, the channel can be interpreted as a virtual two-user interference channel with transmitter cooperation, in which the channel gains of desired and interfering links are determined by the precoding vectors and the channel vectors. In this equivalent interference channel, the quasi-optimal transmission strategy is the Han-Kobayashi scheme [12], i.e., splitting messages into common and private parts and allocating the power according to the relative channel gains between the interfering and the desired link [13]. Motivated by this, RSMA mimics this near capacity-achieving strategy in downlink MIMO.

To reap the spectral efficiency gains by using RSMA in downlink MIMO, it is significant to find the optimal linear precoding solution; yet it is challenging to find such a precoding vector. Unlike the sum spectral efficiency maximization problem for SDMA relying on private messages only, the problem for RSMA has an additional unique challenge induced by the common message rate, which is the minimum of all the achievable rates for the common message at the users. This minimum function is non-smooth, making the sum-spectral efficiency maximization problem for RSMA challenging to solve. In this work, we put forth a new approach for designing a precoder to maximize the sum spectral efficiency of multi-antenna RSMA with imperfect CSIT.

## A. Related Works

Recently, to cope with imperfect CSIT in downlink MIMO systems, the idea of rate-splitting has been actively re-explored as a multiple access technique, i.e., RSMA [14]. In [15], it was shown that RSMA provides sum degrees-of-freedom gains in multi-antenna broadcast channels where erroneous CSIT is given. Exploiting the idea of [15], in [16], the achievable spectral efficiency was analyzed while fixing the precoder for the common message as a random precoder and the precoder for the private messages as ZF.

Besides the theoretical analysis, there exist several prior works that have developed practical linear precoding designs for RSMA multi-antenna systems. In [7], a linear precoding design was proposed based on the weighted minimum mean square (WMMSE) approach [3]. Specifically, a non-convex original problem was transformed into a quadratically constrained quadratic program (QCQP) by using the equivalence between the sum spectral efficiency maximization and the sum mean square error (MSE) minimization problem. Subsequently, an interior point method was used to solve the QCQP. Employing the same idea, in [8], [17], a max-min fairness problem with RSMA was addressed. In [10], a linear precoding method for general RSMA was proposed by exploiting a concave-convex procedure (CCCP) that successively approximates the original

problem into convex forms. To evaluate the performance of RSMA in practical settings, e.g., finite constellation and implementable channel coding, [18] performed a link-level simulation in RSMA downlink MIMO systems. In [19], considering a single-antenna downlink channel, a power control method was proposed by incorporating the SIC constraint. In [20], considering downlink massive MIMO, it was shown that hierarchical RSMA that uses multiple-layer partial common messages can be well-harmonized in massive MIMO systems thanks to its spatial covariance separability [21]. Beyond the sum spectral efficiency maximization in downlink, other variants also exist. For example, RSMA for energy efficiency maximization [22], RSMA with hardware impairments [23], RSMA in joint MIMO radar and communication system [24], and RSMA in uplink channels [25] have been studied in the context of optimization for RSMA. Further, multi-antenna RSMA in interference channels [26] was also presented.

A key obstacle of the RSMA linear precoding design arises from the common message rate that should be determined as the minimum of all the achievable rates. To resolve this, the conventional methods use convex relaxation. Namely, an original non-convex problem is relaxed into a convex problem first, and then this convexified problem is put into an off-the-shelf optimization toolbox such as CVX to obtain a solution. A limitation of such approaches is that the optimization toolbox is hard to implement in practical hardware due to its extremely high complexity [27]. For this reason, the existing precoding optimization methods for RSMA are hardly used in practice. Therefore, this paper proposes a new optimization framework for MIMO RSMA that outperforms the existing methods in terms of complexity and also performance.

### B. Contributions

This paper proposes a new approach for linear precoding optimization in downlink MIMO with RSMA. The contributions of this paper are listed as follows.

- Considering an imperfect CSIT model, in which the CSI error statistic is modeled as complex Gaussian with zero-mean and a certain covariance matrix, we derive a lower bound on the instantaneous sum spectral efficiency for RSMA. In contrast to the sum spectral efficiency maximization for SDMA with imperfect CSIT, this lower bound entails the non-smooth minimum function for the common message rate. To convert the non-smooth minimum function into a tractable form, we take the LogSumExp technique, which offers a tight approximation of the minimum function to a smooth function. Then, by representing all optimization variables (precoding vectors) onto a higher dimensional vector, we reformulate

the lower bound of the instantaneous sum spectral efficiency for RSMA into a tractable non-convex function in the form of the log-sum of Rayleigh quotients.

- Using the derived lower bound with the smooth function approximation, we establish the first-order optimality condition for the sum spectral efficiency maximization problem. Remarkably, it is shown that the derived condition is cast as an eigenvector-dependent nonlinear eigenvalue problem [28], where the optimization variable behaves as an eigenvector, and the objective function behaves as an eigenvalue. Accordingly, we reveal that if we find the leading eigenvector that ensures the derived optimality condition, the best local optimal solution is obtained, maximizing the approximate lower bound of the instantaneous sum spectral efficiency for RSMA.

- To obtain the leading eigenvector of the derived condition, we put forth a novel algorithm inspired by a power iteration, referred to as generalized power iteration for rate-splitting (GPI-RS). Adopting the conventional power iteration principle, the idea of GPI-RS is to compute the leading eigenvector iteratively. The solution obtained by GPI-RS jointly provides the precoding directions and power allocation for the common and private messages. Notably, we do not rely on CVX in the proposed algorithm; thereby, it is more beneficial to implement in practical hardware. In addition to this, the computational complexity is less compared to the existing WMMSE-based method [7]. Later, we also generalize the proposed GPI-RS for a case in which multiple-layer RSMA is used. In multiple-layer RSMA, not only common message, but also the partial common message that includes messages of a subset of the users are jointly used. We show that the proposed method is suitably extended to this case.

- Simulation results show that the proposed GPI-RS provides spectral efficiency gains over the existing methods, including the conventional convex relaxation-based WMMSE method [7] in various system environments. To be specific, the proposed GPI-RS provides around 20% sum spectral efficiency gains, while consuming only $6 \sim 7\%$ of the computation time compared to the conventional method. Further, we empirically confirm that the GPI-RS converges well.

*Notation*: The superscripts $(\cdot)^{\mathsf{T}}$, $(\cdot)^{\mathsf{H}}$, and $(\cdot)^{-1}$ denote the transpose, Hermitian, and matrix inversion, respectively. $\mathbf{I}_N$ is the identity matrix of size $N \times N$, Assuming that $\mathbf{A}_1, ..., \mathbf{A}_N \in \mathbb{C}^{K \times K}$, $\mathbf{A} = \mathrm{blkdiag}\,(\mathbf{A}_1, ..., \mathbf{A}_n, ..., \mathbf{A}_N)$ is a block-diagonal matrix concatenating $\mathbf{A}_1, ..., \mathbf{A}_N$.

## II. SYSTEM MODEL

### A. Channel Model

We consider a single-cell downlink MU-MIMO system, where a base station (BS) equipped with $N$ antennas serves $K$ single-antenna users. We denote a user set as $\mathcal{K} = \{1, \cdots, K\}$. The channel vector between the BS and user $k$ is denoted as $\mathbf{h}_k \in \mathbb{C}^N$ for $k \in \mathcal{K}$, where $\mathbf{h}_k$ is generated based on the spatial covariance matrix $\mathbf{R}_k$, i.e., $\mathbf{R}_k = \mathbb{E}\left[\mathbf{h}_k \mathbf{h}_k^{\mathsf{H}}\right]$. For constructing the channel covariance matrix, we adopt the one-ring model [21]. Specifically, we assume that the BS is equipped with uniform circular array with radius $\psi D$ where $\psi$ denotes a signal wavelength and $D = \frac{0.5}{\sqrt{(1-\cos(2\pi/N))^2+\sin^2(2\pi/N)}}$. Then, the channel correlation coefficient between the $n$-th antenna and $m$-th antenna corresponding to user $k$ is defined as

$$[\mathbf{R}_k]_{n,m} = \frac{1}{2\Delta_k} \int_{\theta_k-\Delta_k}^{\theta_k+\Delta_k} e^{-j\frac{2\pi}{\psi}\Psi(x)(\mathbf{r}_n-\mathbf{r}_m)}\mathrm{d}x, \tag{1}$$

where $\theta_k$ is angle-of-arrival (AoA) of user $k$, $\Delta_k$ is the angular spread of user $k$, $\Psi(x) = [\cos(x), \sin(x)]$, and $\mathbf{r}_n$ is the position vector of the $n$-th antenna. By employing the Karhunen-Loeve model as in [20], [21], the channel vector $\mathbf{h}_k$ is represented as

$$\mathbf{h}_k = \mathbf{U}_k \Lambda_k^{\frac{1}{2}} \mathbf{g}_k, \tag{2}$$

where $\Lambda_k \in \mathbb{C}^{r_k \times r_k}$ is a diagonal matrix that contains the non-zero eigenvalues of $\mathbf{R}_k$, $\mathbf{U}_k \in \mathbb{C}^{N \times r_k}$ is a collection of the eigenvectors of $\mathbf{R}_k$ corresponding to the eigenvalues in $\Lambda_k$, and $\mathbf{g}_k \in \mathbb{C}^{r_k}$ is an independent and identically distributed channel vector. We assume that each element of $\mathbf{g}_k$ is drawn from $\mathcal{CN}(0, 1)$. We consider a block fading model, where $\mathbf{g}_k$ keeps constant within one transmission block. Over the two consecutive transmission blocks, $\mathbf{g}_k$ changes independently.

We clarify that the applicability of our method does not depend on particular channel model assumptions. We consider the one-ring model in this paper because it is one of the widely used channel models that suitably captures spatial covariance structures of MIMO channels. The proposed method can be applied in any channel model.

### B. CSIT Acquisition Model

This subsection explains the CSIT estimation and the error model used throughout this paper. We assume perfect channel state information at the receiver (CSIR), which can be achieved via downlink pilots planted in the data packet as described in LTE and 5G NR. In contrast to CSIR, a BS should estimate CSIT, so that only imperfect knowledge of CSIT is allowed. Generally,

two approaches are known to estimate the CSIT, each of which is linear MMSE (LMMSE) and limited feedback, respectively. LMMSE yields the optimum CSIT estimation performance, provided that the channel is distributed as Gaussian. Nonetheless, LMMSE only can be used when the channel reciprocity holds. On the contrary to that, limited feedback can be employed in any environment. In this paper, we focus on LMMSE, but note that our method is also useful with limited feedback.

As mentioned above, LMMSE can be exploited when the BS can use channel reciprocity. Specifically, assuming that the uplink and downlink channels are reciprocal, the BS estimates the CSIT from uplink training sent from the users using the MMSE estimation [29]. For this reason, LMMSE is adequate to use in time division duplex (TDD) systems where the channel reciprocity holds. Using LMMSE, the estimated CSIT is presented as

$$\hat{\mathbf{h}}_k = \mathbf{h}_k - \mathbf{e}_k, \tag{3}$$

where $\mathbf{e}_k$ is the CSIT estimation error vector. Since $\mathbf{h}_k$ is distributed as Gaussian, $\hat{\mathbf{h}}_k$ and $\mathbf{e}_k$ are also Gaussian that is independent to each other. The error covariance is obtained as

$$\mathbb{E}[\mathbf{e}_k \mathbf{e}_k^\mathsf{H}] = \boldsymbol{\Phi}_k = \mathbf{R}_k - \mathbf{R}_k \left( \mathbf{R}_k + \frac{\sigma^2}{\tau_{\mathsf{ul}} p_{\mathsf{ul}}} \right)^{-1} \mathbf{R}_k, \tag{4}$$

where $\tau_{\mathsf{ul}}$ and $p_{\mathsf{ul}}$ are uplink training length and uplink training transmit power. As the uplink training length and power increases to infinity, the error covariance $\boldsymbol{\Phi}_k = \mathbf{0}$ and the CSIT error $\mathbf{e}_k$ also vanishes; then the BS has the perfect CSIT.

### C. RSMA Signal Model

Using RSMA, the message for user $k$ is split into the common message part $s_{\mathsf{c},k}$ and the private message $s_k$. The common message part $s_{\mathsf{c},k}$ from each user is combined to encode the common message $s_{\mathsf{c}}$. The common message $s_{\mathsf{c}}$ is drawn from a public codebook so that any user associated with the BS can decode it. On the contrary, the private message $s_k$ comes from an individual codebook. Therefore it is only decodable to intended users.

One common message and $K$ private messages are linearly precoded and then superimposed, so that the transmit signal $\mathbf{x} \in \mathbb{C}^N$ is given by

$$\mathbf{x} = \mathbf{f}_{\mathsf{c}} s_{\mathsf{c}} + \sum_{i=1}^{K} \mathbf{f}_i s_i, \tag{5}$$

where $\mathbf{f_c} \in \mathbb{C}^N$ and $\mathbf{f}_i \in \mathbb{C}^N$ are the precoding vectors for the common and private messages respectively with the transmit power constraint: $\|\mathbf{f_c}\|^2 + \sum_{i=1}^{K} \|\mathbf{f}_i\|^2 \le 1$. We note that not only the direction of each message, but also the power allocated to each message are controlled by the precoding vectors. For example, if $\|\mathbf{f_c}\| = 0$, then no common message is delivered; so our RSMA signal model reduces to typical SDMA.

The received signal at user $k$ for $k \in \mathcal{K}$ is written as

$$y_k = \mathbf{h}_k^{\mathsf{H}} \mathbf{f_c} s_c + \mathbf{h}_k^{\mathsf{H}} \mathbf{f}_k s_k + \sum_{\ell=1,\ell \neq k}^{K} \mathbf{h}_k^{\mathsf{H}} \mathbf{f}_\ell s_\ell + z_k, \tag{6}$$

where $z_k \sim \mathcal{CN}(0, \sigma^2)$ is additive white Gaussian noise. We also assume that $s_c$ and $s_k$ are drawn from an independent Gaussian codebook, i.e., $s_c, s_k \sim \mathcal{CN}(0, P)$.

### D. Performance Characterization

To characterize the performance of RSMA, we first explain the decoding process performed by each user. Each user first decodes the common message $s_c$ by treating all the other private messages as noise. Once the common message is successfully decoded, using SIC, the users remove the common message from the received signal and decode the private messages with a reduced amount of interference.

To successfully perform SIC, the common message $s_c$ should be decodable to every user without any error. To this end, the code rate of the common message $s_c$ is set as the minimum of the ergodic spectral efficiencies among the users. Accordingly, under the premise that the BS has imperfect CSIT, i.e., $\hat{\mathbf{h}}_k = \mathbf{h}_k + \mathbf{e}_k$ for $k \in \mathcal{K}$, the ergodic spectral efficiency of the common message is obtained as [4], [7]

$$\begin{aligned} R_c &= \min_{k \in \mathcal{K}} \left\{ \mathbb{E}_{\{\mathbf{h}_k\}} \left[ \log_2 \left( 1 + \frac{|\mathbf{h}_k^{\mathsf{H}} \mathbf{f_c}|^2}{\sum_{\ell=1}^{K} |\mathbf{h}_k^{\mathsf{H}} \mathbf{f}_\ell|^2 + \sigma^2/P} \right) \right] \right\} \\ &= \min_{k \in \mathcal{K}} \left\{ \mathbb{E}_{\{\hat{\mathbf{h}}_k\}} \left[ \mathbb{E}_{\{\mathbf{e}_k\}} \left[ \log_2 \left( 1 + \frac{|\mathbf{h}_k^{\mathsf{H}} \mathbf{f_c}|^2}{\sum_{\ell=1}^{K} |\mathbf{h}_k^{\mathsf{H}} \mathbf{f}_\ell|^2 + \sigma^2/P} \right) \bigg| \hat{\mathbf{h}}_k \right] \right] \right\}, \end{aligned} \tag{7}$$

where in (7), the inner expectation is taken over the randomness associated with the CSIT error $(\mathbb{E}_{\{\mathbf{e}_k\}}[\cdot])$ within one particular coherence block and the outer expectation is taken over the randomness associated with the imperfect knowledge of the channel fading process $(\mathbb{E}_{\{\hat{\mathbf{h}}_k\}}[\cdot])$. Assuming that the channel code length spans an infinite number of channel blocks and we set the channel coding rate of the common message $s_c$ less than or equal to $R_c$, no decoding error for $s_c$

occurs. Similar to (7), the ergodic spectral efficiency of the private message $s_k$ after cancelling the common message $s_c$ is obtained as [4], [7]

$$R_k = \mathbb{E}_{\{\hat{\mathbf{h}}_k\}}\left[\mathbb{E}_{\{\mathbf{e}_k\}}\left[\log_2\left(1 + \frac{|\mathbf{h}_k^{\mathsf{H}}\mathbf{f}_k|^2}{\sum_{\ell=1,\ell\neq k}^{K}|\mathbf{h}_k^{\mathsf{H}}\mathbf{f}_\ell|^2 + \sigma^2/P}\right)\bigg|\hat{\mathbf{h}}_k\right]\right]. \tag{8}$$

Since we assume that the common message is successfully eliminated, we observe that there is no interference from the common message in (8). Under the assumption that the channel code length spans an infinite number of channel blocks and the channel coding rate of the private message $s_k$ is less than or equal to $R_k$, the users can successfully decode $s_k$.

Our main goal is to optimize the precoders using imperfect knowledge on CSIT per each fading block. For this reason, we focus on one particular fading block without loss of generality, allowing to assume that $\hat{\mathbf{h}}_k$, $k \in \mathcal{K}$ is given. In a certain fading block, we can define the instantaneous spectral efficiency. Specifically, the instantaneous spectral efficiency of the common message achieved at user $k$ is defined as

$$R_c^{\mathsf{ins.}}(k) = \mathbb{E}_{\{\mathbf{e}_k\}}\left[\log_2\left(1 + \frac{|\mathbf{h}_k^{\mathsf{H}}\mathbf{f}_c|^2}{\sum_{\ell=1}^{K}|\mathbf{h}_k^{\mathsf{H}}\mathbf{f}_\ell|^2 + \sigma^2/P}\right)\bigg|\hat{\mathbf{h}}_k\right]. \tag{9}$$

The instantaneous spectral efficiency differs from the ergodic spectral efficiency. On the one hand, the ergodic spectral efficiency is the long-term performance that can be achieved when the channel code length spans very long channel blocks. On the other hand, the instantaneous spectral efficiency is the short-term rate expression when taking into account the channel estimation error effect per channel realization. Considering multiple fading blocks, the instantaneous spectral efficiency and the ergodic spectral efficiency are connected as $R_c = \min_{k\in\mathcal{K}}\left\{\mathbb{E}_{\{\hat{\mathbf{h}}_k\}}[R_c^{\mathsf{ins.}}(k)]\right\}$.

Unfortunately, however, (9) is not tractable. The main challenge is that no closed-form exists for the expectation on CSIT error. To address this, we characterize a lower bound by adopting a similar approach in [4]. We rewrite the received signal (6) with the CSIT error term as follows:

$$\begin{aligned}y_k &= \mathbf{h}_k^{\mathsf{H}}\mathbf{f}_c s_c + \sum_{\ell=1}^{K}\mathbf{h}_k^{\mathsf{H}}\mathbf{f}_\ell s_k + z_k \\ &\overset{(a)}{=} \hat{\mathbf{h}}_k^{\mathsf{H}}\mathbf{f}_c s_c + \sum_{\ell=1}^{K}\hat{\mathbf{h}}_k^{\mathsf{H}}\mathbf{f}_\ell s_k + \underbrace{\mathbf{e}_k^{\mathsf{H}}\mathbf{f}_c s_c + \sum_{i=1}^{K}\mathbf{e}_k^{\mathsf{H}}\mathbf{f}_i s_i}_{(b)} + z_k,\end{aligned} \tag{10}$$

where (a) follows $\mathbf{h}_k = \hat{\mathbf{h}}_k + \mathbf{e}_k$. The term (b) is correlated with the common message $s_c$, yet it is not tractable due to the CSIT estimation error $\mathbf{e}_k$. To resolve this, employing a concept of generalized mutual information, we treat (b) as independent Gaussian noise, which is the worst

case of mutual information. Then a lower bound on the instantaneous spectral efficiency is made as:

$$
\begin{aligned}
R_{\mathsf{C}}^{\mathsf{ins.}}(k) &\overset{(c)}{\geq} \mathbb{E}_{\{\mathbf{e}_k\}}\left[\log_2\left(1 + \frac{|\hat{\mathbf{h}}_k^{\mathsf{H}}\mathbf{f}_{\mathsf{C}}|^2}{\sum_{\ell=1}^{K}|\hat{\mathbf{h}}_k^{\mathsf{H}}\mathbf{f}_\ell|^2 + |\mathbf{e}_k^{\mathsf{H}}\mathbf{f}_{\mathsf{C}}|^2 + \sum_{\ell=1}^{K}|\mathbf{e}_k^{\mathsf{H}}\mathbf{f}_\ell|^2 + \frac{\sigma^2}{P}}\right)\right] \\
&\overset{(d)}{\geq} \log_2\left(1 + \frac{|\hat{\mathbf{h}}_k^{\mathsf{H}}\mathbf{f}_{\mathsf{C}}|^2}{\sum_{\ell=1}^{K}|\hat{\mathbf{h}}_k^{\mathsf{H}}\mathbf{f}_\ell|^2 + \mathbf{f}_{\mathsf{C}}^{\mathsf{H}}\mathbb{E}\left[\mathbf{e}_k\mathbf{e}_k^{\mathsf{H}}\right]\mathbf{f}_{\mathsf{C}} + \sum_{\ell=1}^{K}\mathbf{f}_\ell^{\mathsf{H}}\mathbb{E}\left[\mathbf{e}_k\mathbf{e}_k^{\mathsf{H}}\right]\mathbf{f}_\ell + \frac{\sigma^2}{P}}\right) \\
&\overset{(e)}{=} \log_2\left(1 + \frac{|\hat{\mathbf{h}}_k^{\mathsf{H}}\mathbf{f}_{\mathsf{C}}|^2}{\sum_{\ell=1}^{K}|\hat{\mathbf{h}}_k^{\mathsf{H}}\mathbf{f}_\ell|^2 + \mathbf{f}_{\mathsf{C}}^{\mathsf{H}}\boldsymbol{\Phi}_k\mathbf{f}_{\mathsf{C}} + \sum_{\ell=1}^{K}\mathbf{f}_\ell^{\mathsf{H}}\boldsymbol{\Phi}_k\mathbf{f}_\ell + \frac{\sigma^2}{P}}\right) = \bar{R}_{\mathsf{C}}^{\mathsf{ins.}}(k), \quad (11)
\end{aligned}
$$

where (c) comes from treating (b) in (10) as independent Gaussian noise, (d) follows Jensen's inequality, and (e) comes from the CSIT error covariance $\mathbb{E}[\mathbf{e}_k\mathbf{e}_k^{\mathsf{H}}] = \boldsymbol{\Phi}_k$. A lower bound on the instantaneous spectral efficiency, denoted as $\bar{R}_{\mathsf{C}}^{\mathsf{ins.}}(k)$, is derived as a closed-form, so that this can be handled easily. Finally, we take another lower bound on the ergodic spectral efficiency. The ergodic spectral efficiency of the common message is represented with $\bar{R}_{\mathsf{C}}^{\mathsf{ins.}}(k)$ as

$$
\begin{aligned}
R_{\mathsf{C}} = \min_{k\in\mathcal{K}}\left\{\mathbb{E}_{\{\hat{\mathbf{h}}_k\}}\left[R_{\mathsf{C}}^{\mathsf{ins.}}(k)\right]\right\} &\geq \min_{k\in\mathcal{K}}\left\{\mathbb{E}_{\{\hat{\mathbf{h}}_k\}}\left[\bar{R}_{\mathsf{C}}^{\mathsf{ins.}}(k)\right]\right\} \\
&\overset{(f)}{\geq} \mathbb{E}_{\{\hat{\mathbf{h}}_{k\in\mathcal{K}}\}}\left[\min_{k\in\mathcal{K}}\left\{\bar{R}_{\mathsf{C}}^{\mathsf{ins.}}(k)\right\}\right] \quad (12)
\end{aligned}
$$

where (f) follows the fact that putting the minimum operator into the expectation does not increase the value. We take (12) as our main objective function for the common message rate.

Next, we characterize a lower bound on the instantaneous spectral efficiency of the private message. We first define the instantaneous spectral efficiency of the private message $s_k$ in a certain fading block as $R_k^{\mathsf{ins.}}$. Using a similar technique to the common message case, we derive a lower bound on $R_k^{\mathsf{ins.}}$ such as

$$
\begin{aligned}
R_k^{\mathsf{ins.}} &\geq \mathbb{E}_{\{\mathbf{e}_k\}}\left[\log_2\left(1 + \frac{|\hat{\mathbf{h}}_k^{\mathsf{H}}\mathbf{f}_k|^2}{\sum_{\ell=1,\ell\neq k}^{K}|\hat{\mathbf{h}}_k^{\mathsf{H}}\mathbf{f}_\ell|^2 + \sum_{\ell=1}^{K}|\mathbf{e}_k^{\mathsf{H}}\mathbf{f}_\ell|^2 + \frac{\sigma^2}{P}}\right)\right] \\
&\geq \log_2\left(1 + \frac{|\hat{\mathbf{h}}_k^{\mathsf{H}}\mathbf{f}_k|^2}{\sum_{\ell=1,\ell\neq k}^{K}|\hat{\mathbf{h}}_k^{\mathsf{H}}\mathbf{f}_\ell|^2 + \sum_{\ell=1}^{K}\mathbf{f}_\ell^{\mathsf{H}}\boldsymbol{\Phi}_k\mathbf{f}_\ell + \frac{\sigma^2}{P}}\right) = \bar{R}_k^{\mathsf{ins.}}, \quad (13)
\end{aligned}
$$

Considering multiple fading blocks, the obtained lower bound on the instantaneous spectral efficiency $\bar{R}_k^{\mathsf{ins.}}$ is connected to the ergodic spectral efficiency as follows:

$$
R_k = \mathbb{E}_{\{\hat{\mathbf{h}}_k\}}[R_k^{\mathsf{ins.}}] \geq \mathbb{E}_{\{\hat{\mathbf{h}}_k\}}[\bar{R}_k^{\mathsf{ins.}}]. \quad (14)
$$

Combining (12) and (14), we finally complete the following lower bound on the ergodic sum spectral efficiency $R_\Sigma$:

$$R_\Sigma \geq \bar{R}_\Sigma = \mathbb{E}_{\{\hat{\mathbf{h}}_{k\in\mathcal{K}}\}}\left[\min_{k\in\mathcal{K}}\{\bar{R}_{\mathsf{c}}^{\mathsf{ins.}}(k)\} + \sum_{k=1}^{K}\bar{R}_k^{\mathsf{ins.}}\right]. \tag{15}$$

Now we are ready to formulate our main problem.

**Remark 1** (Comparison to [10]). Similar to our lower bound, [10] also proposed a lower bound on the instantaneous spectral efficiencies by incorporating the CSIT estimation error. To be specific, [10] assumed a particular scenario of CSIT estimation that $\mathbf{h}_k = \hat{\mathbf{h}}_k + \delta\mathbf{e}_k$, where $\mathbb{E}[\mathbf{e}_k^{\mathsf{H}}\mathbf{e}_k] = \mathbf{I}$ and $\mathbb{E}[\mathbf{e}_k] = 0$. We note that this is a special case of our CSIT estimation model. Under this premise, the instantaneous spectral efficiency of the common message achieved at user $k$, $R_{\mathsf{c}}^{\mathsf{ins.}}(k)$, is expressed as

$$
\begin{aligned}
R_{\mathsf{c}}^{\mathsf{ins.}}(k) &= \mathbb{E}_{\{\mathbf{e}_k\}}\left[\left.\log_2\left(1 + \frac{\mathbf{h}_k^{\mathsf{H}}\mathbf{Q}_{\mathsf{c}}\mathbf{h}_k}{\sum_{\ell=1}^{K}\mathbf{h}_k^{\mathsf{H}}\mathbf{Q}_\ell\mathbf{h}_k + \sigma^2/P}\right)\right| \hat{\mathbf{h}}_k\right] \\
&= \mathbb{E}_{\{\mathbf{e}_k\}}\left[\log_2\left(\mathbf{h}_k^{\mathsf{H}}\mathbf{Q}_{\mathsf{c}}\mathbf{h}_k + \sum_{\ell=1}^{K}\mathbf{h}_k^{\mathsf{H}}\mathbf{Q}_\ell\mathbf{h}_k + \sigma^2/P\right)\right] - \mathbb{E}_{\{\mathbf{e}_k\}}\left[\log_2\left(\sum_{\ell=1}^{K}\mathbf{h}_k^{\mathsf{H}}\mathbf{Q}_\ell\mathbf{h}_k + \sigma^2/P\right)\right],
\end{aligned}
\tag{16}
$$

where $\mathbf{Q}_{\mathsf{c}} = \mathbf{f}_{\mathsf{c}}\mathbf{f}_{\mathsf{c}}^{\mathsf{H}}$ and $\mathbf{Q}_k = \mathbf{f}_k\mathbf{f}_k^{\mathsf{H}}$. By using Jensen's inequality and the fact that the CSIT estimation error has zero-mean, we obtain an upper bound on the second term in (16) as

$$\mathbb{E}_{\{\mathbf{e}_k\}}\left[\log_2\left(\sum_{\ell=1}^{K}\mathbf{h}_k^{\mathsf{H}}\mathbf{Q}_\ell\mathbf{h}_k + \sigma^2/P\right)\right] \leq \log_2\left(\sum_{\ell=1}^{K}(\hat{\mathbf{h}}_k^{\mathsf{H}}\mathbf{Q}_\ell\hat{\mathbf{h}}_k + \delta^2\mathsf{tr}(\mathbf{Q}_\ell)) + \sigma^2/P\right) \tag{17}$$

Additionally, [10] proved that the first term in (16) is non-decreasing with $\delta$, so we get a lower bound by putting $\delta = 0$, yielding

$$\mathbb{E}_{\{\mathbf{e}_k\}}\left[\log_2\left(\mathbf{h}_k^{\mathsf{H}}\mathbf{Q}_{\mathsf{c}}\mathbf{h}_k + \sum_{\ell=1}^{K}\mathbf{h}_k^{\mathsf{H}}\mathbf{Q}_\ell\mathbf{h}_k + \sigma^2/P\right)\right] \geq \log_2\left(\hat{\mathbf{h}}_k^{\mathsf{H}}\mathbf{Q}_{\mathsf{c}}\hat{\mathbf{h}}_k + \sum_{\ell=1}^{K}\hat{\mathbf{h}}_k^{\mathsf{H}}\mathbf{Q}_\ell\hat{\mathbf{h}}_k + \sigma^2/P\right). \tag{18}$$

Combining (17) and (18), [10] claimed that we can make a lower bound as

$$R_{\mathsf{c}}^{\mathsf{ins.}}(k) \geq \log_2\left(1 + \frac{\hat{\mathbf{h}}_k^{\mathsf{H}}\mathbf{Q}_{\mathsf{c}}\hat{\mathbf{h}}_k}{\sum_{\ell=1}^{K}(\hat{\mathbf{h}}_k^{\mathsf{H}}\mathbf{Q}_\ell\hat{\mathbf{h}}_k + \delta^2\mathsf{tr}(\mathbf{Q}_\ell)) + \sigma^2/P}\right) = \tilde{R}_{\mathsf{c}}^{\mathsf{ins.}}(k). \tag{19}$$

The lower bound (19) is closely related to our lower bound. Rewriting (19), we have

$$\tilde{R}_{\mathsf{c}}^{\mathsf{ins.}}(k) = \log_2\left(1 + \frac{|\hat{\mathbf{h}}_k^{\mathsf{H}}\mathbf{f}_{\mathsf{c}}|^2}{\sum_{\ell=1}^{K}|\hat{\mathbf{h}}_k^{\mathsf{H}}\mathbf{f}_\ell|^2 + \sum_{\ell=1}^{K}\mathbf{f}_\ell^{\mathsf{H}}\cdot\delta^2\mathbf{I}\cdot\mathbf{f}_\ell + \sigma^2/P}\right). \tag{20}$$

Comparing $\tilde{R}_{\mathsf{c}}^{\mathsf{ins.}}(k)$ in (20) to our lower bound $\bar{R}_{\mathsf{c}}^{\mathsf{ins.}}(k)$ in (11) under the assumption that $\mathbf{\Phi}_k = \delta^2 \mathbf{I}$, (20) seems to be tighter since the denominator of the SINR in (20) does not include the term $\mathbf{f}_{\mathsf{c}}^{\mathsf{H}} \cdot \delta^2 \mathbf{I} \cdot \mathbf{f}_{\mathsf{c}}$. Nonetheless, the lower bound (20) is limited in its applicability. This is because, the lower bound technique to derive (20) cannot be applied when $\mathbf{\Phi}_k \neq \delta^2 \mathbf{I}$, i.e., the CSIT error is spatially correlated. Since it is usual that MIMO channels have particular spatial correlation structures, our lower bound is more proper to use in general cases.

*E. Problem Formulation*

We aim to maximize $\bar{R}_{\Sigma}$ in (15). In our setup, maximizing $\bar{R}_{\Sigma}$ is equivalent to maximizing $\min_{k \in \mathcal{K}} \{\bar{R}_{\mathsf{c}}^{\mathsf{ins.}}(k)\} + \sum_{k=1}^{K} \bar{R}_k^{\mathsf{ins.}}$ per each fading block, wherein the BS is able to calculate $\bar{R}_k^{\mathsf{ins.}}$ and $\min_{k \in \mathcal{K}} \{\bar{R}_{\mathsf{c}}^{\mathsf{ins.}}(k)\}$ in a closed-form by using the estimated CSIT. Accordingly, we formulate an optimization problem as follows:

$$\underset{\mathbf{f}_{\mathsf{c}}, \mathbf{f}_1, \cdots, \mathbf{f}_K}{\text{maximize}} \quad \min_{k \in \mathcal{K}} \{\bar{R}_{\mathsf{c}}^{\mathsf{ins.}}(k)\} + \sum_{k=1}^{K} \bar{R}_k^{\mathsf{ins.}} \tag{21}$$

$$\text{subject to} \quad \|\mathbf{f}_{\mathsf{c}}\|^2 + \sum_{k=1}^{K} \|\mathbf{f}_k\|^2 \leq 1. \tag{22}$$

We tackle (21) as our main problem. Finding the global solution of (21) is infeasible due to its non-convexity and non-smoothness.

### III. EXISTING APPROACH: WMMSE

In this section, we briefly introduce the existing WMMSE approach [7] for the sum spectral efficiency maximization in an RSMA scenario. We focus on two points: *(i)* how to solve a non-convex optimization problem and *(ii)* how to incorporate the CSIT estimation error. Then we explain the main distinguishable points of the proposed method.

We first explain how to relax a non-convex problem. To solve a non-convex sum spectral efficiency maximization problem, [7] adopted a well-known WMMSE relaxation technique [3]. Specifically, we denote that $\hat{s}_{\mathsf{c}}(k)$, $\hat{s}_k$ are estimates for $s_{\mathsf{c}}(k)$ and $s_k$, where $s_{\mathsf{c}}(k)$ is the common message received at user $k$. With scalar equalizers $g_{\mathsf{c}}(k)$ and $g_k$, we define the MSEs of the common message received at user $k$ ($\epsilon_{\mathsf{c}}(k)$) and the private message ($\epsilon_k$) as

$$\epsilon_{\mathsf{c}}(k) = \mathbb{E}[|\hat{s}_{\mathsf{c}}(k) - s_{\mathsf{c}}|^2] = \mathbb{E}[|g_{\mathsf{c}}(k)y_k - s_{\mathsf{c}}|^2]$$

$$= |g_{\mathsf{c}}(k)|^2 T_{\mathsf{c}}(k) - 2\mathsf{Re}\left\{g_{\mathsf{c}}(k)\mathbf{h}_k^{\mathsf{H}}\mathbf{f}_{\mathsf{c}}\right\} + 1, \tag{23}$$

$$\epsilon_k = |g_k|^2 T_k - 2\mathsf{Re}\left\{g_k \mathbf{h}_k^{\mathsf{H}}\mathbf{f}_k\right\} + 1, \tag{24}$$

where $T_{\mathsf{c}}(k) = |\mathbf{h}_k^{\mathsf{H}}\mathbf{f}_{\mathsf{c}}|^2 + \sum_{\ell=1}^{K} |\mathbf{h}_k^{\mathsf{H}}\mathbf{f}_\ell|^2 + \sigma^2$ and $T_k = |\mathbf{h}_k^{\mathsf{H}}\mathbf{f}_k|^2 + \sum_{\ell=1,\ell\neq k}^{K} |\mathbf{h}_k^{\mathsf{H}}\mathbf{f}_\ell|^2 + \sigma^2$. Note that here we exchange the power assumption between the message and the precoding vectors, i.e., $\mathbb{E}[|s_{\mathsf{c}}|^2] = \mathbb{E}[|s_k|^2] = 1$ and $\|\mathbf{f}_{\mathsf{c}}\|^2 + \sum_{k=1}^{K} \|\mathbf{f}_k\|^2 \leq P$ for ease of description. This does not change the SINR. The minimum MSEs are achieved when $g_{\mathsf{c}}(k) = \mathbf{f}_{\mathsf{c}}^{\mathsf{H}}\mathbf{h}_k T_{\mathsf{c}}(k)^{-1}$ and $g_k = \mathbf{f}_k^{\mathsf{H}}\mathbf{h}_k T_k^{-1}$, providing the following MMSE: $\epsilon_{\mathsf{c}}^{\mathsf{MMSE}}(k) = T_{\mathsf{c}}(k)^{-1}(T_{\mathsf{c}}(k) - |\mathbf{h}_k^{\mathsf{H}}\mathbf{f}_{\mathsf{c}}|^2)$ and $\epsilon_k^{\mathsf{MMSE}} = T_k^{-1}(T_k - |\mathbf{h}_k^{\mathsf{H}}\mathbf{f}_k|^2)$.

Then the augmented WMSEs are give by

$$\xi_{\mathsf{c}}(k) = u_{\mathsf{c}}(k)\epsilon_{\mathsf{c}}(k) - \log_2(u_{\mathsf{c}}(k))$$

$$= \mathbf{f}_{\mathsf{c}}^{\mathsf{H}}\left(u_{\mathsf{c}}(k)|g_{\mathsf{c}}(k)|^2\mathbf{h}_k\mathbf{h}_k^{\mathsf{H}}\right)\mathbf{f}_{\mathsf{c}} + \sum_{\ell=1}^{K} \mathbf{f}_\ell^{\mathsf{H}}\left(u_{\mathsf{c}}(k)|g_{\mathsf{c}}(k)|^2\mathbf{h}_k\mathbf{h}_k^{\mathsf{H}}\right)\mathbf{f}_\ell - 2\mathsf{Re}\left\{u_{\mathsf{c}}(k)g_{\mathsf{c}}(k)\mathbf{h}_k^{\mathsf{H}}\mathbf{f}_{\mathsf{c}}\right\}$$

$$+ \sigma^2 u_{\mathsf{c}}(k)|g_{\mathsf{c}}(k)|^2 + u_{\mathsf{c}}(k) - \log_2(u_{\mathsf{c}}(k)), \tag{25}$$

$$\xi_k = u_k\epsilon_k - \log_2(u_k)$$

$$= \mathbf{f}_k^{\mathsf{H}}\left(u_k|g_k|^2\mathbf{h}_k\mathbf{h}_k^{\mathsf{H}}\right)\mathbf{f}_k + \sum_{\ell=1,\ell\neq k}^{K} \mathbf{f}_\ell^{\mathsf{H}}\left(u_k|g_k|^2\mathbf{h}_k\mathbf{h}_k^{\mathsf{H}}\right)\mathbf{f}_\ell - 2\mathsf{Re}\left\{u_k g_k\mathbf{h}_k^{\mathsf{H}}\mathbf{f}_k\right\}$$

$$+ \sigma^2 u_k|g_k|^2 + u_k - \log_2(u_k). \tag{26}$$

We note that the optimal weights to achieve the minimum of $\xi_{\mathsf{c}}(k)$ and $\xi_k$ are obtained as $u_{\mathsf{c}}(k) = 1/\epsilon_{\mathsf{c}}^{\mathsf{MMSE}}(k)$ and $u_k = 1/\epsilon_k^{\mathsf{MMSE}}$. Upon this weight update, by the rate-WMMSE equivalence [3], the sum spectral efficiency is maximized by solving the following WMSE minimization problem:

$$\underset{\mathbf{f}_{\mathsf{c}},\mathbf{f}_1,\cdots,\mathbf{f}_K,\xi_{\mathsf{c}}}{\text{minimize}} \; \xi_{\mathsf{c}} + \sum_{k=1}^{K} \xi_k \tag{27}$$

$$\text{subject to} \; \xi_{\mathsf{c}}(k) \leq \xi_{\mathsf{c}}, \; \forall k \in \mathcal{K}, \tag{28}$$

$$\|\mathbf{f}_{\mathsf{c}}\|^2 + \sum_{k=1}^{K} \|\mathbf{f}_k\|^2 \leq P. \tag{29}$$

The problem (27) is QCQP, which can be solved by using CVX. We compute the optimal weight, equalizer, and the precoding vector in an alternating fashion. We repeat this process until a certain termination criterion is met.

The presented WMMSE approach assumes the perfect CSIT. To take the CSIT estimation error into account, [7] adopted the sample average approximation (SAA) technique. In the SAA technique, we produce $M$ number of samples of the augmented WMSEs by randomly generating

channel vector $\mathbf{h} = \hat{\mathbf{h}}_k + \mathbf{e}_k$ ($\hat{\mathbf{h}}_k$ is given, $\mathbf{e}_k$ is randomly generated) then calculate an empirical average of these samples as follows:

$$\bar{\xi}_{\mathsf{c}}(k) \leftarrow \frac{1}{M} \sum_{m=1}^{M} \xi_{\mathsf{c}}^{(m)}(k), \tag{30}$$

$$\bar{\xi}_k \leftarrow \frac{1}{M} \sum_{m=1}^{M} \xi_{\mathsf{c}}^{(m)}(k). \tag{31}$$

We replace $\xi_{\mathsf{c}}(k)$ and $\xi_k$ by the empirical average augmented WMSE $\bar{\xi}_{\mathsf{c}}(k)$ and $\bar{\xi}_k$ in (27).

Now we clarify the distinguishable points of our method compared to the WMMSE approach. In the WMMSE approach, CVX is required to solve the WMSE minimization problem (27). We need additional efforts to implement CVX in FPGA hardware since CVX is not designed to run in real-time hardware [27]. In addition to this, we observe that the common message rate is controlled by the WMSE constraints (28), where the number of the constraints is equal to the number of users $K$. For this reason, the associated computational complexity of the WMMSE approach scales with $K^{3.5}$ [8], [30]; resulting in the huge computational complexity is caused when there are a large number of users. Compared to this, as we will show later, our method does not require CVX to obtain a solution. Further, since we control the common message rate by cleverly approximating the minimum function, the computational complexity scales with $K$. For this reason, our method is much more beneficial when there are many users.

## IV. PRECODER OPTIMIZATION WITH GENERALIZED POWER ITERATION

In this section, we explain the key ideas to solve the optimization problem (21). We first approximate the non-smooth minimum function as a smooth function using the LogSumExp technique. Subsequently, we represent the optimization variable onto a higher dimensional vector to reformulate the problem (21) into a tractable non-convex optimization problem expressed as a function of Rayleigh quotients. By deriving the first-order KKT condition for the reformulated problem, we show that the first-order optimality condition is cast as an eigenvector-dependent nonlinear eigenvalue problem (NEPv) [28], and finding the leading eigenvector is equivalent to finding the best local optimal point of the reformulated problem. Consequently, to find the leading eigenvector, we propose a computationally efficient generalized power iteration algorithm.
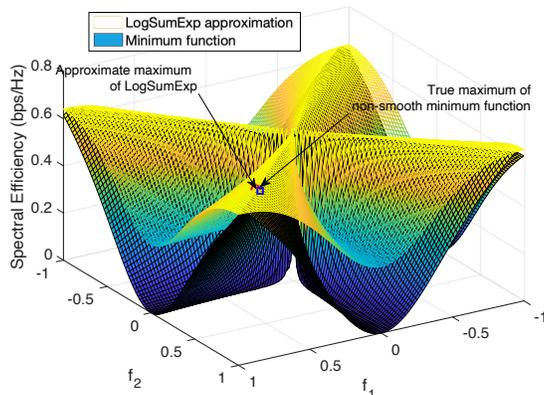
Fig. 1. An illustration of the comparison between the approximate maximum using the LogSumExp and the true maximum of the non-smooth minimum function.

## A. Reformulation to a Tractable Form

At first, we approximate the non-smooth minimum function by using the LogSumExp technique. With the LogSumExp, the minimum function is approximated as [31]

$$\min_{i=1,\ldots,N}\{x_i\} \approx -\alpha \log\left(\frac{1}{N}\sum_{i=1}^{N}\exp\left(\frac{x_i}{-\alpha}\right)\right), \tag{32}$$

where the approximation becomes tight as $\alpha \to +0$. Leveraging (32), we approximate

$$\min_{k\in\mathcal{K}}\{\bar{R}_{\mathsf{c}}^{\mathsf{ins.}}(k)\} \approx -\alpha \log\left(\frac{1}{K}\sum_{k=1}^{K}\exp\left(\frac{\bar{R}_{\mathsf{c}}^{\mathsf{ins.}}(k)}{-\alpha}\right)\right). \tag{33}$$

To help understand the LogSumExp approximation technique, we draw an illustration in Fig. 1. In Fig. 1, assuming the $N = 1$, $K = 2$, a landscape of the minimum spectral efficiency between two users is depicted. In addition to that, the approximate minimum spectral efficiency using the LogSumExp technique is also presented. As shown in the figure, the true maximum value of the non-smooth minimum function is tightly approximated by the LogSumExp technique.

Now we rewrite the precoding vectors $\mathbf{f}_{\mathsf{c}}, \mathbf{f}_1, \cdots, \mathbf{f}_K$ in a higher dimensional vector $\bar{\mathbf{f}}$ by stacking each vector as $\bar{\mathbf{f}} = [\mathbf{f}_{\mathsf{c}}^{\mathsf{T}}, \mathbf{f}_1^{\mathsf{T}}, \cdots, \mathbf{f}_K^{\mathsf{T}}]^{\mathsf{T}} \in \mathbb{C}^{N(K+1)\times 1}$. With this, we express the instantaneous spectral efficiencies regarding the common message $s_{\mathsf{c}}$ as

$$\bar{R}_{\mathsf{c}}^{\mathsf{ins.}}(k) = \log_2\left(\frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(k)\bar{\mathbf{f}}}\right), \tag{34}$$

where

$$\mathbf{A}_{\mathsf{c}}(k) = \mathrm{blkdiag}\left((\hat{\mathbf{h}}_k\hat{\mathbf{h}}_k^{\mathsf{H}} + \mathbf{\Phi}_k), \cdots, (\hat{\mathbf{h}}_k\hat{\mathbf{h}}_k^{\mathsf{H}} + \mathbf{\Phi}_k)\right) + \mathbf{I}_{N(K+1)}\frac{\sigma^2}{P}, \tag{35}$$

$$\mathbf{B}_{\mathsf{c}}(k) = \mathbf{A}_{\mathsf{c}}(k) - \mathrm{blkdiag}\left(\hat{\mathbf{h}}_k\hat{\mathbf{h}}_k^{\mathsf{H}}, \mathbf{0}, \cdots, \mathbf{0}\right). \tag{36}$$

Note that we implicitly assume $\|\bar{\mathbf{f}}\|^2 = 1$ to have (34). This assumption does not hurt the optimality since the spectral efficiency is monotonically increasing with the transmit power. It is also worthwhile to mention that (34) is presented as a function of Rayleigh quotient terms. Similar to this, we also write the spectral efficiency of the private message $s_k$ as

$$\bar{R}_k^{\mathsf{ins.}} = \log_2\left(\frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_k\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_k\bar{\mathbf{f}}}\right), \tag{37}$$

where

$$\mathbf{A}_k = \mathrm{blkdiag}\left(\mathbf{0}, (\hat{\mathbf{h}}_k\hat{\mathbf{h}}_k^{\mathsf{H}} + \mathbf{\Phi}_k), \cdots, (\hat{\mathbf{h}}_k\hat{\mathbf{h}}_k^{\mathsf{H}} + \mathbf{\Phi}_k)\right) + \mathbf{I}_{N(K+1)}\frac{\sigma^2}{P}, \tag{38}$$

$$\mathbf{B}_k = \mathbf{A}_k - \mathrm{blkdiag}\left(\mathbf{0}, \cdots, \mathbf{0}, \underbrace{\hat{\mathbf{h}}_k\hat{\mathbf{h}}_k^{\mathsf{H}}}_{\text{the } (k+1)\text{th block}}, \mathbf{0}, \cdots, \mathbf{0}\right). \tag{39}$$

With this Rayleigh quotients representation, the problem (21) is transformed to

$$\underset{\bar{\mathbf{f}}}{\mathrm{maximize}} \quad \log\left(\frac{1}{K}\sum_{k=1}^{K}\exp\left(\log_2\left(\frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(k)\bar{\mathbf{f}}}\right)^{-\frac{1}{\alpha}}\right)\right)^{-\alpha} + \sum_{k=1}^{K}\log_2\left(\frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_k\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_k\bar{\mathbf{f}}}\right) \tag{40}$$

$$\mathrm{subject\ to} \quad \|\bar{\mathbf{f}}\|^2 = 1. \tag{41}$$

We note that in (40), the obtained precoding vector $\bar{\mathbf{f}}$ can always be normalized by dividing the numerator and the denominator of each Rayleigh quotient with $\|\bar{\mathbf{f}}\|$, while not affecting the objective function. Thanks to this feature, the constraint $\|\mathbf{f}\|^2 = 1$ can vanish from (67). Now we are ready to tackle the problem (40).

### B. First-Order Optimality Condition

To approach the solution of the transformed problem (40), we derive a first-order optimality condition of (40). The following lemma shows the main result in this subsection.

**Lemma 1.** *The first-order optimality condition of the optimization problem* (40) *is satisfied if the following holds:*

$$\mathbf{B}_{\mathsf{KKT}}^{-1}(\bar{\mathbf{f}})\mathbf{A}_{\mathsf{KKT}}(\bar{\mathbf{f}})\bar{\mathbf{f}} = \lambda(\bar{\mathbf{f}})\bar{\mathbf{f}}, \tag{42}$$

*where*

$$\mathbf{A}_{\text{KKT}}(\bar{\mathbf{f}}) = \lambda_{\text{num}}(\bar{\mathbf{f}}) \times \sum_{k=1}^{K} \left[ \frac{\exp\left(\frac{1}{-\alpha} \log_2 \left(\frac{\bar{\mathbf{f}}^{\text{H}} \mathbf{A}_{\text{c}}(k) \bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\text{H}} \mathbf{B}_{\text{c}}(k) \bar{\mathbf{f}}}\right)\right)}{\sum_{\ell=1}^{K} \exp\left(\frac{1}{-\alpha} \log_2 \left(\frac{\bar{\mathbf{f}}^{\text{H}} \mathbf{A}_{\text{c}}(\ell) \bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\text{H}} \mathbf{B}_{\text{c}}(\ell) \bar{\mathbf{f}}}\right)\right)} \frac{\mathbf{A}_{\text{c}}(k)}{\bar{\mathbf{f}}^{\text{H}} \mathbf{A}_{\text{c}}(k) \bar{\mathbf{f}}} + \frac{\mathbf{A}_k}{\bar{\mathbf{f}}^{\text{H}} \mathbf{A}_k \bar{\mathbf{f}}} \right], \tag{43}$$

$$\mathbf{B}_{\text{KKT}}(\bar{\mathbf{f}}) = \lambda_{\text{den}}(\bar{\mathbf{f}}) \times \sum_{k=1}^{K} \left[ \frac{\exp\left(\frac{1}{-\alpha} \log_2 \left(\frac{\bar{\mathbf{f}}^{\text{H}} \mathbf{A}_{\text{c}}(k) \bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\text{H}} \mathbf{B}_{\text{c}}(k) \bar{\mathbf{f}}}\right)\right)}{\sum_{\ell=1}^{K} \exp\left(\frac{1}{-\alpha} \log_2 \left(\frac{\bar{\mathbf{f}}^{\text{H}} \mathbf{A}_{\text{c}}(\ell) \bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\text{H}} \mathbf{B}_{\text{c}}(\ell) \bar{\mathbf{f}}}\right)\right)} \frac{\mathbf{B}_{\text{c}}(k)}{\bar{\mathbf{f}}^{\text{H}} \mathbf{B}_{\text{c}}(k) \bar{\mathbf{f}}} + \frac{\mathbf{B}_k}{\bar{\mathbf{f}}^{\text{H}} \mathbf{B}_k \bar{\mathbf{f}}} \right], \tag{44}$$

*with*

$$\lambda(\bar{\mathbf{f}}) = \left\{ \frac{1}{K} \sum_{k=1}^{K} \exp\left( \log_2 \left(\frac{\bar{\mathbf{f}}^{\text{H}} \mathbf{A}_{\text{c}}(k) \bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\text{H}} \mathbf{B}_{\text{c}}(k) \bar{\mathbf{f}}}\right)^{-\frac{1}{\alpha}} \right) \right\}^{-\frac{\alpha}{\log_2 e}} \times \prod_{k=1}^{K} \left( \frac{\bar{\mathbf{f}}^{\text{H}} \mathbf{A}_k \bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\text{H}} \mathbf{B}_k \bar{\mathbf{f}}} \right) = \frac{\lambda_{\text{num}}(\bar{\mathbf{f}})}{\lambda_{\text{den}}(\bar{\mathbf{f}})}. \tag{45}$$

*Proof.* See Appendix A. □

Now we interpret the derived optimality condition (42). We first observe if a precoding vector $\bar{\mathbf{f}}$ satisfies the condition (42), then it also satisfies the first-order optimality condition, which means that the corresponding $\bar{\mathbf{f}}$ is a stationary point of the problem (40) whose gradient is zero. If the problem (40) has multiple stationary points, it is possible to exist multiple $\bar{\mathbf{f}}$ satisfying (42). Next, we see that (42) is presented as a form of the eigenvector problem for the matrix $\mathbf{B}_{\text{KKT}}^{-1}(\bar{\mathbf{f}}) \mathbf{A}_{\text{KKT}}(\bar{\mathbf{f}})$. More rigorously, (42) is cast as a NEPv [28]. As described in [28], NEPv is a generalized version of an eigenvalue problem, in that a matrix can be changed depending on an eigenvector in a nonlinear fashion. In our case, the matrix $\mathbf{B}_{\text{KKT}}^{-1}(\bar{\mathbf{f}}) \mathbf{A}_{\text{KKT}}(\bar{\mathbf{f}})$ is a nonlinear function of the eigenvector $\bar{\mathbf{f}}$. Crucially, in the formulated NEPv (42), the eigenvalue $\lambda(\bar{\mathbf{f}})$ is equivalent to the objective function of the problem (40). Accordingly, if we find the leading eigenvector of the NEPv (42), then it maximizes the objective function among multiple eigenvectors. Eventually, since (42) holds for any eigenvector, finding the leading eigenvector of the NEPv (42) is equivalent to finding the local optimal point that maximizes the objective function of (40) and has zero gradient. This leads to the following proposition.

**Proposition 1.** *Denoting that the local optimal point for the problem* (40) *as* $\bar{\mathbf{f}}^{\star}$, $\bar{\mathbf{f}}^{\star}$ *is the leading eigenvector of* $\mathbf{B}_{\text{KKT}}^{-1}(\bar{\mathbf{f}}^{\star}) \mathbf{A}_{\text{KKT}}(\bar{\mathbf{f}}^{\star})$ *satisfying*

$$\mathbf{B}_{\text{KKT}}^{-1}(\bar{\mathbf{f}}^{\star}) \mathbf{A}_{\text{KKT}}(\bar{\mathbf{f}}^{\star}) \bar{\mathbf{f}}^{\star} = \lambda^{\star} \bar{\mathbf{f}}^{\star}, \tag{46}$$

*where* $\lambda^{\star}$ *is the corresponding eigenvalue.*

---

**Algorithm 1** GPI-RS

---

    **initialize**: $\bar{\mathbf{f}}_{(0)}$ = MRT

    Set the iteration count $t = 1$.

    **while** $\left\| \bar{\mathbf{f}}_{(t)} - \bar{\mathbf{f}}_{(t-1)} \right\| > \epsilon$ **do**

        Construct the matrices $\mathbf{A}_{\mathsf{KKT}}(\bar{\mathbf{f}}_{(t-1)})$ and $\mathbf{B}_{\mathsf{KKT}}(\bar{\mathbf{f}}_{(t-1)})$ by using (43) and (44).

        Update $\bar{\mathbf{f}}_{(t)} \leftarrow \frac{\mathbf{B}_{\mathsf{KKT}}(\bar{\mathbf{f}}_{(t-1)})^{-1}\mathbf{A}_{\mathsf{KKT}}(\bar{\mathbf{f}}_{(t-1)})\bar{\mathbf{f}}_{(t-1)}}{\|\mathbf{B}_{\mathsf{KKT}}(\bar{\mathbf{f}}_{(t-1)})^{-1}\mathbf{A}_{\mathsf{KKT}}(\bar{\mathbf{f}}_{(t-1)})\bar{\mathbf{f}}_{(t-1)}\|}$.

        $t \leftarrow t + 1$.

    **end while**

---

Finding $\bar{\mathbf{f}}^{\star}$ is, however, not straightforward due to the intertwined nature of the problem. In the next subsection, we propose a novel method called GPI-RS. GPI-RS is able to obtain the leading eigenvector of the matrix $\mathbf{B}_{\mathsf{KKT}}(\bar{\mathbf{f}})^{-1}\mathbf{A}_{\mathsf{KKT}}(\bar{\mathbf{f}})$ in a computationally efficient fashion.

## C. Generalized Power Iteration for Rate-Splitting

The basic process of the proposed GPI-RS follows that of the conventio power iteration. Given $\bar{\mathbf{f}}_{(t-1)}$ obtained in the $(t-1)$th iteration, we construct the matrices $\mathbf{B}_{\mathsf{KKT}}(\bar{\mathbf{f}}_{(t-1)})$ and $\mathbf{A}_{\mathsf{KKT}}(\bar{\mathbf{f}}_{(t-1)})$ using (43) and (44). Then, we update the precoding vector for the current iteration as

$$\bar{\mathbf{f}}_{(t)} \leftarrow \frac{\mathbf{B}_{\mathsf{KKT}}^{-1}(\bar{\mathbf{f}}_{(t-1)})\mathbf{A}_{\mathsf{KKT}}(\bar{\mathbf{f}}_{(t-1)})\bar{\mathbf{f}}_{(t-1)}}{\|\mathbf{B}_{\mathsf{KKT}}^{-1}(\bar{\mathbf{f}}_{(t-1)})\mathbf{A}_{\mathsf{KKT}}(\bar{\mathbf{f}}_{(t-1)})\bar{\mathbf{f}}_{(t-1)}\|}. \tag{47}$$

We repeat this process until the convergence criterion is met. In this paper, we use $\left\| \bar{\mathbf{f}}_{(t)} - \bar{\mathbf{f}}_{(t-1)} \right\| < \epsilon$ for small enough $\epsilon$. We summarize this process in Algorithm 1. For an initial point $\bar{\mathbf{f}}_{(0)}$, we use maximum ratio transmission (MRT), which works well in the later simulations.

**Remark 2.** (Joint power control and beamforming) The proposed GPI-RS identifies the leading eigenvector $\bar{\mathbf{f}}^{\star}$ that maximizes (42). Since the vector $\bar{\mathbf{f}}$ is constructed by stacking all the precoding vectors corresponding to each message, the power allocation and the beamforming direction of each message are jointly identified within the found vector. For example, if $\|\bar{\mathbf{f}}^{\star}(1:N)\| = 0$, this means $\|\mathbf{f}_{\mathsf{c}}\| = 0$ in the obtained solution. Then the common message is not assigned any transmit power, therefore we do not use RSMA and go back to use classical SDMA. Thanks to this feature, the proposed GPI-RS automatically determines the message setups depending on channel conditions; so that there is no need to employ a separate process to determine whether to use a common message.

**Remark 3.** (Algorithm complexity) The total computational complexity of the proposed GPI-RS is dominated by the calculation of $\mathbf{B}_{\mathsf{KKT}}^{-1}(\bar{\mathbf{f}})$. The matrix $\mathbf{B}_{\mathsf{KKT}}^{-1}(\bar{\mathbf{f}})$ is the sum of the block-diagonal matrices as presented in (44). Specifically, $K+1$ number of $N \times N$ submatrices are concatenated, so that the total size is $(K+1)N \times (K+1)N$. For this reason, the inverse matrix $\mathbf{B}_{\mathsf{KKT}}^{-1}(\bar{\mathbf{f}})$ is obtained by computing the inverse of each submatrix, and this requires the complexity with the order of $O(\frac{1}{3}(K+1)N^3)$. This results in that the complexity of the proposed GPI-RS per iteration is with the order of $O(\frac{1}{3}KN^3)$ when $N$ and $K$ increase with the same order. We note that this is substantially small compared to the existing methods. For example, the conventional WMMSE methods based on QCQP [7], [32] need the complexity order of $O((KN)^{3.5})$ [8], [30]. Further, the CCCP based method [10] is associated with the complexity order of $O(N^6 K^{0.5} 2^{3.5K})$. In particular, it is noteworthy that the proposed GPI-RS has the linear order complexity with the number of user $K$; which makes the proposed GPI-RS advantageous when there are a large number of users. Additionally, our algorithm is easy to implement in practice in that CVX is not needed to use to obtain a solution. We note that the computational complexity of the proposed method can be further reduced by adopting matrix inversion approximation techniques such as Chebyshev iteration [33] or Neumann series [34], or by limiting the maximum number of iterations in the proposed method.

**Remark 4.** (Selection of the parameter $\alpha$) Even though small $\alpha$ is desirable since it provides accurate approximation in the LogSumExp technique, using too small $\alpha$ may cause the algorithm to diverge. Analytically identifying the optimal $\alpha$, however, is very challenging. To find the proper $\alpha$ numerically, we can modify the GPI-RS to obtain the smallest $\alpha$ that makes the GPI-RS algorithm converge. Specifically, we start the GPI-RS with a small $\alpha$. If the iteration loop of the GPI-RS does not converge within the predetermined number of iterations, then we enforce to terminate the loop, increase $\alpha$, and newly start the algorithm again. We repeat this process until the algorithm converges before the predetermined number. We can empirically adapt the starting $\alpha$ value and the increasing ratio depending on the system configuration to reduce the algorithm time.

**Remark 5.** (Principle of GPI-RS) We explain the principle of the GPI-RS algorithm through the conventional power iteration. In the conventional power iteration, we obtain the leading eigenvector of a matrix $\mathbf{M} \in \mathbb{C}^{n \times n}$ by iteratively calculating $\mathbf{q}_{(t+1)} = \frac{\mathbf{M}^t \mathbf{q}_{(0)}}{\|\mathbf{M}^t \mathbf{q}_{(0)}\|}$. A rationale that

the power iteration converges to the leading eigenvector of $\mathbf{M}$ is as follows. Since a set of eigenvectors form a set of basis, we can represent $\mathbf{q}_{(0)} = \sum_{i=1}^{n} \alpha_i \mathbf{x}_i$, where $\mathbf{x}_i$ indicates the $i$-th eigenvector and $\alpha_i$ is the corresponding weight. Denoting that $\lambda_i$ is the $i$-th eigenvalue that $|\lambda_1| > |\lambda_2| \geq \cdots \geq |\lambda_n|$, we use $\mathbf{M} = \sum_{i=1}^{n} \lambda_i \mathbf{x}_i$ to derive

$$\mathbf{M}\mathbf{q}_{(t)} = \sum_{i=1}^{n} \alpha_i \lambda_i^t \mathbf{x}_i = \alpha_1 \lambda_1^t \left( \mathbf{x}_1 + \underbrace{\sum_{i=2}^{n} \frac{\alpha_i}{\alpha_1} \left( \frac{\lambda_i}{\lambda_1} \right)^t \mathbf{x}_i}_{(a)} \right). \tag{48}$$

As $t \to \infty$, (a) vanishes; thereby the remaining term converges to the leading eigenvector $\mathbf{x}_1$.

As presented in Proposition 1, our problem generalizes a conventional eigenvalue problem by considering an eigenvector-dependent matrix $\mathbf{M}(\mathbf{x})$, known as NEPv [28]. That is to say, we aim to identify $\mathbf{x}_1$ that fulfills $\mathbf{M}(\mathbf{x}_1)\mathbf{x}_1 = \lambda_1 \mathbf{x}_1$, where $\lambda_1$ is the maximum eigenvalue and $\mathbf{M}(\mathbf{x}) = \mathbf{B}_{\mathsf{KKT}}(\mathbf{x})^{-1} \mathbf{A}_{\mathsf{KKT}}(\mathbf{x})$ in our case. For convenience, denote $\mathbf{M}(\mathbf{x}_i)\mathbf{x}_i = g(\mathbf{x}_i)$. Then with an arbitrary vector $\mathbf{x}$, the Taylor expansion at $\mathbf{x}_1$ for $g(\mathbf{x})$ leads to

$$g(\mathbf{x})^{\mathsf{H}} \mathbf{x}_i = g(\mathbf{x}_1)^{\mathsf{H}} \mathbf{x}_i + (\mathbf{x} - \mathbf{x}_1)^{\mathsf{H}} \nabla g(\mathbf{x}_1) \mathbf{x}_i + o(\|\mathbf{x} - \mathbf{x}_1\|). \tag{49}$$

On one hand, we have

$$\left( g(\mathbf{x})^{\mathsf{H}} \mathbf{x}_1 \right)^2 = (\lambda_1 + o(\|\mathbf{x} - \mathbf{x}_1\|))^2 \tag{50}$$

due to the fact that $\mathbf{M}(\mathbf{x}_1)\mathbf{x}_1 = \lambda_1 \mathbf{x}_1$. On the other hand, assuming that $\{\mathbf{v}_1, \cdots, \mathbf{v}_n\}$ is a set of orthonormal basis where $\mathbf{v}_1 = \mathbf{x}_1$, we also have

$$\sum_{i=2}^{n} \left( g(\mathbf{x})^{\mathsf{H}} \mathbf{v}_i \right)^2 \leq \sum_{i=2}^{n} \left[ \lambda_i^2 (\mathbf{x}^{\mathsf{H}} \mathbf{v}_i)^2 + 2\lambda_i (\mathbf{x}^{\mathsf{H}} \mathbf{v}_i) o(\|\mathbf{x} - \mathbf{x}_1\|) + o(\|\mathbf{x} - \mathbf{x}_1\|)^2 \right] \tag{51}$$

$$\leq (\lambda_2 \|\mathbf{x} - \mathbf{x}_1\| + o(\|\mathbf{x} - \mathbf{x}_1\|))^2 \tag{52}$$

Under a premise that $|\lambda_1| > |\lambda_2| \geq |\lambda_i|, \forall i \neq 1, 2$, by iteratively projecting $\mathbf{x}$ onto $\mathbf{M}(\mathbf{x})$ with the GPI-RS algorithm, each component corresponding to non-leading eigenvectors $\mathbf{x}_2, \cdots, \mathbf{x}_n$ vanishes. Accordingly, the GPI-RS converges to the leading eigenvector $\mathbf{x}_1$.

## V. GENERALIZATION TO MULTIPLE-LAYER RATE SPLITTING

This section discusses how to generalize the proposed framework to multiple-layer RSMA. As mentioned in [9], if groups of users are located within multiple clusters, it is beneficial to

exploit a partial common message that includes the messages of a subset of the users. Employing multiple-layer RSMA, the transmit signal $\mathbf{x}$ is given by

$$\mathbf{x} = \mathbf{f}_{\mathsf{c}} s_{\mathsf{c}} + \sum_{i=1}^{G} \mathbf{f}_{\mathsf{c},\mathcal{K}_i} s_{\mathsf{c},\mathcal{K}_i} + \sum_{k=1}^{K} \mathbf{f}_k s_k, \tag{53}$$

where the partial common message $s_{\mathsf{c},\mathcal{K}_i}$ is decoded for the users included in $\mathcal{K}_i \subset \mathcal{K}$. We denote $|\mathcal{K}_i| = K_i$. If $G = 0$, i.e., there is no partial common message, then (53) reduces to the 1-layer RSMA (5). Assuming $\mathcal{K}_i \cap \mathcal{K}_j = \emptyset$, user $k$ for $k \in \mathcal{K}_i$ decodes the messages with the following order: $s_{\mathsf{c}} \rightarrow s_{\mathsf{c},\mathcal{K}_i} \rightarrow s_k$. To guarantee that the partial common message $s_{\mathsf{c},\mathcal{K}_i}$ is successfully decoded for the users in $\mathcal{K}_i$, the information rate of $s_{\mathsf{c},\mathcal{K}_i}$ is determined as

$$R_{\mathsf{c},\mathcal{K}_i} = \min_{k \in \mathcal{K}_i} \left\{ \mathbb{E}_{\{\hat{\mathbf{h}}_k\}} \left[ R_{\mathsf{c},\mathcal{K}_i}^{\mathsf{ins.}}(k) \right] \right\}. \tag{54}$$

By using the same technique presented in Section II-D, we derive a lower bound as

$$R_{\mathsf{c},\mathcal{K}_i} \geq \mathbb{E}_{\{\hat{\mathbf{h}}_{k \in \mathcal{K}_i}\}} \left[ \min_{k \in \mathcal{K}_i} \left\{ \bar{R}_{\mathsf{c},\mathcal{K}_i}^{\mathsf{ins.}}(k) \right\} \right], \tag{55}$$

where

$$\bar{R}_{\mathsf{c},\mathcal{K}_i}^{\mathsf{ins.}}(k) = \log_2 \left( 1 + \frac{|\hat{\mathbf{h}}_k^{\mathsf{H}} \mathbf{f}_{\mathsf{c},\mathcal{K}_i}|^2}{\sum_{j=1,j\neq i}^{G} |\hat{\mathbf{h}}_k^{\mathsf{H}} \mathbf{f}_{\mathsf{c},\mathcal{K}_j}|^2 + \sum_{\ell=1}^{K} |\hat{\mathbf{h}}_k^{\mathsf{H}} \mathbf{f}_\ell|^2 + \sum_{j=1}^{G} \mathbf{f}_{\mathsf{c},\mathcal{K}_j}^{\mathsf{H}} \boldsymbol{\Phi}_k \mathbf{f}_{\mathsf{c},\mathcal{K}_j} + \sum_{\ell=1}^{K} \mathbf{f}_\ell^{\mathsf{H}} \boldsymbol{\Phi}_k \mathbf{f}_\ell + \frac{\sigma^2}{P}} \right). \tag{56}$$

We observe that the SINR in (56) does not have interference from the common message $s_{\mathsf{c}}$ since we assume that the common message is already decoded and eliminated via SIC. Similar to this, we also characterize a lower bound on the instantaneous spectral efficiency for the common message and the private message as follows:

$$\bar{R}_{\mathsf{c}}^{\mathsf{ins.}}(k)$$

$$= \log_2 \left( 1 + \frac{|\hat{\mathbf{h}}_k^{\mathsf{H}} \mathbf{f}_{\mathsf{c}}|^2}{\sum_{j=1}^{G} |\hat{\mathbf{h}}_k^{\mathsf{H}} \mathbf{f}_{\mathsf{c},\mathcal{K}_j}|^2 + \sum_{\ell=1}^{K} |\hat{\mathbf{h}}_k^{\mathsf{H}} \mathbf{f}_\ell|^2 + \mathbf{f}_{\mathsf{c}}^{\mathsf{H}} \boldsymbol{\Phi}_k \mathbf{f}_{\mathsf{c}} + \sum_{j=1}^{G} \mathbf{f}_{\mathsf{c},\mathcal{K}_j}^{\mathsf{H}} \boldsymbol{\Phi}_k \mathbf{f}_{\mathsf{c},\mathcal{K}_j} + \sum_{\ell=1}^{K} \mathbf{f}_\ell^{\mathsf{H}} \boldsymbol{\Phi}_k \mathbf{f}_\ell + \frac{\sigma^2}{P}} \right), \tag{57}$$

$$\bar{R}_k = \log_2 \left( 1 + \frac{|\hat{\mathbf{h}}_k^{\mathsf{H}} \mathbf{f}_k|^2}{\sum_{j=1,j\neq i}^{G} |\hat{\mathbf{h}}_k^{\mathsf{H}} \mathbf{f}_{\mathsf{c},\mathcal{K}_j}|^2 + \sum_{\ell=1,\ell\neq k}^{K} |\hat{\mathbf{h}}_k^{\mathsf{H}} \mathbf{f}_\ell|^2 + \sum_{j=1,j\neq i}^{G} \mathbf{f}_{\mathsf{c},\mathcal{K}_j}^{\mathsf{H}} \boldsymbol{\Phi}_k \mathbf{f}_{\mathsf{c},\mathcal{K}_j} + \sum_{\ell=1}^{K} \mathbf{f}_\ell^{\mathsf{H}} \boldsymbol{\Phi}_k \mathbf{f}_\ell + \frac{\sigma^2}{P}} \right). \tag{58}$$

Accordingly, we formulate the sum spectral efficiency maximization problem for multiple-layer RSMA as

$$\underset{\mathbf{f}_c,\mathbf{f}_{c,\mathcal{K}_1},\cdots,\mathbf{f}_{c,\mathcal{K}_G},\mathbf{f}_1,\cdots,\mathbf{f}_K}{\text{maximize}} \quad \min_{k\in\mathcal{K}}\{\bar{R}_c^{\text{ins.}}(k)\} + \sum_{i=1}^{G}\min_{k\in\mathcal{K}_i}\{\bar{R}_{c,\mathcal{K}_i}^{\text{ins.}}(k)\} + \sum_{k=1}^{K}\bar{R}_k^{\text{ins.}} \tag{59}$$

$$\text{subject to} \quad \|\mathbf{f}_c\|^2 + \sum_{k=1}^{K}\|\mathbf{f}_k\|^2 \leq 1. \tag{60}$$

To obtain a solution of (59), we reformulate (59) by following our approach used in the 1-layer RSMA. We first approximate $\min_{k\in\mathcal{K}_i}\left\{\bar{R}_{c,\mathcal{K}_i}^{\text{ins.}}(k)\right\}$ by using the LogSumExp technique as

$$\min_{k\in\mathcal{K}_i}\{\bar{R}_{c,\mathcal{K}_i}^{\text{ins.}}(k)\} \approx -\alpha \log\left(\frac{1}{K_i}\sum_{k\in\mathcal{K}_i}\exp\left(\frac{\bar{R}_{c,\mathcal{K}_i}^{\text{ins.}}(k)}{-\alpha}\right)\right). \tag{61}$$

By following our approach, we define a high dimensional precoding vector $\bar{\mathbf{f}}$ as

$$\bar{\mathbf{f}} = [\mathbf{f}_c^{\mathsf{T}}, \mathbf{f}_{c,\mathcal{K}_1}^{\mathsf{T}}, \cdots, \mathbf{f}_{c,\mathcal{K}_G}^{\mathsf{T}}, \mathbf{f}_1^{\mathsf{T}}, \cdots, \mathbf{f}_K^{\mathsf{T}}]^{\mathsf{T}} \in \mathbb{C}^{N(K+G+1)\times 1}. \tag{62}$$

The higher dimensional precoding vector $\bar{\mathbf{f}}$ in (62) allows to represent the spectral efficiency expression of the partial common message $s_{c,\mathcal{K}_i}$ as a Rayleigh quotient form.

$$\bar{R}_{c,\mathcal{K}_i}^{\text{ins.}}(k) = \log_2\left(\frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{c,\mathcal{K}_i}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{c,\mathcal{K}_i}(k)\bar{\mathbf{f}}}\right), \tag{63}$$

where

$$\mathbf{A}_{c,\mathcal{K}_i}(k) = \text{blkdiag}\left(\mathbf{0}, (\hat{\mathbf{h}}_k\hat{\mathbf{h}}_k^{\mathsf{H}} + \mathbf{\Phi}_k), \cdots, (\hat{\mathbf{h}}_k\hat{\mathbf{h}}_k^{\mathsf{H}} + \mathbf{\Phi}_k)\right) + \mathbf{I}_{N(K+G+1)}\frac{\sigma^2}{P}, \tag{64}$$

$$\mathbf{B}_{c,\mathcal{K}_i}(k) = \mathbf{A}_{c,\mathcal{K}_i}(k) - \text{blkdiag}\left(\mathbf{0}, \cdots, \underbrace{\hat{\mathbf{h}}_k\hat{\mathbf{h}}_k^{\mathsf{H}}}_{(1+i)\text{th block}}, \cdots, \mathbf{0}\right). \tag{65}$$

We also represent $\bar{R}_c^{\text{ins.}}(k)$ and $\bar{R}_k^{\text{ins.}}$ adequately by considering the partial common message. With this representation, the problem (59) is converted to

$$\underset{\bar{\mathbf{f}}}{\text{maximize}} \quad \log\left(\frac{1}{K}\sum_{k=1}^{K}\exp\left(\log_2\left(\frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_c(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_c(k)\bar{\mathbf{f}}}\right)^{-\frac{1}{\alpha}}\right)\right)^{-\alpha}$$

$$+ \sum_{i=1}^{G}\log\left(\frac{1}{K_i}\sum_{k\in\mathcal{K}_i}\exp\left(\log_2\left(\frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{c,\mathcal{K}_i}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{c,\mathcal{K}_i}(k)\bar{\mathbf{f}}}\right)^{-\frac{1}{\alpha}}\right)\right)^{-\alpha} + \sum_{k=1}^{K}\log_2\left(\frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_k\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_k\bar{\mathbf{f}}}\right) \tag{66}$$

$$\text{subject to} \quad \left\|\bar{\mathbf{f}}\right\|^2 = 1. \tag{67}$$

To apply the GPI-RS algorithm for (66), we derive the optimality condition for (66) in the following corollary.

**Corollary 1.** *The first-order optimality condition of the optimization problem* (66) *is satisfied if the following holds:*

$$\mathbf{B}_{\mathsf{KKT}}^{-1}(\bar{\mathbf{f}})\mathbf{A}_{\mathsf{KKT}}(\bar{\mathbf{f}})\bar{\mathbf{f}} = \lambda(\bar{\mathbf{f}})\bar{\mathbf{f}}, \tag{68}$$

*where*

$$
\begin{aligned}
\mathbf{A}_{\mathsf{KKT}}(\bar{\mathbf{f}}) = \lambda_{\mathsf{num}}(\bar{\mathbf{f}}) \times & \left[ \sum_{k=1}^{K} \left\{ \frac{\exp\left( \frac{1}{-\alpha} \log_2 \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(k)\bar{\mathbf{f}}} \right) \right)}{\sum_{\ell=1}^{K} \exp\left( \frac{1}{-\alpha} \log_2 \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(\ell)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(\ell)\bar{\mathbf{f}}} \right) \right)} \frac{\mathbf{A}_{\mathsf{c}}(k)}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(k)\bar{\mathbf{f}}} + \frac{\mathbf{A}_k}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_k\bar{\mathbf{f}}} \right\} \right. \\
& \left. + \sum_{i=1}^{G} \left\{ \sum_{k\in\mathcal{K}_i} \left( \frac{\exp\left( \frac{1}{-\alpha} \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c},\mathcal{K}_i}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c},\mathcal{K}_i}(k)\bar{\mathbf{f}}} \right)}{\sum_{j\in\mathcal{K}_i} \exp\left( \frac{1}{-\alpha} \log_2 \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c},\mathcal{K}_i}(j)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c},\mathcal{K}_i}(j)\bar{\mathbf{f}}} \right) \right)} \frac{\mathbf{A}_{\mathsf{c},\mathcal{K}_i}(k)}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c},\mathcal{K}_i}(k)\bar{\mathbf{f}}} \right) \right\} \right],
\end{aligned} \tag{69}
$$

$$
\begin{aligned}
\mathbf{B}_{\mathsf{KKT}}(\bar{\mathbf{f}}) = \lambda_{\mathsf{den}}(\bar{\mathbf{f}}) \times & \left[ \sum_{k=1}^{K} \left\{ \frac{\exp\left( \frac{1}{-\alpha} \log_2 \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(k)\bar{\mathbf{f}}} \right) \right)}{\sum_{\ell=1}^{K} \exp\left( \frac{1}{-\alpha} \log_2 \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(\ell)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(\ell)\bar{\mathbf{f}}} \right) \right)} \frac{\mathbf{B}_{\mathsf{c}}(k)}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(k)\bar{\mathbf{f}}} + \frac{\mathbf{B}_k}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_k\bar{\mathbf{f}}} \right\} \right. \\
& \left. + \sum_{i=1}^{G} \left\{ \sum_{k\in\mathcal{K}_i} \left( \frac{\exp\left( \frac{1}{-\alpha} \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c},\mathcal{K}_i}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c},\mathcal{K}_i}(k)\bar{\mathbf{f}}} \right)}{\sum_{j\in\mathcal{K}_i} \exp\left( \frac{1}{-\alpha} \log_2 \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c},\mathcal{K}_i}(j)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c},\mathcal{K}_i}(j)\bar{\mathbf{f}}} \right) \right)} \frac{\mathbf{B}_{\mathsf{c},\mathcal{K}_i}(k)}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c},\mathcal{K}_i}(k)\bar{\mathbf{f}}} \right) \right\} \right],
\end{aligned} \tag{70}
$$

*with*

$$
\begin{aligned}
\lambda(\bar{\mathbf{f}}) = & \left\{ \frac{1}{K} \sum_{k=1}^{K} \exp\left( \log_2 \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(k)\bar{\mathbf{f}}} \right)^{-\frac{1}{\alpha}} \right) \right\}^{-\frac{\alpha}{\log_2 e}} \times \prod_{k=1}^{K} \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_k\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_k\bar{\mathbf{f}}} \right) \\
& \times \prod_{i=1}^{G} \left\{ \frac{1}{K_i} \sum_{k\in\mathcal{K}_i} \exp\left( \log_2 \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c},\mathcal{K}_i}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c},\mathcal{K}_i}(k)\bar{\mathbf{f}}} \right)^{-\frac{1}{\alpha}} \right) \right\}^{-\frac{\alpha}{\log_2 e}} = \frac{\lambda_{\mathsf{num}}(\bar{\mathbf{f}})}{\lambda_{\mathsf{den}}(\bar{\mathbf{f}})}.
\end{aligned} \tag{71}
$$

*Proof.* It can be proven by extending Lemma 1. □

With Corollary 1, we apply the GPI-RS described in Algorithm 1 by using (69) and (70) instead of (43) and (44).

## VI. Numerical Results

In this section, we evaluate the sum spectral efficiency performance to demonstrate the proposed GPI-RS. For the baseline methods, we consider the followings:

- **MRT**: The precoding vectors is designed by matching the estimated channel vector. Specifically, we have $\mathbf{f}_k = \hat{\mathbf{h}}_k$, $k \in \mathcal{K}$, and $\mathbf{f}_c = \mathbf{0}$.

- **RZF** : The precoding vectors are designed by following the ZF rule, while regularizing it depending on SNR:

$$\mathbf{f}_k = \left( \hat{\mathbf{H}}\hat{\mathbf{H}}^{\mathsf{H}} + \mathbf{I}\frac{\sigma^2}{P} \right)^{-1} \hat{\mathbf{h}}_k^{\mathsf{H}}, \ k \in \mathcal{K}, \ \text{and} \ \mathbf{f}_c = \mathbf{0}. \tag{72}$$

  As SNR goes to infinity, RZF becomes equal to ZF.

- **Sum SE Max with no RS**: In this method, we use the method proposed in [4] to maximize the sum spectral efficiency using classical SDMA without considering RSMA. Note that we do not incorporate the CSIT estimation error into this method, i.e., we treat the estimated channel vector as the true channel in this method.

- **WMMSE-SAA**: This case indicates the WMMSE method with the SAA technique [7]. We use 1000 samples for the SAA technique. A detailed description of this approach is presented in Section III.

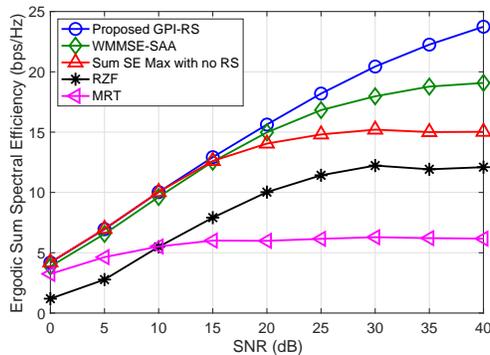In what follows, we present the simulation results.

**Ergodic Sum Spectral Efficiency per SNR**: First, we compare the ergodic sum spectral efficiency of the proposed GPI-RS and the other baseline methods. The basic simulation setups are explained in the caption of Fig. 2. For updating $\alpha$, we set the initial $\alpha$ value as 0.1 if $\mathsf{SNR} <$ 15dB, and 0.5 for the rest of the cases. We note that this initial setup is designed empirically. If the GPI-RS loop is not terminated within 50 iterations, we increase $\alpha$ by 0.5 and repeat the algorithm. As shown in Fig. 2, the proposed GPI-RS provides meaningful spectral efficiency gains over the baseline methods in both of the $6 \times 4$ and the $12 \times 8$ cases. In particular, compared to the WMMSE-SAA method at $\mathsf{SNR} = 40$dB, the GPI-RS obtains around 24% gains in the $6 \times 4$ case and 19% gains in the $12 \times 8$ case. We observe that considerable gains are achieved in both cases by using the proposed method. The rationales of the performance gains are two folds. First, the GPI-RS can reach the best local optimal point by NEPv principle, while the WMMSE-SAA approach cannot guarantee the best local optimum. Second, we incorporate the CSIT estimation error into our performance characterization in a rigorous way, while the WMMSE-SAA relies on randomly generated samples. These two features bring the gains of the proposed GPI-RS method.

Especially, the second point sheds light on a reason why the proposed GPI-RS obtains more significant performance gains in the high SNR regime. In the high SNR regime, the performance is mainly determined by the interference induced from the inaccurate CSIT as the noise becomes negligible. For this reason, to achieve high spectral efficiency performance, rigorous treatment of CSIT error is required. The SAA, used in the WMMSE-SAA, relies on randomly generated samples instead of incorporating the CSIT estimation error effects into the spectral efficiency performance. Compared to this, the GPI-RS specifically incorporates the CSIT estimation error effects into its design by deriving a rigorous lower bound. This difference causes more significant performance gaps in the high SNR regime. Additionally, by comparing the sum SE max with no RS case, we see that the spectral efficiency gains of the proposed GPI-RS increase as SNR increases. This indicates that the common message rate portion increases in the high SNR regime, which matches our intuition on RSMA.
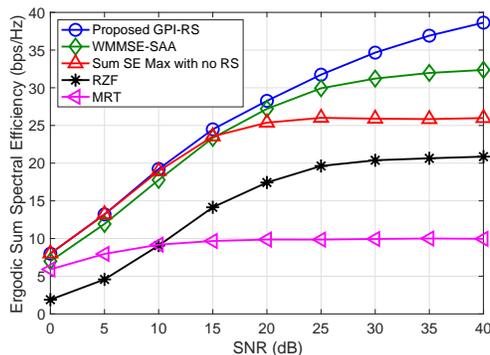
In Fig. 3, we also illustrate the ergodic sum spectral efficiency by assuming the the $32 \times 16$ case. As in the above case, we also observe that the GPI-RS provides significant gains over the baseline methods in Fig. 3. Specifically, at $\mathsf{SNR} = 40$dB, the GPI-RS offers about 30% gains compared to the sum spectral efficiency maximization method without RSMA. This demonstrates that the proposed GPI-RS works well in the regimes of large BS antennas and users.

**Ergodic Sum Spectral Efficiency per CSIT Accuracy**: Now, we investigate the sum spectral efficiency depending on the CSIT accuracy. We depict the sum spectral efficiency by increasing $\tau_{\mathsf{ul}} p_{\mathsf{ul}}$ in Fig. 4. Note that the CSIT estimation accuracy increases as $\tau_{\mathsf{ul}} p_{\mathsf{ul}}$ increases. In Fig. 4, we observe that the relative performance gains of the GPI-RS over the WMMSE-SAA increase as the CSIT becomes accurate. For instance, if $\tau_{\mathsf{ul}} p_{\mathsf{ul}} = 0.1$, the WMMSE-SAA outperforms the GPI-RS by 5%, while if $\tau_{\mathsf{ul}} p_{\mathsf{ul}} = 8$, the GPI-RS outperforms the WMMSE-SAA by 27%. This is because, as the CSIT is more accurate, the regularization term of our lower bound also becomes accurate; resulting in that our lower bound becomes tight. Then the CSIT estimation error is suitably reflected into the GPI-RS design.

**Convergence**: We depict the convergence behavior. The $\alpha$ update process is equal as above: the initial $\alpha$ value is 0.1 if $\mathsf{SNR} < 15$dB and 0.5 for the rest of the cases. We update $\alpha$ by 0.5 per every 50 iterations. We recall that this is for finding the smallest $\alpha$ that guarantees convergence. In Fig. 5, we observe that the GPI-RS converges well with small $\alpha$ in low SNR. On contrary to this, in the high SNR, we need to tune $\alpha$ to make the GPI-RS converge. For instance, when $\mathsf{SNR} = 20$dB, $\alpha$ needs to be updated one time until convergence, while when $\mathsf{SNR} = 40$dB,

Fig. 2. The sum spectral efficiency comparison per SNR. The simulation setup is: In (a), $N = 6$, $K = 4$, $\tau_{ul}p_{ul} = 4$ and $\sigma^2 = 1$. In (b) $N = 12$, $K = 8$, $\tau_{ul}p_{ul} = 4$ and $\sigma^2 = 1$. In both cases, each user's location is randomly determined, so that the AoA $\theta_k$ is drawn from the uniform distribution. The angular spread is fixed as $\Delta_k = \pi/6$ for $k \in \mathcal{K}$.
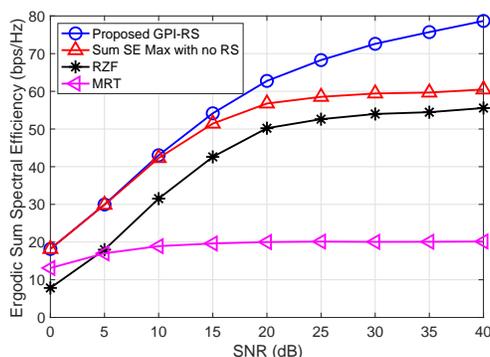


Fig. 3. The sum spectral efficiency comparison per SNR assuming $N = 32$, $K = 16$. The other setups are same with Fig. 2.

$\alpha$ needs to be updated two times until convergence. Through this observation, we numerically confirm that using the presented $\alpha$ update method, the convergence of the proposed GPI-RS is
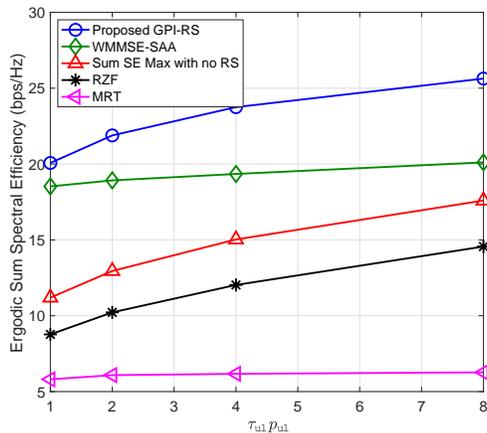
Fig. 4. The sum spectral efficiency comparison per CSIT accuracy. The simulation setup is: $N = 6$, $K = 4$, and $\sigma^2 = 1$. The AoA $\theta_k$ is drawn from the uniform distribution. The angular spread is $\Delta_k = \pi/6$ for $k \in \mathcal{K}$.
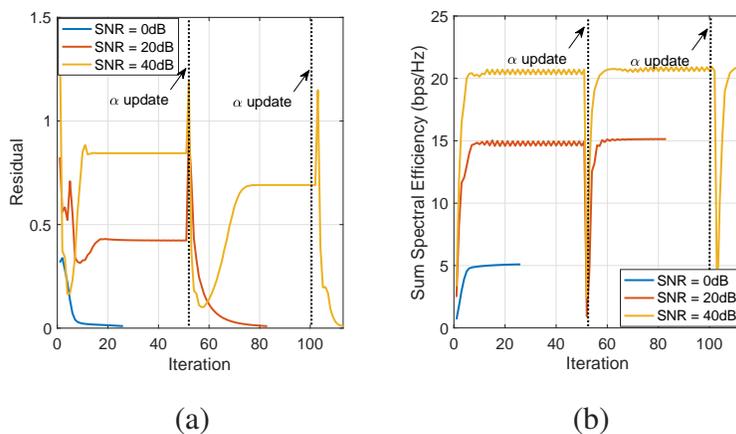


(a)          (b)

Fig. 5. The convergence behavior of the proposed GPI-RS per iteration. In (a) the residual is depicted. In (b). the sum spectral efficiency is drawn. The simulation setup is: $N = 6$, $K = 4$, $\tau_{\mathsf{ul}} p_{\mathsf{ul}} = 4$ and $\sigma^2 = 1$. The AoA $\theta_k$ is drawn from the uniform distribution and the angular spread is $\Delta_k = \pi/6$ for $k \in \mathcal{K}$. Residual is defined as $\|\bar{\mathbf{f}}_{(t)} - \bar{\mathbf{f}}_{(t-1)}\|$.

TABLE I

AVERAGE MATLAB CPU TIME (SEC)

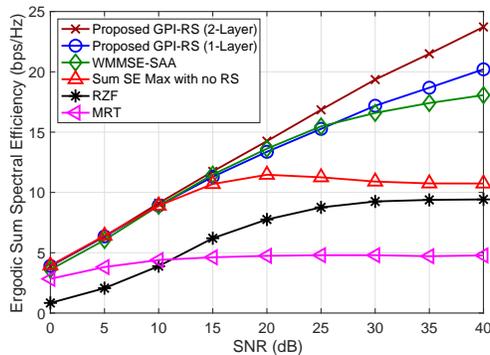| Setup | **Proposed GPI-RS** | **WMMSE-SAA** | Comparison (%) |
|-------|---------------------|---------------|----------------|
| $6 \times 4$ | 2.62 | 35.12 | 7.4% |
| $12 \times 8$ | 11.24 | 170.29 | 6.6% |

guaranteed well.

Fig. 6. The sum spectral efficiency comparison per CSIT accuracy. The simulation setup is: $N = 6$, $K = 4$, $\tau_{ul} p_{ul} = 4$, and $\sigma^2 = 1$. We assume that user 1, 2 and user 3, 4 are clustered in the same location: $\theta_k = \pi/3$ for $k \in \{1, 2\}$ and $\theta_k = 2\pi/3$ for $k \in \{3, 4\}$. The angular spread is fixed as $\Delta_k = \pi/6$ for $k \in \mathcal{K}$.

**Computation Time**: As a complement to the complexity analysis in Remark 3, we compare numerical MATLAB computation times between the proposed GPI-RS and the WMMSE-SAA in Table I. The simulation setups are equivalent to those used to produce Fig. 2. As shown in Table I, the proposed GPI-RS consumes only 7.4% of the computation time compared to the WMMSE-SAA in the $6 \times 4$ case and 6.6% in the $12 \times 8$ case on average. This dramatic complexity reduction comes from two sources: First, we approximate the non-smooth minimum function via a LogSumExp technique, by which we avoid having $K$ distinct constraints on the common message rate in the optimization problem. Second, by using the GPI-RS, we do not rely on an off-the-shelf optimization toolbox, including CVX. This result indicates that the proposed method is beneficial in terms of computational complexity, not only in an analytical (big-$\mathcal{O}$) sense but also in a practical (numerical computation time) sense. We note that Table I is only a rough measure of relative computational complexity.

**Multiple-layer RSMA** : Next, we evaluate the sum spectral efficiency performance of the multiple-layer RSMA using the generalized GPI-RS. In particular, we assume the 2-layer RSMA scenario in the $6 \times 4$ case, wherein one common message, two partial common messages, and private messages are transmitted. Specifically, the partial common messages are $s_{c,\mathcal{K}_1}$ and $s_{c,\mathcal{K}_2}$, and $\mathcal{K}_1 = \{1, 2\}$ and $\mathcal{K}_2 = \{3, 4\}$; therefore the partial common message $s_{c,\mathcal{K}_1}$ is intended to user 1 and 2 and the partial common message $s_{c,\mathcal{K}_2}$ is intended to user 3 and 4. We construct a favorable channel environment for this message setup, in which users 1, 2, and user 3, 4 are clustered in the same location, respectively. The performance result is illustrated in Fig. 6,

whose caption includes detailed simulation setups. As shown in Fig. 6, the GPI-RS for 2-layer RSMA achieves the best sum spectral efficiency performance. Specifically, the GPI-RS for 2-layer RSMA has around 17% gains over the GPI-RS for 1-layer RSMA and round 31% gains over the WMMSE-SAA. This confirms the observation of [20] more in a more rigorous way. We also show that the proposed framework is well extended to multiple-layer RSMA.

## VII. CONCLUSION

In this paper, we have proposed a novel precoding optimization method for downlink MIMO with RSMA. Aiming to maximize the sum spectral efficiency of the considered system, we have formulated an optimization problem, while a sum spectral efficiency maximization problem regarding linear precoding vectors is infeasible to solve due to its non-convexity and non-smoothness. To resolve this, we have approximated a non-smooth minimum function using the LogSumExp technique and reformulated the problem into a tractable form. We have shown that the first-order optimality condition of the reformulated problem is cast as a NEPv. In order to find the leading eigenvector for the derived condition, we have proposed the GPI-RS. We also have extended the GPI-RS to the multiple-layer RSMA scenario. The simulations have demonstrated that the GPI-RS brings significant spectral efficiency gains in various environments while the associated complexity is small compared to the existing WMMSE-SAA method.

For future work, it is promising to consider the finite blocklength regime where a non-zero decoding probability is induced [35], [36]. In particular, decoding failure on a common message can cause significant interference in the SINR of private messages; therefore, precoders need to be designed carefully to maximize the spectral efficiency. In addition, design for physical layer security [37], [38] with RSMA is of interest. Considering RSMA in terahertz line-of-sight MIMO environments [39] is also promising.

## APPENDIX A
### PROOF OF LEMMA 1

We first derive the KKT condition of the problem (40). The corresponding Lagrangian function is defined as

$$L(\bar{\mathbf{f}}) = \log\left(\frac{1}{K}\sum_{k=1}^{K}\exp\left(\log_2\left(\frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(k)\bar{\mathbf{f}}}\right)^{-\frac{1}{\alpha}}\right)\right)^{-\alpha} + \sum_{k=1}^{K}\log_2\left(\frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_k\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_k\bar{\mathbf{f}}}\right). \tag{73}$$

To find a stationary point, we take the partial derivatives of $L(\bar{\mathbf{f}})$ with respect to $\bar{\mathbf{f}}$ and set it to zero. For simplicity, we denote the first and the second part of the Lagrangian function as $L_1(\bar{\mathbf{f}})$ and $L_2(\bar{\mathbf{f}})$, respectively. Since we have

$$\partial \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}\bar{\mathbf{f}}} \right) / \partial \bar{\mathbf{f}}^{\mathsf{H}} = \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}\bar{\mathbf{f}}} \right) \left[ \frac{\mathbf{A}\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}\bar{\mathbf{f}}} - \frac{\mathbf{B}\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}\bar{\mathbf{f}}} \right], \tag{74}$$

the partial derivative of $L_1(\bar{\mathbf{f}})$ is obtained using the above calculation:

$$\frac{\partial L_1(\bar{\mathbf{f}})}{\partial \bar{\mathbf{f}}^{\mathsf{H}}} = \sum_{k=1}^{K} \left[ \frac{\exp \left( \frac{1}{-\alpha} \log_2 \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(k)\bar{\mathbf{f}}} \right) \right)}{\sum_{\ell=1}^{K} \exp \left( \frac{1}{-\alpha} \log_2 \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(\ell)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(\ell)\bar{\mathbf{f}}} \right) \right)} \times \partial \left( \log_2 \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(k)\bar{\mathbf{f}}} \right) \right) / \partial \mathbf{f}^{\mathsf{H}} \right]$$

$$= \frac{1}{\log 2} \sum_{k=1}^{K} \left[ \frac{\exp \left( \frac{1}{-\alpha} \log_2 \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(k)\bar{\mathbf{f}}} \right) \right)}{\sum_{\ell=1}^{K} \exp \left( \frac{1}{-\alpha} \log_2 \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(\ell)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(\ell)\bar{\mathbf{f}}} \right) \right)} \times \left\{ \frac{\mathbf{A}_{\mathsf{c}}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(k)\bar{\mathbf{f}}} - \frac{\mathbf{B}_{\mathsf{c}}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(k)\bar{\mathbf{f}}} \right\} \right] \tag{75}$$

Similar to this, we calculate $L_2(\bar{\mathbf{f}})/\partial \bar{\mathbf{f}}^{\mathsf{H}}$ as follows.

$$\frac{\partial L_2(\bar{\mathbf{f}})}{\partial \bar{\mathbf{f}}^{\mathsf{H}}} = \frac{1}{\log 2} \sum_{k=1}^{K} \left[ \frac{\mathbf{A}_k\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_k\bar{\mathbf{f}}} - \frac{\mathbf{B}_k\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_k\bar{\mathbf{f}}} \right]. \tag{76}$$

The first-order KKT condition holds when

$$\frac{\partial L_1(\bar{\mathbf{f}})}{\partial \bar{\mathbf{f}}^{\mathsf{H}}} + \frac{\partial L_2(\bar{\mathbf{f}})}{\partial \bar{\mathbf{f}}^{\mathsf{H}}} = 0 \tag{77}$$

$$\Leftrightarrow \sum_{k=1}^{K} \left[ \frac{\exp \left( \frac{1}{-\alpha} \log_2 \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(k)\bar{\mathbf{f}}} \right) \right)}{\sum_{\ell=1}^{K} \exp \left( \frac{1}{-\alpha} \log_2 \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(\ell)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(\ell)\bar{\mathbf{f}}} \right) \right)} \times \left\{ \frac{\mathbf{A}_{\mathsf{c}}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(k)\bar{\mathbf{f}}} - \frac{\mathbf{B}_{\mathsf{c}}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(k)\bar{\mathbf{f}}} \right\} \right] + \sum_{k=1}^{K} \left[ \frac{\mathbf{A}_k\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_k\bar{\mathbf{f}}} - \frac{\mathbf{B}_k\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_k\bar{\mathbf{f}}} \right] = 0. \tag{78}$$

Defining $\mathbf{A}_{\mathsf{KKT}}(\bar{\mathbf{f}})$, $\mathbf{B}_{\mathsf{KKT}}(\bar{\mathbf{f}})$, and $\lambda(\bar{\mathbf{f}})$ as

$$\mathbf{A}_{\mathsf{KKT}}(\bar{\mathbf{f}}) = \lambda_{\mathsf{num}}(\bar{\mathbf{f}}) \times \sum_{k=1}^{K} \left[ \frac{\exp \left( \frac{1}{-\alpha} \log_2 \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(k)\bar{\mathbf{f}}} \right) \right)}{\sum_{\ell=1}^{K} \exp \left( \frac{1}{-\alpha} \log_2 \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(\ell)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(\ell)\bar{\mathbf{f}}} \right) \right)} \frac{\mathbf{A}_{\mathsf{c}}(k)}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(k)\bar{\mathbf{f}}} + \frac{\mathbf{A}_k}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_k\bar{\mathbf{f}}} \right], \tag{79}$$

$$\mathbf{B}_{\mathsf{KKT}}(\bar{\mathbf{f}}) = \lambda_{\mathsf{den}}(\bar{\mathbf{f}}) \times \sum_{k=1}^{K} \left[ \frac{\exp \left( \frac{1}{-\alpha} \log_2 \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(k)\bar{\mathbf{f}}} \right) \right)}{\sum_{\ell=1}^{K} \exp \left( \frac{1}{-\alpha} \log_2 \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(\ell)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(\ell)\bar{\mathbf{f}}} \right) \right)} \frac{\mathbf{B}_{\mathsf{c}}(k)}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(k)\bar{\mathbf{f}}} + \frac{\mathbf{B}_k}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_k\bar{\mathbf{f}}} \right], \tag{80}$$

$$\lambda(\bar{\mathbf{f}}) = \left\{ \frac{1}{K} \sum_{k=1}^{K} \exp \left( \log_2 \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_{\mathsf{c}}(k)\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_{\mathsf{c}}(k)\mathbf{f}} \right)^{-\frac{1}{\alpha}} \right) \right\}^{-\frac{\alpha}{\log_2 e}} \times \prod_{k=1}^{K} \left( \frac{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{A}_k\bar{\mathbf{f}}}{\bar{\mathbf{f}}^{\mathsf{H}}\mathbf{B}_k\bar{\mathbf{f}}} \right) = \frac{\lambda_{\mathsf{num}}(\bar{\mathbf{f}})}{\lambda_{\mathsf{den}}(\bar{\mathbf{f}})}, \tag{81}$$

the first-order KKT condition is rearranged as

$$\mathbf{A}_{\mathsf{KKT}}(\bar{\mathbf{f}})\bar{\mathbf{f}} = \lambda(\bar{\mathbf{f}})\mathbf{B}_{\mathsf{KKT}}(\bar{\mathbf{f}})\bar{\mathbf{f}} \Leftrightarrow \mathbf{B}_{\mathsf{KKT}}(\bar{\mathbf{f}})^{-1}\mathbf{A}_{\mathsf{KKT}}(\bar{\mathbf{f}})\bar{\mathbf{f}} = \lambda(\bar{\mathbf{f}})\bar{\mathbf{f}}. \tag{82}$$

This completes the proof. $\qquad\square$

## REFERENCES

[1] J. Park, J. Choi, N. Lee, W. Shin, and H. V. Poor, "Sum spectral efficiency optimization for rate splitting in downlink MU-MISO: A generalized power iteration approach," in *Proc. IEEE Wireless Commun. and Netw. Conf. Workshop*, 2021, pp. 1–6.

[2] G. Caire and S. Shamai, "On the achievable throughput of a multiantenna Gaussian broadcast channel," *IEEE Trans. Inf. Theory*, vol. 49, no. 7, pp. 1691–1706, Jul. 2003.

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

[4] J. Choi, N. Lee, S. Hong, and G. Caire, "Joint user selection, power allocation, and precoding design with imperfect CSIT for multi-cell MU-MIMO downlink systems," *IEEE Trans. Wireless Commun.*, vol. 19, no. 1, pp. 162–176, 2020.

[5] N. Jindal, "MIMO broadcast channels with finite-rate feedback," *IEEE Trans. Inf. Theory*, vol. 52, no. 11, pp. 5045–5060, 2006.

[6] J. Park, N. Lee, J. G. Andrews, and R. W. Heath, "On the optimal feedback rate in interference-limited multi-antenna cellular systems," *IEEE Trans. Wireless Commun.*, vol. 15, no. 8, pp. 5748–5762, 2016.

[7] H. Joudeh and B. Clerckx, "Sum-rate maximization for linearly precoded downlink multiuser MISO systems with partial CSIT: A rate-splitting approach," *IEEE Trans. Commun.*, vol. 64, no. 11, pp. 4847–4861, 2016.

[8] ——, "Robust transmission in downlink multiuser MISO systems: A rate-splitting approach," *IEEE Trans. Signal Process.*, vol. 64, no. 23, pp. 6227–6242, 2016.

[9] M. Dai, B. Clerckx, D. Gesbert, and G. Caire, "A rate splitting strategy for massive MIMO with imperfect CSIT," *IEEE Trans. Wireless Commun.*, vol. 15, no. 7, pp. 4611–4624, 2016.

[10] Z. Li, C. Ye, Y. Cui, S. Yang, and S. Shamai, "Rate splitting for multi-antenna downlink: Precoder design and practical implementation," *IEEE J. Sel. Areas Commun.*, vol. 38, no. 8, pp. 1910–1924, 2020.

[11] Y. Mao and B. Clerckx, "Beyond dirty paper coding for multi-antenna broadcast channel with partial CSIT: A rate-splitting approach," *IEEE Trans. Commun.*, vol. 68, no. 11, pp. 6775–6791, 2020.

[12] Te Han and K. Kobayashi, "A new achievable rate region for the interference channel," *IEEE Trans. Inf. Theory*, vol. 27, no. 1, pp. 49–60, 1981.

[13] R. H. Etkin, D. N. C. Tse, and H. Wang, "Gaussian interference channel capacity to within one bit," *IEEE Trans. Inf. Theory*, vol. 54, no. 12, pp. 5534–5562, 2008.

[14] B. Clerckx, H. Joudeh, C. Hao, M. Dai, and B. Rassouli, "Rate splitting for MIMO wireless networks: A promising PHY-layer strategy for LTE evolution," *IEEE Commun. Mag.*, vol. 54, no. 5, pp. 98–105, 2016.

[15] S. Yang, M. Kobayashi, D. Gesbert, and X. Yi, "Degrees of freedom of time correlated MISO broadcast channel with delayed CSIT," *IEEE Trans. Inf. Theory*, vol. 59, no. 1, pp. 315–328, 2013.

[16] C. Hao, Y. Wu, and B. Clerckx, "Rate analysis of two-receiver MISO broadcast channel with finite rate feedback: A rate-splitting approach," *IEEE Trans. Commun.*, vol. 63, no. 9, pp. 3232–3246, 2015.

[17] A. Z. Yalçın and Y. Yapıcı, "Max-min fair beamforming for cooperative multigroup multicasting with rate-splitting," *IEEE Trans. Wireless Commun.*, vol. 20, no. 1, pp. 254–268, 2021.

[18] O. Dizdar, Y. Mao, W. Han, and B. Clerckx, "Rate-splitting multiple access for downlink multi-antenna communications: Physical layer design and link-level simulations," in *Proc. IEEE Int. Symp. Pers., Indoor Mobile Radio Commun.*, 2020, pp. 1–6.

[19] Z. Yang, M. Chen, W. Saad, and M. Shikh-Bahaei, "Downlink sum-rate maximization for rate splitting multiple access (RSMA)," in *Proc. IEEE Int. Conf. Comm.*, 2020, pp. 1–6.

[20] M. Dai, B. Clerckx, D. Gesbert, and G. Caire, "A rate splitting strategy for massive MIMO with imperfect CSIT," *IEEE Trans. Wireless Commun.*, vol. 15, no. 7, pp. 4611–4624, Jul. 2016.

[21] A. Adhikary, J. Nam, J. Ahn, and G. Caire, "Joint spatial division and multiplexing - The large-scale array regime," *IEEE Trans. Inf. Theory*, vol. 59, no. 10, pp. 6441–6463, Oct. 2013.

[22] Y. Mao, B. Clerckx, and V. O. K. Li, "Rate-splitting for multi-antenna non-orthogonal unicast and multicast transmission: Spectral and energy efficiency analysis," *IEEE Trans. Commun.*, vol. 67, no. 12, pp. 8754–8770, 2019.

[23] A. Papazafeiropoulos, B. Clerckx, and T. Ratnarajah, "Rate-splitting to mitigate residual transceiver hardware impairments in massive MIMO systems," *IEEE Trans. Veh. Technol.*, vol. 66, no. 9, pp. 8196–8211, 2017.

[24] C. Xu, B. Clerckx, S. Chen, Y. Mao, and J. Zhang, "Rate-splitting multiple access for multi-antenna joint radar and communications," *IEEE J. Sel. Areas Commun.*, pp. 1–1, 2021.

[25] J. Zeng, T. Lv, W. Ni, R. P. Liu, N. C. Beaulieu, and Y. J. Guo, "Ensuring max–min fairness of UL SIMO-NOMA: A rate splitting approach," *IEEE Trans. Veh. Technol.*, vol. 68, no. 11, pp. 11 080–11 093, 2019.

[26] C. Hao and B. Clerckx, "MISO networks with imperfect CSIT: A topological rate-splitting approach," *IEEE Trans. Commun.*, vol. 65, no. 5, pp. 2164–2179, 2017.

[27] J. Krivochiza, J. Merlano Duncan, S. Andrenacci, S. Chatzinotas, and B. Ottersten, "FPGA acceleration for computationally efficient symbol-level precoding in multi-user multi-antenna communication systems," *IEEE Access*, vol. 7, pp. 15 509–15 520, 2019.

[28] Y. Cai, L.-H. Zhang, Z. Bai, and R.-C. Li, "On an eigenvector-dependent nonlinear eigenvalue problem," *SIAM J. Matrix Anal. Appl.*, vol. 39, no. 3, pp. 1360–1382, 2018.

[29] H. Yin, D. Gesbert, M. Filippou, and Y. Liu, "A coordinated approach to channel estimation in large-scale multiple-antenna systems," *IEEE J. Sel. Areas Commun.*, vol. 31, no. 2, pp. 264–273, 2013.

[30] P. Patil, B. Dai, and W. Yu, "Hybrid data-sharing and compression strategy for downlink cloud radio access network," *IEEE Trans. Commun.*, vol. 66, no. 11, pp. 5370–5384, Nov. 2018.

[31] C. Shen and H. Li, "On the dual formulation of boosting algorithms," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 12, pp. 2216–2231, 2010.

[32] 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.

[33] G. Peng, L. Liu, P. Zhang, S. Yin, and S. Wei, "Low-computing-load, high-parallelism detection method based on Chebyshev iteration for massive MIMO systems with VLSI architecture," *IEEE Trans. Signal Process.*, vol. 65, no. 14, pp. 3775–3788, 2017.

[34] D. Zhu, B. Li, and P. Liang, "On the matrix inversion approximation based on Neumann series in massive MIMO systems," in *IEEE International Conf. on Commun. (ICC)*, 2015, pp. 1763–1769.

[35] J. Choi and J. Park, "MIMO design for Internet-of-Things: Joint optimization of spectral efficiency and error probability in finite blocklength regime," *IEEE Internet of Things J.*, pp. 1–1, 2021.

[36] Y. Polyanskiy, H. V. Poor, and S. Verdu, "Channel coding rate in the finite blocklength regime," *IEEE Trans. Inf. Theory*, vol. 56, no. 5, pp. 2307–2359, May 2010.

[37] K. Lee, J. Choi, D. K. Kim, and J. Park, "Secure transmission for hierarchical information accessibility in downlink MU-MIMO," *ArXiv*, 2021. [Online]. Available: https://arxiv.org/abs/2109.07727

[38] J. Choi and J. Park, "Sum secrecy spectral efficiency maximization in downlink MU-MIMO: Colluding eavesdroppers," *IEEE Trans. Veh. Technol.*, vol. 70, no. 1, pp. 1051–1056, 2021.

[39] H. Do, S. Cho, J. Park, H.-J. Song, N. Lee, and A. Lozano, "Terahertz line-of-sight MIMO communication: Theory and practical challenges," *IEEE Commun. Mag.*, vol. 59, no. 3, pp. 104–109, 2021.