

# Multiuser Uplink MIMO Communications Assisted by Multiple Reconfigurable Intelligent Surfaces

Yang Lv, *Member, IEEE*, Zhiqiang He, *Member, IEEE*, and Yue Rong, *Senior Member, IEEE*

**Abstract**—This letter focuses on the design of multiuser uplink multiple-input multiple-output (MIMO) communications assisted by multiple reconfigurable intelligent surfaces (RISs), where we consider both linear minimal mean-squared error (MMSE) and nonlinear MMSE-decision feedback equalization (DFE) receivers. The RISs, consisting of massive low-cost passive elements, can reflect the incident signals with either continuous or discrete phase shifts. Adopting the sum mean-squared error (MSE) minimization criterion, we jointly optimize the RIS phase shifts, the linear or nonlinear receiving matrices, and the source precoding matrices in iterative manner. Simulation results show that, when the signals in direct user-access point (AP) links suffer severe attenuation, the introduction of indirect user-RIS-AP links can significantly improve the system performance.

**Index Terms**—Reconfigurable intelligent surfaces, MIMO communications, multiuser, MMSE, DFE, semidefinite relaxation.

## I. INTRODUCTION

IN recent years, for the ability to considerably enhance the quality of wireless communications along with high spectrum and energy efficiency, reconfigurable intelligent surfaces (RISs) [1], also called intelligent reflecting surfaces [2], etc., have drawn wide attention in academia and industry. An RIS is generally a planar array of low-cost passive elements, each of which can independently induce an adjustable phase shift and/or amplitude change to its incident electromagnetic wave, thereby reconfiguring the characteristics of wireless channels and making the concept of “smart radio environments” emerge.

Considering a downlink single-user multiple-input single-output (MISO) system, [3] investigated the joint design of the active beamforming at the access point (AP) and the passive beamforming at the RIS. References [4] and [5] then extended the single-user scenario in [3] to the multiuser ones with continuous and discrete RIS phase shifts, respectively, and both of them focused on minimizing the transmission power of the AP. Targeting at maximizing all users’ weighted sum-rate, [6] studied both perfect and imperfect channel state information (CSI) circumstances. In [7], the direct links between the base station and users were neglected due to the poor conditions for signal propagation, and two energy efficiency maximization algorithms were presented. Towards a single-user uplink single-input multiple-output (SIMO) system, [8] developed a passive beamforming and information transfer (PBIT) technique where

the RIS utilizes the on/off states of passive elements to carry its private data. Such technique was further expanded in [9] to apply to multiuser multiple-input multiple-output (MIMO) systems. To guard against the eavesdroppers, [10] addressed the issue of secure communications with the aid of RISs.

This letter investigates the design of multiple RISs-assisted multiuser uplink MIMO communications with both continuous and discrete RIS phase shifts take into consideration, where either the linear minimal mean-squared error (MMSE) receiver or the nonlinear MMSE-decision feedback equalization (DFE) receiver [11] can be adopted at the AP. Our design criterion is to minimize the sum mean-squared error (MSE) for the signal waveform estimation of all data streams and the transmission power of each user is limited. To solve the formulated nonconvex problem, we employ the block coordinate decent (BCD) method [12, Sec. 2.7] to iteratively optimize the signal phase shifts at all RISs, the linear or nonlinear receiving matrices at the AP, and the source precoding matrices at all users, where the semidefinite relaxation (SDR) technique [13], the QR factorization [14, Sec. 2.1], and the Karush-Kuhn-Tucker (KKT) conditions [15] are respectively used. Simulation results show that, when direct user-AP communication links are unfavorable, in comparison to the systems without RISs, significant performance improvement can be observed for RISs-assisted systems, and such advantage becomes more prominent as the number of passive elements in each RIS increases.

*Notations:*  $\mathbb{C}^n$  is the space of complex column vectors with dimension  $n$  and  $\mathbb{C}^{m \times n}$  is that of  $m$ -by- $n$  complex matrices.  $\triangleq$  denotes “defined as”.  $\sim$  signifies “distributed as”.  $j \triangleq \sqrt{-1}$  is the imaginary unit.  $(a)_{m \times n}$  denotes an  $m$ -by- $n$  matrix with all its elements being scalar  $a$ .  $(\cdot)^*$ ,  $|\cdot|$ , and  $\text{angle}(\cdot)$  stand for the complex conjugate, modulus, and phase angle of a scalar, respectively.  $\mathbf{I}_n$  is the  $n$ th-order identity matrix.  $(\cdot)^T$  and  $(\cdot)^H$  represent the transpose and Hermitian transpose of a vector or matrix, respectively.  $\text{vec}(\mathbf{X})$  is a column vector made up of stacked columns of matrix  $\mathbf{X}$  [16].  $[\mathbf{a}]_{1:n}$  denotes a subvector of vector  $\mathbf{a}$ , containing its first  $n$  elements.  $[\mathbf{X}]_{m,n}$  indicates the  $m$ th row and  $n$ th column element of  $\mathbf{X}$ .  $\text{diag}(\cdot)$  and  $\text{bd}(\cdot)$  are a diagonal and a block diagonal matrix, respectively, with the diagonal elements given in parentheses.  $\otimes$  is the Kronecker product operator.  $\text{tr}(\cdot)$ ,  $(\cdot)^{-1}$ ,  $(\cdot)^\dagger$ , and  $(\cdot)^{\frac{1}{2}}$  stand for the trace, inverse, Moore-Penrose inverse, and square root of a matrix [14], respectively.  $\mathbf{X} \succeq \mathbf{0}$  means that  $\mathbf{X}$  is Hermitian positive semidefinite (PSD). The distribution of a circularly symmetric complex Gaussian random variable (vector) having zero mean and variance  $a$  (covariance matrix  $\mathbf{X}$ ) is denoted by  $\mathcal{CN}(0, a)$  ( $\mathcal{CN}(\mathbf{0}, \mathbf{X})$ ).  $\mathbb{E}[\cdot]$  stands for the statistical expectation.

This work was supported by the National Key Research and Development Program of China under Grant 2018YFB1801101.

Yang Lv and Zhiqiang He are with the Key Laboratory of Universal Wireless Communications, Ministry of Education, Beijing University of Posts and Telecommunications, Beijing 100876, China (e-mail: lvyangcn@foxmail.com; hezq@bupt.edu.cn).

Yue Rong is with the School of Electrical Engineering, Computing and Mathematical Sciences, Curtin University, Bentley, WA 6102, Australia (e-mail: y.rong@curtin.edu.au).

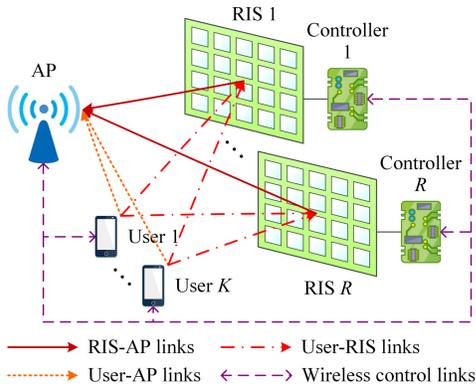


Fig. 1. Multiuser uplink MIMO communications assisted by multiple RISs.

## II. SYSTEM MODEL AND PROBLEM FORMULATION

As shown in Fig. 1, we consider an  $R$  RISs-assisted uplink MIMO communication system that serves  $K$  users. For  $k = 1, \dots, K$ , the  $k$ th user transmits  $N_k$  data streams via  $N_{s,k}$  antennas ( $N_{s,k} \geq N_k$ ) with the source precoding matrix  $\mathbf{B}_k \in \mathbb{C}^{N_{s,k} \times N_k}$ . By assuming that the modulated signal vector  $\mathbf{s}_k \in \mathbb{C}^{N_k}$  satisfies  $\mathbb{E}[\mathbf{s}_k] = (\mathbf{0})_{N_k \times 1}$  and  $\mathbb{E}[\mathbf{s}_k \mathbf{s}_k^H] = \mathbf{I}_{N_k}$ , the mean power of the transmitted signal vector  $\mathbf{x}_k = \mathbf{B}_k \mathbf{s}_k$  is derived as  $\text{tr}(\mathbb{E}[\mathbf{x}_k \mathbf{x}_k^H]) = \text{tr}(\mathbf{B}_k \mathbf{B}_k^H)$ , which should not be greater than the power budget  $p_k$ . Additionally, the AP is equipped with  $N_d$  antennas. For  $r = 1, \dots, R$ , the  $r$ th RIS consists of  $M_r$  passive elements, each behaving like a keyhole [3], which receives superposed multi-path signals at a physical point and then scatters the combined signal as if from a point source [2]. This kind of RIS is called the reflectarray-based RIS in [1]. Let the passive reflection matrix of the  $r$ th RIS be written as

$$\Theta_r \triangleq \text{diag}(\beta_{r,1} e^{j\theta_{r,1}}, \dots, \beta_{r,M_r} e^{j\theta_{r,M_r}}) \quad (1)$$

where for  $m_r = 1, \dots, M_r$ ,  $\beta_{r,m_r} \in [0, 1]$  and  $\theta_{r,m_r} \in [0, 2\pi)$  are, respectively, the amplitude coefficient and the phase shift generated by the  $m_r$ th passive element for its combined signal. As in [5], [9], [10], for simplicity and maximizing the signal reflection power of RISs, we set  $\beta_{r,m_r} = 1, \forall r, m_r$ . Besides,  $\theta_{r,m_r}, \forall r, m_r$  can take any values in  $[0, 2\pi)$  for continuous RIS phase shifts, and if there are only  $b$  bits of discrete RIS phase shifts due to the limitation of hardware costs, we have  $\theta_{r,m_r} \in \mathcal{L} \triangleq \left\{0, \frac{2\pi}{L}, \dots, \frac{2\pi(L-1)}{L}\right\}, \forall r, m_r$  with  $L \triangleq 2^b$ . Note that, compared to a single RIS, multiple RISs can provide more passive elements, along with additional spatial diversity, and thus can further improve the system performance.

Each RIS is connected to a controller that has multiple responsibilities. For example, the  $r$ th controller can (a) exchange information with other devices via dedicated wireless control links to coordinate communications, (b) regulate the switching between two working modes of the  $r$ th RIS, i.e., the reflecting mode for data transmission and the receiving mode for CSI estimation, and (c) adjust the phase shifts of the  $r$ th RIS once  $\Theta_r$  is updated [3]–[5]. The baseband equivalent channels for the  $k$ th user-AP, the  $k$ th user- $r$ th RIS, and the  $r$ th RIS-AP data links are denoted by  $\mathbf{H}_{sd,k} \in \mathbb{C}^{N_d \times N_{s,k}}$ ,  $\mathbf{H}_{s,rk} \in \mathbb{C}^{M_r \times N_{s,k}}$ , and  $\mathbf{H}_{d,r} \in \mathbb{C}^{N_d \times M_r}$ , respectively, which are assumed to be quasi-static flat-fading and perfectly known through various CSI acquisition means as discussed in [1]–[4]. Here, we can

equip the RISs with receiving radio frequency (RF) chains. After all users send orthogonal pilot signals, the AP and the RIS controllers can estimate  $\mathbf{H}_{sd,k}$  and  $\mathbf{H}_{s,rk}$ , respectively. After the AP sends a pilot signal, the RIS controllers can estimate  $\mathbf{H}_{d,r}$  by exploiting the channel reciprocity. Via control links, all CSI can be collected at one of the RIS controllers, which is selected to execute our proposed algorithms and deliver the optimized system parameters to their corresponding devices.

According to [6], [8], [9], the received signal at the AP is the superposition of different components coming from two types of communication links, i.e., the direct links and the RIS reflection links, among which the time delay difference is assumed to be negligible. Specifically, for  $k = 1, \dots, K$  and  $r = 1, \dots, R$ ,  $\mathbf{H}_{sd,k} \mathbf{x}_k$  is the signal component received directly from the  $k$ th user, and  $\mathbf{H}_{d,r} \Theta_r \mathbf{H}_{s,rk} \mathbf{x}_k$  is that received indirectly from the  $k$ th user via the reflection of the  $r$ th RIS. Hence, the total signal received at the AP is given by  $\mathbf{y} = \sum_{k=1}^K (\mathbf{H}_{sd,k} + \sum_{r=1}^R \mathbf{H}_{d,r} \Theta_r \mathbf{H}_{s,rk}) \mathbf{B}_k \mathbf{s}_k + \mathbf{n}$ , where  $\mathbf{n} \sim \mathcal{CN}(\mathbf{0}, \sigma^2 \mathbf{I}_{N_d})$  denotes the additive white Gaussian noise (AWGN) vector at the AP. Moreover, via setting  $N \triangleq \sum_{k=1}^K N_k$ ,  $N_s \triangleq \sum_{k=1}^K N_{s,k}$ , and  $M \triangleq \sum_{r=1}^R M_r$  as well as  $\mathbf{s} \triangleq [\mathbf{s}_1^T, \dots, \mathbf{s}_K^T]^T \in \mathbb{C}^N$ ,  $\mathbf{B} \triangleq \text{bd}(\mathbf{B}_1, \dots, \mathbf{B}_K) \in \mathbb{C}^{N_s \times N}$ ,  $\mathbf{H}_{sd} \triangleq [\mathbf{H}_{sd,1}, \dots, \mathbf{H}_{sd,K}] \in \mathbb{C}^{N_d \times N_s}$ ,  $\mathbf{H}_d \triangleq [\mathbf{H}_{d,1}, \dots, \mathbf{H}_{d,R}] \in \mathbb{C}^{N_d \times M}$ ,  $\Theta \triangleq \text{bd}(\Theta_1, \dots, \Theta_R) \in \mathbb{C}^{M \times M}$ ,  $\mathbf{H}_{s,k} \triangleq [\mathbf{H}_{s,1k}^T, \dots, \mathbf{H}_{s,Rk}^T]^T \in \mathbb{C}^{M \times N_{s,k}}$  for  $k = 1, \dots, K$ , and  $\mathbf{H}_s \triangleq [\mathbf{H}_{s,1}, \dots, \mathbf{H}_{s,K}] \in \mathbb{C}^{M \times N_s}$ , we have

$$\mathbf{y} = (\mathbf{H}_{sd} + \mathbf{H}_d \Theta \mathbf{H}_s) \mathbf{B} \mathbf{s} + \mathbf{n} = \mathbf{A} \mathbf{s} + \mathbf{n} \quad (2)$$

with  $\mathbf{A} \triangleq (\mathbf{H}_{sd} + \mathbf{H}_d \Theta \mathbf{H}_s) \mathbf{B} \in \mathbb{C}^{N_d \times N}$ . For convenience, we further make  $\Theta = \text{diag}(e^{j\theta_1}, \dots, e^{j\theta_M})$  and  $\theta \triangleq [\theta_1, \dots, \theta_M]^T$ , where for  $m = 1, \dots, M$ ,  $\theta_m \in [0, 2\pi)$  or  $\theta_m \in \mathcal{L}$ , and for  $r = 1, \dots, R$  and  $m_r = 1, \dots, M_r$ ,  $\theta_{r,m_r}$  is just the  $(\sum_{\zeta=0}^{r-1} M_\zeta + m_r)$ th element of  $\theta$  with  $M_0 = 0$ .

At the AP, when we use the linear MMSE receiver, the source signal vector  $\mathbf{s}$  is estimated as  $\bar{\mathbf{s}} = \bar{\mathbf{W}}^H \mathbf{y}$ , where  $\bar{\mathbf{W}} \in \mathbb{C}^{N_d \times N}$  is the linear receiving matrix and  $\bar{\mathbf{s}}$  is the linearly estimated signal vector. Thus the sum MSE for estimating all symbols in  $\mathbf{s}$  can be expressed as  $\bar{E}_s \triangleq \text{tr}(\mathbb{E}[(\bar{\mathbf{s}} - \mathbf{s})(\bar{\mathbf{s}} - \mathbf{s})^H]) = \text{tr}((\bar{\mathbf{W}}^H \mathbf{A} - \mathbf{I}_N)(\bar{\mathbf{W}}^H \mathbf{A} - \mathbf{I}_N)^H + \sigma^2 \bar{\mathbf{W}}^H \bar{\mathbf{W}})$ . When we use the nonlinear MMSE-DFE receiver instead,  $\mathbf{s}$  is estimated as  $\hat{\mathbf{s}} = \mathbf{W}^H \mathbf{y} - \mathbf{D} \hat{\mathbf{s}}$ , where  $\mathbf{W} \in \mathbb{C}^{N_d \times N}$  is the decision feed-forward matrix with its  $i$ th column vector denoted by  $\mathbf{w}_i$  for  $i = 1, \dots, N$ ,  $\mathbf{D} \in \mathbb{C}^{N \times N}$  is the decision feedback matrix which is strictly upper-triangular,  $\tilde{\mathbf{s}} = [\tilde{s}_1, \dots, \tilde{s}_N]^T$  is the detected signal vector, and  $\hat{\mathbf{s}} = [\hat{s}_1, \dots, \hat{s}_N]^T$  is the nonlinearly estimated signal vector. Specifically, we first estimate the  $N$ th symbol of  $\mathbf{s}$  as  $\hat{s}_N = \mathbf{w}_N^H \mathbf{y}$ , which is then detected as  $\tilde{s}_N$ . Following that, we estimate the  $i$ th symbol of  $\mathbf{s}$  for  $i = N-1, \dots, 1$ , one by one, as  $\hat{s}_i = \mathbf{w}_i^H \mathbf{y} - \sum_{l=i+1}^N [\mathbf{D}]_{i,l} \tilde{s}_l$ , which is then detected as  $\tilde{s}_i$ . Besides, like in [11], during our mathematical derivations, we assume that there is no error propagation within the MMSE-DFE receiver, i.e.,  $\tilde{\mathbf{s}} = \mathbf{s}$ , thus we obtain  $\hat{\mathbf{s}} = \mathbf{W}^H \mathbf{y} - \mathbf{D} \mathbf{s}$ . Consequently, the sum MSE is now given by  $E_s \triangleq \text{tr}(\mathbb{E}[(\hat{\mathbf{s}} - \mathbf{s})(\hat{\mathbf{s}} - \mathbf{s})^H]) = \text{tr}((\mathbf{W}^H \mathbf{A} - \mathbf{U})(\mathbf{W}^H \mathbf{A} - \mathbf{U})^H + \sigma^2 \mathbf{W}^H \mathbf{W})$ , where  $\mathbf{U} \triangleq \mathbf{D} + \mathbf{I}_N$  is also called the decision feedback matrix hereafter.

At this point, with the above nonlinear receiving structure, we can formulate the system optimization problem as follows.

$$\min_{\{\theta_m\}, \{\mathbf{W}, \mathbf{U}\}, \{\mathbf{B}_k\}} \text{tr} \left( (\mathbf{W}^H \mathbf{A} - \mathbf{U}) \times (\mathbf{W}^H \mathbf{A} - \mathbf{U})^H + \sigma^2 \mathbf{W}^H \mathbf{W} \right) \quad (3)$$

$$\text{s.t. } \theta_m \in [0, 2\pi) \text{ or } \theta_m \in \mathcal{L}, \quad m = 1, \dots, M, \quad (4)$$

$$[\mathbf{U}]_{i,l} = \begin{cases} 0, & i > l, \\ 1, & i = l, \end{cases} \quad i, l = 1, \dots, N, \quad (5)$$

$$\text{tr}(\mathbf{B}_k \mathbf{B}_k^H) \leq p_k, \quad k = 1, \dots, K \quad (6)$$

where  $\{\theta_m\} \triangleq \{\theta_1, \dots, \theta_M\}$  and  $\{\mathbf{B}_k\} \triangleq \{\mathbf{B}_1, \dots, \mathbf{B}_K\}$ . Noteworthily, when  $\mathbf{W}$  is substituted by  $\bar{\mathbf{W}}$  and  $\mathbf{U}$  is set to be  $\mathbf{I}_N$ , the problem (3)–(6) turns into the one with linear receiving structure. Besides, since only the phase angles of the diagonal elements of  $\Theta_r$  are adjustable and the RISs can not increase the power of their reflected signals, it is more challenging to design the systems with RISs than with amplify-and-forward relays [11]. In the next section, following the BCD method, we shall develop iterative algorithms to solve the problem (3)–(6).

### III. ALGORITHM DESIGN

Here, the three blocks of variables for the problem (3)–(6), i.e.,  $\{\theta_m\}$ ,  $\{\mathbf{W}, \mathbf{U}\}$ , and  $\{\mathbf{B}_k\}$ , are iteratively optimized one by one with the other blocks fixed. At the beginning of such iterative procedures, we initialize  $\mathbf{W} = (\mathbf{1})_{N_d \times N}$ ,  $\mathbf{U} = \mathbf{I}_N$ , and  $\mathbf{B}_k = [\sqrt{p_k/N_k} \mathbf{I}_{N_k}, (\mathbf{0})_{N_k \times (N_s, k - N_k)}]^T$  for  $k = 1, \dots, K$ .

Firstly, with fixed  $\{\mathbf{W}, \mathbf{U}\}$  and  $\{\mathbf{B}_k\}$ , we optimize  $\{\theta_m\}$  by solving the problem as below.

$$\min_{\{\theta_m\}} \text{tr} \left( (\mathbf{W}^H \mathbf{A} - \mathbf{U})(\mathbf{W}^H \mathbf{A} - \mathbf{U})^H \right) \quad (7)$$

$$\text{s.t. } \theta_m \in [0, 2\pi) \text{ or } \theta_m \in \mathcal{L}, \quad m = 1, \dots, M \quad (8)$$

which, via  $\text{tr}(\mathbf{X}\mathbf{X}^H) = \text{vec}(\mathbf{X})^H \text{vec}(\mathbf{X})$  and  $\text{vec}(\mathbf{X}\mathbf{Y}\mathbf{Z}) = (\mathbf{Z}^T \otimes \mathbf{X}) \text{vec}(\mathbf{Y})$  [16, Lemma 4.3.1], can be rewritten as

$$\min_{\{\theta_m\}} \text{vec}(\Theta)^H \mathbf{C} \text{vec}(\Theta) + \text{vec}(\Theta)^H \mathbf{c} + \mathbf{c}^H \text{vec}(\Theta) \quad (9)$$

$$\text{s.t. } \theta_m \in [0, 2\pi) \text{ or } \theta_m \in \mathcal{L}, \quad m = 1, \dots, M \quad (10)$$

with  $\mathbf{C} \triangleq [(\mathbf{H}_s \mathbf{B})^T \otimes (\mathbf{W}^H \mathbf{H}_d)]^H [(\mathbf{H}_s \mathbf{B})^T \otimes (\mathbf{W}^H \mathbf{H}_d)]$  and  $\mathbf{c} \triangleq [(\mathbf{H}_s \mathbf{B})^T \otimes (\mathbf{W}^H \mathbf{H}_d)]^H \text{vec}(\mathbf{W}^H \mathbf{H}_{sd} \mathbf{B} - \mathbf{U})$ . Now, by setting  $\mathbf{v} \triangleq [v_1, \dots, v_M]^T$  along with  $v_m \triangleq e^{j\theta_m}$  for  $m = 1, \dots, M$  and  $\mathcal{V} \triangleq \{1, e^{j\frac{2\pi}{L}}, \dots, e^{j\frac{2\pi(L-1)}{L}}\}$ , we are able to reformulate the problem (9)–(10) as

$$\min_{\mathbf{v}} \mathbf{v}^H \mathbf{C}_0 \mathbf{v} + \mathbf{v}^H \mathbf{c}_0 + \mathbf{c}_0^H \mathbf{v} \quad (11)$$

$$\text{s.t. } |v_m| = 1 \text{ or } v_m \in \mathcal{V}, \quad m = 1, \dots, M \quad (12)$$

where  $\mathbf{C}_0 \in \mathbb{C}^{M \times M}$  is made up of all the elements at the intersections of the  $[(i-1)M+i]$ th rows and the  $[(l-1)M+l]$ th columns of matrix  $\mathbf{C}$  for  $i, l = 1, \dots, M$  and  $\mathbf{c}_0 \in \mathbb{C}^M$  is made up of all the  $[(m-1)M+m]$ th elements of vector  $\mathbf{c}$  for  $m = 1, \dots, M$ . To homogenize the above nonconvex quadratic programming problem, we bring in an auxiliary variable  $\alpha$  and gives the following problem with  $\hat{\mathbf{v}} \triangleq [\hat{v}_1, \dots, \hat{v}_M]^T$ .

$$\min_{\hat{\mathbf{v}}, \alpha} \hat{\mathbf{v}}^H \mathbf{C}_0 \hat{\mathbf{v}} + \hat{\mathbf{v}}^H \mathbf{c}_0 \alpha + \alpha^* \mathbf{c}_0^H \hat{\mathbf{v}} \quad (13)$$

$$\text{s.t. } |\hat{v}_m| = 1, |\alpha| = 1 \text{ or } \hat{v}_m, \alpha \in \mathcal{V}, \quad m = 1, \dots, M. \quad (14)$$

It can be readily observed that, if  $\hat{\mathbf{v}} \triangleq [\hat{v}^T, \alpha]^T$  solves the problem (13)–(14), then  $\mathbf{v} = \alpha^* \hat{\mathbf{v}}$  solves the problem (11)–(12). Via further setting  $\hat{\mathbf{C}}_0 \triangleq \begin{bmatrix} \mathbf{C}_0 & \mathbf{c}_0 \\ \mathbf{c}_0^H & 0 \end{bmatrix}$  and  $\hat{\mathbf{v}} = [\hat{v}_1, \dots, \hat{v}_M, \hat{v}_{M+1}]^T$ , we can rewrite the problem (11)–(12) as

$$\min_{\hat{\mathbf{v}}} \hat{\mathbf{v}}^H \hat{\mathbf{C}}_0 \hat{\mathbf{v}} \quad (15)$$

$$\text{s.t. } |\hat{v}_m| = 1 \text{ or } \hat{v}_m \in \mathcal{V}, \quad m = 1, \dots, M, M+1 \quad (16)$$

of which the objective  $\hat{\mathbf{v}}^H \hat{\mathbf{C}}_0 \hat{\mathbf{v}}$  is equivalent to  $\text{tr}(\hat{\mathbf{C}}_0 \hat{\mathbf{v}} \hat{\mathbf{v}}^H) = \text{tr}(\hat{\mathbf{C}}_0 \hat{\mathbf{V}})$  with  $\hat{\mathbf{V}} \triangleq \hat{\mathbf{v}} \hat{\mathbf{v}}^H \in \mathbb{C}^{(M+1) \times (M+1)}$  required to satisfy  $\hat{\mathbf{V}} \succeq \mathbf{0}$  and  $\text{rank}(\hat{\mathbf{V}}) = 1$ . Since the rank-one constraint of  $\hat{\mathbf{V}}$  is nonconvex, by applying the SDR technique, we leave it out to acquire the following relaxed version of the problem (15)–(16) for both continuous and discrete RIS phase shifts.

$$\min_{\hat{\mathbf{V}}} \text{tr}(\hat{\mathbf{C}}_0 \hat{\mathbf{V}}) \quad (17)$$

$$\text{s.t. } \hat{\mathbf{V}} \succeq \mathbf{0}, [\hat{\mathbf{V}}]_{m,m} = 1, \quad m = 1, \dots, M, M+1. \quad (18)$$

The above is a convex semidefinite programming (SDP) problem that can be optimally solved via, e.g., the convex optimization software: CVX [17]. Because the optimal solution of the problem (17)–(18) may not meet  $\text{rank}(\hat{\mathbf{V}}) = 1$ , the optimal objective value of the problem (17)–(18) only provides a lower bound of that of the problem (15)–(16). Here, to get  $\hat{\mathbf{v}}$  from  $\hat{\mathbf{V}}$ , we adopt the Gaussian randomization method [13] as below.

On the basis of the eigenvalue decomposition (EVD) of  $\hat{\mathbf{V}}$ , i.e.,  $\hat{\mathbf{V}} = \mathbf{Y} \mathbf{\Lambda} \mathbf{Y}^H$  with  $\mathbf{Y}, \mathbf{\Lambda} \in \mathbb{C}^{(M+1) \times (M+1)}$  being a unitary and a diagonal matrix, respectively, we generate  $T$  random vectors as  $\mathbf{z}^{(t)} = \mathbf{Y} \mathbf{\Lambda}^{\frac{1}{2}} \boldsymbol{\beta}^{(t)} = [z_1^{(t)}, \dots, z_M^{(t)}, z_{M+1}^{(t)}]^T$ , where  $\boldsymbol{\beta}^{(t)} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}_{M+1})$  for  $t = 1, \dots, T$ . Then, for constructing the feasible points  $\hat{\mathbf{v}}^{(t)} = [\hat{v}_1^{(t)}, \dots, \hat{v}_M^{(t)}, \hat{v}_{M+1}^{(t)}]^T$ , in the case of continuous RIS phase shifts, we set  $\hat{v}_m^{(t)} = z_m^{(t)} / |z_m^{(t)}|$  for  $m = 1, \dots, M, M+1$ , meanwhile, in the case of discrete RIS phase shifts, we set  $\hat{v}_m^{(t)} = e^{j2\pi\xi/L}$  if  $\text{angle}(z_m^{(t)}) \in [(2\xi-1)\pi/L, (2\xi+1)\pi/L]$  with  $\xi \in \{0, 1, \dots, L-1\}$ . Among all  $\hat{\mathbf{v}}^{(t)}$ , we determine  $\hat{t} = \text{argmin}_{t=1, \dots, T} (\hat{\mathbf{v}}^{(t)})^H \hat{\mathbf{C}}_0 \hat{\mathbf{v}}^{(t)}$  and set  $\hat{\mathbf{v}} = \hat{\mathbf{v}}^{(\hat{t})}$  as the approximate solution of the problem (15)–(16). Following that, the original problem (11)–(12) can thus be approximately solved by  $\mathbf{v} = \hat{v}_{M+1}^* [\hat{\mathbf{v}}]_{1:M}$ . As pointed out in [13], the accuracy of the obtained solution depends on the number of randomizations  $T$ . Via simulation tests, setting  $T = 100(M+1)$  is deemed appropriate here.

Secondly, with fixed  $\{\mathbf{B}_k\}$  and  $\{\theta_m\}$ , we optimize  $\{\mathbf{W}, \mathbf{U}\}$  by using the QR factorization  $\begin{bmatrix} \sqrt{1/\sigma^2} \mathbf{A} \\ \mathbf{I}_N \end{bmatrix} = \mathbf{Q} \mathbf{R} = \begin{bmatrix} \hat{\mathbf{Q}} \\ \hat{\mathbf{Q}} \end{bmatrix} \mathbf{R}$ ,

where  $\mathbf{Q} \in \mathbb{C}^{(N_d+N) \times N}$  consists of orthonormal columns and its first  $N_d$  rows and last  $N$  rows are denoted by  $\hat{\mathbf{Q}}$  and  $\hat{\mathbf{Q}}$ , respectively, and  $\mathbf{R} \in \mathbb{C}^{N \times N}$  is upper-triangular with positive diagonal elements. From Theorem 1 in [11], along with  $\mathbf{D}_R \triangleq \text{diag}([\mathbf{R}]_{1,1}, \dots, [\mathbf{R}]_{N,N})$ , the optimal  $\mathbf{W}$  and  $\mathbf{U}$  are given by  $\mathbf{W} = \sqrt{1/\sigma^2} \hat{\mathbf{Q}} \mathbf{D}_R^{-H}$  and  $\mathbf{U} = \mathbf{D}_R^{-1} \mathbf{R}$ .

Thirdly, with fixed  $\{\theta_m\}$  and  $\{\mathbf{W}, \mathbf{U}\}$ , we optimize  $\{\mathbf{B}_k\}$  via solving the problem as below.

$$\min_{\{\mathbf{B}_k\}} \text{tr}((\mathbf{F} \mathbf{B} - \mathbf{U})(\mathbf{F} \mathbf{B} - \mathbf{U})^H) \quad (19)$$

$$\text{s.t. } \text{tr}(\mathbf{B}_k \mathbf{B}_k^H) \leq p_k, \quad k = 1, \dots, K \quad (20)$$

TABLE I  
PROCEDURES OF THE RIS(-b)-DFE ALGORITHM

- 1) Set  $n=0$ . Initialize  $\{\mathbf{W}, \mathbf{U}\}^{\langle n \rangle}$  and  $\{\mathbf{B}_k\}^{\langle n \rangle}$  as  $\mathbf{W} = (\mathbf{1})_{N_d \times N}$ ,  $\mathbf{U} = \mathbf{I}_N$ , and  $\mathbf{B}_k = [\sqrt{p_k/N_k} \mathbf{I}_{N_k}, (\mathbf{0})_{N_k \times (N_{s,k} - N_k)}]^T$ ,  $k = 1, \dots, K$ .
- 2) Obtain  $\{\theta_m\}^{\langle n+1 \rangle}$  with fixed  $\{\mathbf{W}, \mathbf{U}\}^{\langle n \rangle}$  and  $\{\mathbf{B}_k\}^{\langle n \rangle}$  by solving the problem (7)–(8) for either continuous or discrete RIS phase shifts.
- 3) Obtain  $\{\mathbf{W}, \mathbf{U}\}^{\langle n+1 \rangle}$  with fixed  $\{\mathbf{B}_k\}^{\langle n \rangle}$  and  $\{\theta_m\}^{\langle n+1 \rangle}$  by computing  $\mathbf{W} = \sqrt{1/\sigma^2} \dot{\mathbf{Q}} \mathbf{D}_R^{-H}$  and  $\mathbf{U} = \mathbf{D}_R^{-1} \mathbf{R}$ .
- 4) Obtain  $\{\mathbf{B}_k\}^{\langle n+1 \rangle}$  with fixed  $\{\theta_m\}^{\langle n+1 \rangle}$  and  $\{\mathbf{W}, \mathbf{U}\}^{\langle n+1 \rangle}$  by solving the problem (19)–(20).
- 5) Let  $n = n+1$ . Continue to Step 2) if  $n < 20$ . Otherwise, end the algorithm.

where  $\mathbf{F} \triangleq \mathbf{W}^H (\mathbf{H}_{sd} + \mathbf{H}_d \boldsymbol{\Theta} \mathbf{H}_s) \in \mathbb{C}^{N \times N_s}$ . Now, we further set  $\mathbf{F} = [\mathbf{F}_1, \dots, \mathbf{F}_K]$  and  $\mathbf{U} = [\mathbf{U}_1, \dots, \mathbf{U}_K]$ . Here, for  $k = 1, \dots, K$ ,  $N_{s,k}$  columns of  $\mathbf{F}$ , from the  $(\sum_{i=0}^{k-1} N_{s,i} + 1)$ th to the  $(\sum_{i=0}^k N_{s,i})$ th one with  $N_{s,0} = 0$ , together make up  $\mathbf{F}_k \in \mathbb{C}^{N \times N_{s,k}}$ , and  $N_k$  columns of  $\mathbf{U}$ , from the  $(\sum_{i=0}^{k-1} N_i + 1)$ th to the  $(\sum_{i=0}^k N_i)$ th one with  $N_0 = 0$ , make up  $\mathbf{U}_k \in \mathbb{C}^{N \times N_k}$ . Thus, we can decompose the problem (19)–(20) into

$$\min_{\mathbf{B}_k} \text{tr}((\mathbf{F}_k \mathbf{B}_k - \mathbf{U}_k)(\mathbf{F}_k \mathbf{B}_k - \mathbf{U}_k)^H) \quad (21)$$

$$\text{s.t. } \text{tr}(\mathbf{B}_k \mathbf{B}_k^H) \leq p_k \quad (22)$$

for  $k = 1, \dots, K$ . With the Lagrange multiplier  $\mu_k$ , the KKT conditions of the above problem are  $\mu_k \geq 0$ ,  $\text{tr}(\mathbf{B}_k \mathbf{B}_k^H) \leq p_k$ ,  $\mu_k [\text{tr}(\mathbf{B}_k \mathbf{B}_k^H) - p_k] = 0$ , and  $(\mathbf{F}_k^H \mathbf{F}_k + \mu_k \mathbf{I}_{N_{s,k}}) \mathbf{B}_k = \mathbf{F}_k^H \mathbf{U}_k$ , which can be solved under two possible cases. One case is  $\mu_k = 0$ , where we have  $\mathbf{B}_k = (\mathbf{F}_k^H \mathbf{F}_k)^\dagger \mathbf{F}_k^H \mathbf{U}_k$  and it is just the optimal solution so long as  $\text{tr}(\mathbf{B}_k \mathbf{B}_k^H) \leq p_k$  is satisfied. Otherwise, we turn to the other case  $\mu_k > 0$ , where we have  $\mathbf{B}_k = (\mathbf{F}_k^H \mathbf{F}_k + \mu_k \mathbf{I}_{N_{s,k}})^{-1} \mathbf{F}_k^H \mathbf{U}_k$  and it should be substituted into  $\text{tr}(\mathbf{B}_k \mathbf{B}_k^H) = p_k$  for determining  $\mu_k$ . Given that  $\text{tr}(\mathbf{B}_k \mathbf{B}_k^H)$  is a decreasing function with respect to  $\mu_k$ , we can obtain the optimal  $\mu_k$  via the bisection method [18].

So far, one iteration of our nonlinear receiver-based algorithm has been completed. Hereafter, such iterative algorithm is called “the RIS-DFE algorithm” for continuous RIS phase shifts or “the RIS-b-DFE algorithm” for  $b$  bits of discrete RIS phase shifts, whose procedures are listed in Table I with the output of the  $n$ th iteration marked by superscript  $\langle n \rangle$ . When the linear MMSE receiver is used, so long as we replace  $\mathbf{W}$  with  $\bar{\mathbf{W}}$ , which is also initialized by  $(\mathbf{1})_{N_d \times N}$ , and set  $\mathbf{U} = \mathbf{I}_N$ , those optimization processes of  $\{\theta_m\}$  and  $\{\mathbf{B}_k\}$  remain unaffected. For optimizing  $\bar{\mathbf{W}}$  with fixed  $\{\theta_m\}$  and  $\{\mathbf{B}_k\}$ , by referring to Theorem 4 in [11], the optimal  $\bar{\mathbf{W}}$  is derived as  $\bar{\mathbf{W}} = (\mathbf{A} \mathbf{A}^H + \sigma^2 \mathbf{I}_{N_d})^{-1} \mathbf{A}$ . Such linear receiver-based iterative algorithm is hereafter called “the RIS(-b)-L algorithm”.

When the primal-dual path-following interior-point method is employed to solve the SDP problem (17)–(18) [13], via utilizing the analytical approach in [11], we are able to obtain the computational complexity order (CCO) for one iteration of the RIS-L/DFE algorithm as  $\mathcal{O}(M^{4.5} + M^4 N_m^2 + N_m^3 + TM^3)$ , where  $N_m \triangleq \max\{N, N_s, N_d\}$ , and the CCO for one iteration of the RIS-b-L/DFE algorithm as  $\mathcal{O}(M^{4.5} + M^4 N_m^2 + N_m^3 + TM^3 + 2^b TM)$ . From simulation tests, after each iteration of the algorithms proposed above, the value of the objective,  $E_s$  or  $\bar{E}_s$ , monotonically decreases, which is also bounded below

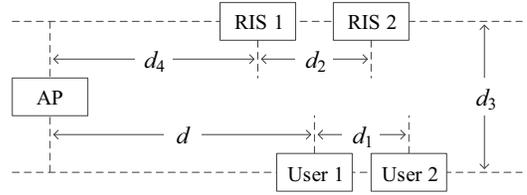


Fig. 2. System architecture configuration for numerical simulations.

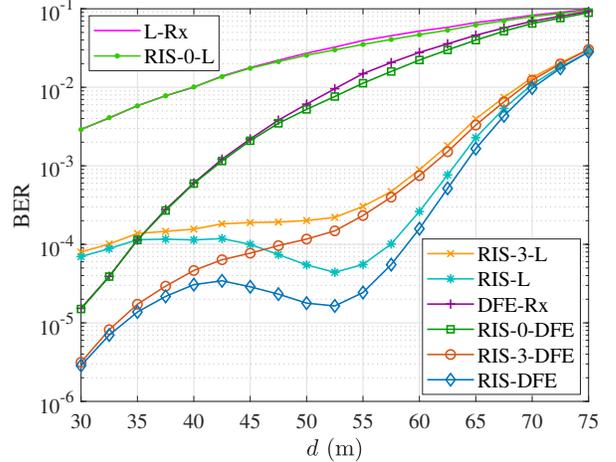


Fig. 3. BER versus  $d$  performance comparisons with  $\tilde{M} = 80$ .

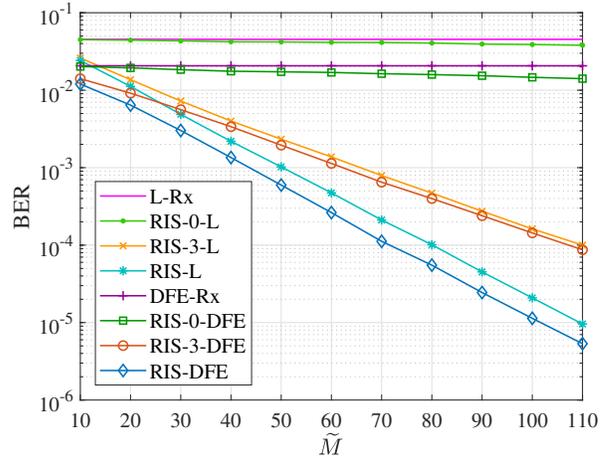


Fig. 4. BER versus  $\tilde{M}$  performance comparisons with  $d = 57.5$  m.

by 0, so the convergence of our algorithms is guaranteed due to Proposition A.3 in [12]. Besides, it is recommended to execute the iterations of our algorithms for 20 times, since after that amount, the performance gains are hardly visible.

#### IV. NUMERICAL SIMULATIONS

This section verifies the bit-error-rate (BER) performance of the proposed algorithms via Monte Carlo simulations. Here, we make comparisons among (a) the RIS(-b)-DFE and RIS(-b)-L algorithms, (b) the nonlinear and linear receiver-based algorithms that do not optimize the RIS phase shifts and fix  $\theta_m = 0$  for  $m = 1, \dots, M$ , called “the RIS-0-DFE algorithm” and “the RIS-0-L algorithm”, and (c) the nonlinear and linear receiver-based algorithms without the assistance of RISs, i.e., fixing  $\boldsymbol{\Theta} = (\mathbf{0})_{M \times M}$ , called “the DFE-Rx algorithm” and “the L-Rx algorithm”, where the algorithms (b) and (c) follow the

same way to optimize  $\{\mathbf{B}_k\}$  and  $\{\mathbf{W}, \mathbf{U}\}$  or  $\bar{\mathbf{W}}$  as in Section III and also execute their iterations for 20 times to ensure the fairness of comparisons. All the numerical simulation results are acquired through the average of 1000 independent channel realizations, and for each of them, half a million information bits per data stream are transmitted from every user with the QPSK modulation mode. Besides, in the MMSE-DFE receiver, the symbols fed back are the practical ones affected by the phenomenon of error propagation, which are regenerated from the information bits detected before with detection errors.

Similar to [3]–[5], the system architecture configuration for our simulations is illustrated by Fig. 2 with  $K = R = 2$ , where there are two users located  $d_1 = 2.5$  m apart on a horizontal line, which is parallel to the line that has two RISs located  $d_2 = 5$  m apart, and the AP is equidistant from these two lines which are vertically  $d_3 = 8$  m apart. In addition, the horizontal distance between the AP and the first RIS is  $d_4 = 55$  m, and that between the AP and the first user is denoted by  $d$ . Like in [1], [7], [8], we set the small-scale fading of the channels  $\mathbf{H}_{s,d,k}$ ,  $\mathbf{H}_{s,r,k}$ , and  $\mathbf{H}_{d,r}$  for  $k = 1, \dots, K$  and  $r = 1, \dots, R$  to follow the Rayleigh fading with their elements independently generated from  $\mathcal{CN}(0, 1)$ . The large-scale fading is modelled by  $\Gamma(\tilde{d}) = Q_0(\tilde{d}/d_0)^{-\tilde{\gamma}}$  where  $Q_0 = -30$  dB stands for the path loss at the reference distance  $d_0 = 1$  m,  $\tilde{d}$  is the distance of signal propagation, and  $\tilde{\gamma}$  represents the path loss exponent (PLE) [4], [5], [9]. Here, due to numerous obstacles, the direct user-AP links are assumed to have severe signal attenuation, so we set the PLEs of  $\mathbf{H}_{s,d,k}$  as  $\gamma_{s,d,k} = \gamma_{sd} = 3.7$ . Meanwhile, the user-RIS links are assumed to be in good condition with the PLEs of  $\mathbf{H}_{s,r,k}$  set as  $\gamma_{s,r,k} = \gamma_s = 2.6$ . Since the RISs are usually placed to ensure high-quality RIS-AP links, we set the PLEs of  $\mathbf{H}_{d,r}$  as  $\gamma_{d,r} = \gamma_d = 2.4$ . Furthermore, we also set  $N_k = \tilde{N} = 4$ ,  $N_{s,k} = \tilde{N}_s = 4$ ,  $p_k = P = 13$  dBm, and  $M_r = \tilde{M}$ , together with  $N_d = 8$ ,  $\sigma^2 = -90$  dBm, and  $b = 3$ .

For all the considered algorithms, Fig. 3 exhibits their BER performance with  $\tilde{M} = 80$  and  $d \in [30 \text{ m}, 75 \text{ m}]$ , additionally, Fig. 4 exhibits that with  $d = 57.5$  m and  $\tilde{M} \in [10, 110]$ . It can be seen that the nonlinear receiver-based algorithms perform better than their linear receiver-based counterparts. Compared to the algorithms without using the RISs, the algorithms that fix  $\{\theta_m\}$  to be zeros can merely provide little gain, while those algorithms that have  $\{\theta_m\}$  optimized can significantly improve the performance, among which, the algorithms with continuous  $\{\theta_m\}$  outperform their counterparts with discrete  $\{\theta_m\}$ . In Fig. 3, as  $d$  increases, the BER performance of all algorithms gets worse in general due to the growth of path loss in direct links. However, when  $d$  ranges from 42.5 m to 52.5 m, the performance of the RIS-L/DFE algorithm becomes better, which is because the users are approaching the RISs that enhance the quality of communications via the optimization of continuous  $\{\theta_m\}$ . In Fig. 4, since the L/DFE-Rx algorithm does not utilize the RISs, its performance is unrelated to  $\tilde{M}$ , while the performance of the other algorithms becomes better as  $\tilde{M}$  increases. When  $\tilde{M} \geq 100$ , compared to the L/DFE-Rx and RIS-0-L/DFE algorithms, the RIS-3-L/DFE algorithm can improve the performance by nearly two orders of magnitude, and such superiority expands to three orders of magnitude for the RIS-L/DFE algorithm.

## V. CONCLUSION

Aiming at enhancing the performance of multiuser uplink MIMO communications via multiple RISs, this letter developed iterative algorithms to jointly optimize the continuous or discrete RIS phase shifts, the MMSE or MMSE-DFE receiving matrices, and the source precoding matrices. Simulations confirmed the effectiveness of our proposed algorithms.

## REFERENCES

- [1] M. A. ElMossallamy, H. Zhang, L. Song, K. G. Seddik, Z. Han, and G. Y. Li, "Reconfigurable intelligent surfaces for wireless communications: Principles, challenges, and opportunities," *IEEE Trans. Cognitive Commun. Netw.*, vol. 6, no. 3, pp. 990–1002, Sep. 2020.
- [2] Q. Wu and R. Zhang, "Towards smart and reconfigurable environment: Intelligent reflecting surface aided wireless network," *IEEE Commun. Mag.*, vol. 58, no. 1, pp. 106–112, Jan. 2020.
- [3] Q. Wu and R. Zhang, "Intelligent reflecting surface enhanced wireless network: Joint active and passive beamforming design," in *Proc. 2018 IEEE Global Commun. Conf. (GLOBECOM)*, Abu Dhabi, United Arab Emirates, Dec. 9–13, 2018, pp. 1–6.
- [4] Q. Wu and R. Zhang, "Intelligent reflecting surface enhanced wireless network via joint active and passive beamforming," *IEEE Trans. Wireless Commun.*, vol. 18, no. 11, pp. 5394–5409, Nov. 2019.
- [5] Q. Wu and R. Zhang, "Beamforming optimization for wireless network aided by intelligent reflecting surface with discrete phase shifts," *IEEE Trans. Commun.*, vol. 68, no. 3, pp. 1838–1851, Mar. 2020.
- [6] H. Guo, Y.-C. Liang, J. Chen, and E. G. Larsson, "Weighted sum-rate maximization for reconfigurable intelligent surface aided wireless networks," *IEEE Trans. Wireless Commun.*, vol. 19, no. 5, pp. 3064–3076, May 2020.
- [7] C. Huang, A. Zappone, G. C. Alexandropoulos, M. Debbah, and C. Yuen, "Reconfigurable intelligent surfaces for energy efficiency in wireless communication," *IEEE Trans. Wireless Commun.*, vol. 18, no. 8, pp. 4157–4170, Aug. 2019.
- [8] W. Yan, X. Yuan, and X. Kuai, "Passive beamforming and information transfer via large intelligent surface," *IEEE Wireless Commun. Lett.*, vol. 9, no. 4, pp. 533–537, Apr. 2020.
- [9] W. Yan, X. Yuan, Z.-Q. He, and X. Kuai, "Passive beamforming and information transfer design for reconfigurable intelligent surfaces aided multiuser MIMO systems," *IEEE J. Sel. Areas Commun.*, vol. 38, no. 8, pp. 1793–1808, Aug. 2020.
- [10] X. Yu, D. Xu, Y. Sun, D. W. K. Ng, and R. Schober, "Robust and secure wireless communications via intelligent reflecting surfaces," *IEEE J. Sel. Areas Commun.*, vol. 38, no. 11, pp. 2637–2652, Nov. 2020.
- [11] Y. Lv, Z. He, and Y. Rong, "Two-way AF MIMO multi-relay system design using MMSE-DFE techniques," *IEEE Trans. Wireless Commun.*, vol. 20, no. 1, pp. 389–405, Jan. 2021.
- [12] D. P. Bertsekas, *Nonlinear Programming*, 2nd ed. Belmont, Massachusetts: Athena Scientific, 1999.
- [13] Z.-Q. Luo, W.-K. Ma, A. M.-C. So, Y. Ye, and S. Zhang, "Semidefinite relaxation of quadratic optimization problems," *IEEE Signal Process. Mag.*, vol. 27, no. 3, pp. 20–34, May 2010.
- [14] R. A. Horn and C. R. Johnson, *Matrix Analysis*, 2nd ed. New York, USA: Cambridge Univ. Press, 2013.
- [15] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, UK: Cambridge Univ. Press, 2004, sec. 5.5.3, pp. 243–246.
- [16] R. A. Horn and C. R. Johnson, *Topics in Matrix Analysis*. Cambridge, UK: Cambridge Univ. Press, 1991.
- [17] M. Grant and S. Boyd, "CVX: Matlab software for disciplined convex programming, version 2.2," CVX Research, Inc., Jan. 2020. [Online]. Available: <http://cvxr.com/cvx>
- [18] W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery, *Numerical Recipes: The Art of Scientific Computing*, 3rd ed. Cambridge, UK: Cambridge Univ. Press, 2007.