

Tensor-Based Channel Estimation and Reflection Design for RIS-Aided Millimeter-Wave MIMO Communication Systems

Sepideh Gherekhloo, Khaled Ardah, André L. F. de Almeida, and Martin Haardt

Abstract—In this work, we consider both channel estimation and reflection coefficient design problems in point-to-point reconfigurable intelligent surface (RIS)-aided millimeter-wave (mmWave) MIMO communication systems. First, we show that by exploiting the low-rank nature of mmWave MIMO channels, the received training signals can be written as a low-rank multi-way tensor admitting a canonical polyadic (CP) decomposition. Utilizing such a structure, a tensor-based RIS channel estimation method (termed TenRICE) is proposed, wherein the tensor factor matrices are estimated using an alternating least squares method. Using TenRICE, the transmitter-to-RIS and the RIS-to-receiver channels are efficiently and separately estimated, up to a trivial scaling factor. After that, we formulate the beamforming and RIS reflection coefficient design as a spectral efficiency maximization task. Due to its non-convexity, we propose a heuristic non-iterative two-step method, where the RIS reflection vector is obtained in a closed form using a Frobenius-norm maximization (FroMax) strategy. Our numerical results show that TenRICE has a superior performance, compared to benchmark methods, approaching the Cramér–Rao lower bound with a low training overhead. Moreover, we show that FroMax achieves a comparable performance to benchmark methods with a lower complexity.

Index Terms—Reconfigurable intelligent surface, channel estimation, RIS reflection design, CP tensor decomposition.

I. INTRODUCTION

Reconfigurable intelligent surfaces (RISs) have been proposed recently as a cost-effective technology for reconfiguring the propagation channels in wireless communication systems [1]. An RIS is a 2D surface equipped with a large number of tunable units that can be realized using, e.g., inexpensive antennas or metamaterials and controlled in real-time to influence the communication channels without generating its own signals. Among its many applications, an RIS can be utilized as a solution to the signal-blockage problem in millimeter-wave (mmWave)-based communications by providing alternative and tunable RIS-aided channels.

Recently, RIS-aided communications have attracted great attention, due to their potential of improving the efficiency of wireless mobile communications. RIS reflection design, in particular, have been extensively investigated under various setups

The authors gratefully acknowledge the support of the German Research Foundation (DFG) under contracts no. HA 2239/14-1 and no. HA 2239/6-2 and the support of CAPES/PRINT (Grant no. 88887.311965/2018-00). The research of André L. F. de Almeida is partially supported by the CNPq (Grant no. 306616/2016-5).

S. Gherekhloo, K. Ardah, and M. Haardt are with the Communications Research Laboratory (CRL), TU Ilmenau, Ilmenau, Germany (e-mail: {khaled.ardah,sepideh.gherekhloo, martin.haardt}@tu-ilmenau.de). A. L. F. de Almeida is with the Wireless Telecom Research Group (GTEL), Federal University of Ceará, Fortaleza, Brazil (e-mail: andre@gtel.ufc.br).

and objectives, see [2]–[5] and reference therein. However, due to the non-convexity of the involved problems, relaxations and alternating optimization techniques are commonly used to obtain a locally optimal solution. For example, the authors in [2] considered the capacity maximization and proposed an alternating optimization approach to find a locally optimal solution by iteratively optimizing the transmit covariance matrix or one of the RIS reflection coefficients with the others being fixed. However, such an alternating approach increases the computational complexity and becomes a limiting factor in practice, especially in a massive RIS setup.

The vast majority of the existing works assume perfect channel state information (CSI) at the transceivers, see [2]–[5], which can never be obtained in practice. Recently, RIS-aided channel estimation (CE) methods have been proposed, e.g., in [6]–[9]. These works, however, require that the number of training subframes is, at least, equal to the number of RIS reflection units to obtain an accurate CSI estimate, which increases the training overhead and complexity. To overcome these issues, several approaches have been studied, e.g., by exploiting the low-rank nature of mmWave channels and the multidimensional (i.e., tensor) structure of the received signals. The former allows the CE to be formulated as a sparse-recovery problem and solved using compressed sensing (CS) tools [10]–[12], which are known to require a few measurements to have an accurate estimate, see [13]–[15]. In [13], by exploiting the low-rank nature of the mmWave channels, we have proposed the TRICE framework, which formulates the CE in RIS-aided mmWave MIMO systems as a two-stage multidimensional sparse-recovery problem. On the other hand, tensor-based signal modeling and processing methods offer fundamental advantages over their bilinear (matrix) counterparts, since they have the ability to improve the identifiability of the parameters due to the powerful uniqueness properties of tensor decompositions [16]. In [17], it is shown that the received signals in RIS-aided MIMO communication systems can be written as a 3-way tensor admitting a canonical polyadic (CP) decomposition. However, the proposed method in [17] assumes sub-6 GHz systems and, thus, requires a large number of training subframes, similarly to [6]–[9].

In this paper, we extend our TRICE framework in [13] and propose a CP **Tensor** decomposition method for **RIS**-aided **CE** in mmWave MIMO systems, termed TenRICE, by jointly exploiting the tensor structure of the received signals and the low-rank nature of mmWave channels. Using

the TenRICE method, the transmitter-to-RIS and the RIS-to-receiver channels can be estimated separately, up to a trivial scaling factor. After that, we formulate the beamforming and the RIS reflection coefficient design as a spectral efficiency (SE) maximization problem. Due to its non-convexity, we propose a heuristic non-iterative two-step solution, where the RIS reflection vector is obtained, in contrast to [2], in a closed form using a **Frobenius-norm Maximization (FroMax)** strategy. Our numerical results show that TenRICE has a superior performance, compared to the TRICE framework, approaching the Cramér–Rao bound (CRB). Moreover, we show that FroMax achieves a comparable performance to benchmark methods with a lower complexity.

II. SYSTEM MODEL

In this paper¹, we consider an RIS-aided mmWave MIMO communication system as depicted in Fig. 1, where a transmitter (TX) with M_T antennas is communicating with a receiver (RX) with M_R antennas via an RIS-aided MIMO channel. The direct channel between the TX and the RX is assumed unavailable or too weak, e.g., due to blockage. The RIS has M_S inexpensive reflecting elements arranged uniformly with half-wavelength inter-element spacing on a rectangular surface with M_S^v vertical and M_S^h horizontal elements such that $M_S = M_S^v \cdot M_S^h$.

Let $\mathbf{H}_T \in \mathbb{C}^{M_S \times M_T}$ be the TX to RIS channel and $\mathbf{H}_R \in \mathbb{C}^{M_R \times M_S}$ be the RIS to RX channel with $\mathbb{E}\{\|\mathbf{H}_T\|_F^2\} = M_S M_T$ and $\mathbb{E}\{\|\mathbf{H}_R\|_F^2\} = M_S M_R$. We assume a block-fading channel scenario, where \mathbf{H}_T and \mathbf{H}_R remain constant during every channel coherence block and change from block to block. We assume that every block is divided into two sub-blocks: one for CE and another for data transmission (DT), see Fig. 2.

In the CE phase, we conduct a channel training procedure that occupies $K = K_T \cdot K_S$ subframes. The received signal at the RX at the (s, t) th subframe is given as

$$\mathbf{y}_{s,t} = \mathbf{W}^H \mathbf{H}_R \text{diag}\{\phi_s\} \mathbf{H}_T \tilde{\mathbf{f}}_t s_t + \mathbf{W}^H \mathbf{z}_{s,t} \in \mathbb{C}^{K_R}, \quad (1)$$

where $\mathbf{W} \in \mathbb{C}^{M_R \times K_R}$ is a fixed training decoding matrix with K_R beams, $\tilde{\mathbf{f}}_t \in \mathbb{C}^{M_T}$ is the t th training vector of the TX with $\|\tilde{\mathbf{f}}_t\|_2^2 = 1$, $t \in \{1, \dots, K_T\}$, $\phi_s \in \mathbb{C}^{M_S}$ is the s th training vector of the RIS with $|\phi_s[m]| = \frac{1}{\sqrt{M_S}}, \forall m$, $s \in \{1, \dots, K_S\}$, $s_t \in \mathbb{C}$ is the unit-norm pilot symbol, and $\mathbf{z}_{s,t} \in \mathbb{C}^{M_R}$ is the additive white Gaussian noise vector

¹**Notation:** The transpose, the conjugate transpose (Hermitian), the Moore-Penrose pseudoinverse, the Kronecker product, and the Khatri-Rao product are denoted as \mathbf{A}^T , \mathbf{A}^H , \mathbf{A}^+ , \otimes , and \diamond , respectively. Moreover, $\mathbf{1}_N$ is the all ones vector of length N , \mathbf{I}_N is the $N \times N$ identity matrix, $\text{diag}\{\mathbf{a}\}$ forms a diagonal matrix \mathbf{A} by putting the entries of the input vector \mathbf{a} in its main diagonal, $\text{undia}\{\mathbf{A}\}$ is the reverse of the diag operator, $\text{vec}\{\mathbf{A}\}$ forms a vector by stacking the columns of \mathbf{A} over each other, and the n -mode product of a tensor $\mathcal{A} \in \mathbb{C}^{I_1 \times I_2 \times \dots \times I_N}$ with a matrix $\mathbf{B} \in \mathbb{C}^{J \times I_n}$ is denoted as $\mathcal{A} \times_n \mathbf{B}$. Throughout this paper, we assume that the singular values of a given diagonal singular matrix are arranged in a decreasing order. Moreover, the following properties are used: *Property 1:* $\text{vec}\{\mathbf{ABC}\} = (\mathbf{C}^T \otimes \mathbf{A}) \text{vec}\{\mathbf{B}\}$. *Property 2:* $\mathbf{AB} \diamond \mathbf{CD} = (\mathbf{A} \otimes \mathbf{C})(\mathbf{B} \diamond \mathbf{D})$. *Property 3:* $(\mathbf{A} \otimes \mathbf{C})(\mathbf{B} \otimes \mathbf{D}) = \mathbf{AB} \otimes \mathbf{CD}$. *Property 4:* Let $\mathbf{A}_1 \in \mathbb{C}^{J_1 \times L_1}$ and $\mathbf{A}_2 \in \mathbb{C}^{J_2 \times L_2}$. Then $\mathbf{A}_1 \otimes \mathbf{A}_2 = \mathbf{A}_1 \Omega_1 \diamond \mathbf{A}_2 \Omega_2$, where $\Omega_1 = \mathbf{I}_{L_1} \otimes \mathbf{1}_{L_2}^T$ and $\Omega_2 = \mathbf{1}_{L_1}^T \otimes \mathbf{I}_{L_2}$ so that $\Omega_1 \diamond \Omega_2 = \mathbf{I}_{L_1 L_2}$. *Property 5:* $\text{vec}\{\mathbf{A} \text{diag}\{\mathbf{b}\} \mathbf{C}\} = (\mathbf{C}^T \diamond \mathbf{A}) \mathbf{b}$.

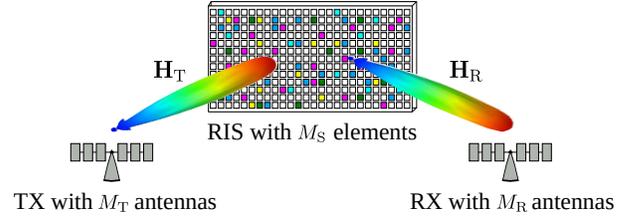


Fig. 1. An RIS-aided MIMO mmWave communication system.

Phase 1: Channel Estimation (CE) No. of subframes: $K = K_T \cdot K_S$				Phase 2: Data Transmission (DT) No. of data-streams: N_s			
$\hat{\mathbf{f}}_1$	\dots	$\hat{\mathbf{f}}_1$	\dots	$\hat{\mathbf{f}}_{K_T}$	\dots	$\hat{\mathbf{f}}_{K_T}$	\dots
ϕ_1	\dots	ϕ_{K_S}	\dots	ϕ_1	\dots	ϕ_{K_S}	\dots

Fig. 2. One channel coherence block.

having zero-mean circularly symmetric complex-valued entries with variance σ^2 . We stack $\{\mathbf{y}_{s,t}\}_{t=1}^{K_T}$ on top of each other as $\mathbf{y}_s = [\mathbf{y}_{s,1}^T, \dots, \mathbf{y}_{s,K_T}^T]^T$ and after that we stack $\{\mathbf{y}_s\}_{s=1}^{K_S}$ next to each other as $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_{K_S}]$. Then, using Properties 2 and 5, the above measurement matrix \mathbf{Y} can be written as

$$\mathbf{Y} = (\mathbf{F}^T \otimes \mathbf{W}^H) \mathbf{H}_c \Phi + \mathbf{Z} \in \mathbb{C}^{K_R K_T \times K_S}, \quad (2)$$

where $\mathbf{H}_c = \mathbf{H}_T^T \diamond \mathbf{H}_R$ represents the cascaded channel matrix, $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_{K_S}]$, $\mathbf{z}_s = [(\mathbf{W}^H \mathbf{z}_{s,1})^T, \dots, (\mathbf{W}^H \mathbf{z}_{s,K_T})^T]^T$, $\mathbf{F} = [\mathbf{f}_1 s_1, \dots, \mathbf{f}_{K_T} s_{K_T}]$, and $\Phi = [\phi_1, \dots, \phi_{K_S}]$. Given the measurement matrix \mathbf{Y} , the main goal of Section III is to obtain an accurate estimate of \mathbf{H}_T and \mathbf{H}_R , while keeping the number of training subframes K as small as possible.

In the DT phase, given the estimated channels $\hat{\mathbf{H}}_R$ and $\hat{\mathbf{H}}_T$, the TX first designs the precoding matrix $\mathbf{P} \in \mathbb{C}^{M_T \times N_s}$, the decoding matrix $\mathbf{Q} \in \mathbb{C}^{M_R \times N_s}$, and the RIS reflection coefficient vector $\omega \in \mathbb{C}^{M_S}$ with $|\omega[m]| = \frac{1}{\sqrt{M_S}}, \forall m$, to transmit the vector $\mathbf{s} \in \mathbb{C}^{N_s}$ of N_s data streams with $\mathbb{E}[\mathbf{s} \mathbf{s}^H] = \mathbf{I}_{N_s}$ to the RX. Therefore, the received signal vector at the RX is given as

$$\mathbf{y} = \mathbf{Q}^H \mathbf{H}_c \mathbf{P} \mathbf{s} + \mathbf{Q}^H \mathbf{z} \in \mathbb{C}^{N_s}, \quad (3)$$

where $\mathbf{H}_e = \mathbf{H}_R \text{diag}\{\omega\} \mathbf{H}_T$ is the effective channel matrix. The system SE is given as

$$\text{SE} = \log_2 \det(\mathbf{I}_{N_s} + \mathbf{R}^{-1} \mathbf{Q}^H \mathbf{H}_c \mathbf{P} \mathbf{P}^H \mathbf{H}_e^H \mathbf{Q}), \quad (4)$$

where $\mathbf{R} = \sigma^2 \mathbf{Q}^H \mathbf{Q}$ is the noise covariance matrix. In Section IV, we propose a non-iterative beamforming and RIS reflection coefficient design method to maximize the SE, where the RIS reflection vector is obtained in a closed form using a FroMax strategy.

Channel model: In mmWave-based communications [18], it was observed that the number of paths L_T and L_R for \mathbf{H}_T and \mathbf{H}_R respectively, are very small compared to the number of antenna elements. This implies that $\text{rank}\{\mathbf{H}_T\} \leq L_T$ and $\text{rank}\{\mathbf{H}_R\} \leq L_R$. Therefore, similarly to [13], by assuming that the TX and the RX employ uniform linear arrays (ULAs)²,

²The extension of the proposed methods to scenarios where the TX and/or the RX are equipped with uniform rectangular arrays (URAs) is straightforward.

\mathbf{H}_T and \mathbf{H}_R follow the geometric channel model, which can be written as

$$\mathbf{H}_T = \frac{1}{\sqrt{L_T}} \sum_{\ell=1}^{L_T} g_{T,\ell} \mathbf{v}_{2D}(\mu_{T,\ell}^v, \mu_{T,\ell}^h) \mathbf{v}_{1D}(\psi_{T,\ell})^T = \mathbf{B}_T \mathbf{G}_T \mathbf{A}_T^T, \quad (5)$$

$$\mathbf{H}_R = \frac{1}{\sqrt{L_R}} \sum_{\ell=1}^{L_R} g_{R,\ell} \mathbf{v}_{1D}(\psi_{R,\ell}) \mathbf{v}_{2D}(\mu_{R,\ell}^v, \mu_{R,\ell}^h)^T = \mathbf{A}_R \mathbf{G}_R \mathbf{B}_R^T,$$

where $g_{X,\ell} \sim \mathcal{CN}(0, 1)$ is the ℓ th path gain, $\psi_{T,\ell} \in [0, 2\pi]$ is the ℓ th direction-of-departure (DoD) spatial frequency from the TX, $\psi_{R,\ell} \in [0, 2\pi]$ is the ℓ th direction-of-arrival (DoA) spatial frequency at the RX, $\mu_{T,\ell}^h \in [0, 2\pi]$ and $\mu_{T,\ell}^v \in [0, \pi]$ are the ℓ th horizontal and vertical DoA spatial frequencies at the RIS, while $\mu_{R,\ell}^h \in [0, 2\pi]$ and $\mu_{R,\ell}^v \in [0, \pi]$ are the ℓ th horizontal and vertical DoD spatial frequencies from the RIS. In (5), the 1D and the 2D array steering vectors are given as $\mathbf{v}_{1D}(\nu) = [1, e^{j\nu}, \dots, e^{j(M-1)\nu}]^T \in \mathbb{C}^M$ and $\mathbf{v}_{2D}(\nu^v, \nu^h) = \mathbf{v}_{1D}(\nu^v) \diamond \mathbf{v}_{1D}(\nu^h)$, respectively, where $\mathbf{v}_{1D}(\nu^v) \in \mathbb{C}^{M^v}$ and $\mathbf{v}_{1D}(\nu^h) \in \mathbb{C}^{M^h}$. Moreover, \mathbf{H}_T and \mathbf{H}_R are written in a compact form by letting $\mathbf{A}_X = [\mathbf{v}_{1D}(\psi_{X,1}), \dots, \mathbf{v}_{1D}(\psi_{X,L_X})] \in \mathbb{C}^{M_X \times L_X}$, $\mathbf{B}_X = \mathbf{B}_X^v \diamond \mathbf{B}_X^h$, $\mathbf{B}_X^v = [\mathbf{v}_{1D}(\mu_{X,1}^v), \dots, \mathbf{v}_{1D}(\mu_{X,L_X}^v)] \in \mathbb{C}^{M_X^v \times L_X}$, and $\mathbf{G}_X = \frac{1}{\sqrt{L_X}} \text{diag}\{g_{X,1}, \dots, g_{X,L_X}\}$ for $X \in \{T, R\}$, $Y \in \{v, h\}$.

III. PHASE 1: THE PROPOSED CE METHOD (TENRICE)

In this section, we propose our **Tensor-based RIS-aided CE** (TenRICE) algorithm by jointly exploiting the low-rank nature of mmWave channels and the tensor structure of received signals. By utilizing the channels model in (5), the cascaded channel matrix $\mathbf{H}_c = \mathbf{H}_T^T \diamond \mathbf{H}_R$ in (2) can be written as

$$\mathbf{H}_c = (\mathbf{A}_T \mathbf{G}_T \mathbf{B}_T^T \diamond \mathbf{A}_R \mathbf{G}_R \mathbf{B}_R^T) \stackrel{(a)}{=} (\mathbf{A}_T \otimes \mathbf{A}_R) \mathbf{G} \mathbf{B}, \quad (6)$$

where $\mathbf{G} = \mathbf{G}_T \otimes \mathbf{G}_R \in \mathbb{C}^{L \times L}$, $\mathbf{B} = \mathbf{B}_T^T \diamond \mathbf{B}_R^T \in \mathbb{C}^{L \times M_s}$, $L = L_R \cdot L_T$, and $\stackrel{(a)}{=}$ is obtained from Property 2. In [13], we have shown that \mathbf{B} can be expressed as $\mathbf{B} = (\mathbf{B}_v \diamond \mathbf{B}_h)^T$, where $\mathbf{B}_v = [\mathbf{v}_{1D}(\mu_1^v), \dots, \mathbf{v}_{1D}(\mu_L^v)] \in \mathbb{C}^{M_s^v \times L}$, $\mathbf{B}_h = [\mathbf{v}_{1D}(\mu_1^h), \dots, \mathbf{v}_{1D}(\mu_L^h)] \in \mathbb{C}^{M_s^h \times L}$, $\mu_n^v = \mu_{T,\ell}^v + \mu_{R,k}^v$, $\mu_n^h = \mu_{T,\ell}^h + \mu_{R,k}^h$, $\ell \in \{1, \dots, L_T\}$, $k \in \{1, \dots, L_R\}$, and $n = (\ell - 1) \cdot L_R + k \in \{1, \dots, L\}$. Then, using Property 2, (6) can be rewritten as

$$\mathbf{H}_c = (\mathbf{A}_T \otimes \mathbf{A}_R) \mathbf{G} (\mathbf{B}_v \diamond \mathbf{B}_h)^T, \quad (7)$$

which is characterized by the following spatial frequency vectors: $\boldsymbol{\psi}_R = [\psi_{R,1}, \dots, \psi_{R,L_R}]^T$, $\boldsymbol{\psi}_T = [\psi_{T,1}, \dots, \psi_{T,L_T}]^T$, $\boldsymbol{\mu}^h = [\mu_1^h, \dots, \mu_L^h]^T$, and $\boldsymbol{\mu}^v = [\mu_1^v, \dots, \mu_L^v]^T$ that define \mathbf{A}_R , \mathbf{A}_T , \mathbf{B}_h , and \mathbf{B}_v , respectively. Therefore, to obtain an estimate of \mathbf{H}_c , it is sufficient to obtain an estimate of the above vectors from the measurement matrix \mathbf{Y} in (2), including the path gains vector $\mathbf{g} = \text{undia}\{\mathbf{G}\}$. In [13], we have proposed a two-stage framework, termed TRICE, which estimates $\boldsymbol{\psi}_R$ and $\boldsymbol{\psi}_T$ in the first stage as well as $\boldsymbol{\mu}^h$, $\boldsymbol{\mu}^v$, and \mathbf{g} in the second stage using any efficient multidimensional sparse-recovery technique, like CS [12] and ESPRIT [19]. To further improve the performance of the TRICE framework, we propose in the following the TenRICE method by exploiting the tensor structure of the measurement matrix \mathbf{Y} .

We assume that the RIS reflection coefficient matrix during the training phase has a Kronecker structure given as $\Phi = \Phi_v \otimes \Phi_h$, where $\Phi_v \in \mathbb{C}^{M_s^v \times K_s^v}$, $\Phi_h \in \mathbb{C}^{M_s^h \times K_s^h}$, and $K_s^v \cdot K_s^h = K_s$. By substituting (6) into (2), the vectorized form of \mathbf{Y} , i.e., $\mathbf{y} = \text{vec}\{\mathbf{Y}\}$ can be written as

$$\begin{aligned} \mathbf{y} &\stackrel{(a)}{=} \text{vec}\{(\mathbf{F}^T \mathbf{A}_T \otimes \mathbf{W}^H \mathbf{A}_R) \mathbf{G} (\mathbf{B}_v \diamond \mathbf{B}_h)^T \Phi\} + \mathbf{z} \\ &\stackrel{(b)}{=} \text{vec}\{(\mathbf{F}^T \mathbf{A}_T \Omega_T \otimes \mathbf{W}^H \mathbf{A}_R \Omega_R) \mathbf{G} (\mathbf{B}_v \diamond \mathbf{B}_h)^T \Phi\} + \mathbf{z} \\ &\stackrel{(c)}{=} (\Phi_v^T \mathbf{B}_v \diamond \Phi_h^T \mathbf{B}_h \diamond \mathbf{F}^T \mathbf{A}_T \Omega_T \otimes \mathbf{W}^H \mathbf{A}_R \Omega_R) \mathbf{g} + \mathbf{z}, \quad (8) \end{aligned}$$

where $\mathbf{z} = \text{vec}\{\mathbf{Z}\}$ and $\mathbf{g} = \text{undia}\{\mathbf{G}\}$. Moreover, $\stackrel{(a)}{=}$, $\stackrel{(b)}{=}$, and $\stackrel{(c)}{=}$ are obtained by applying Properties 1, 2, and 4, where $\Omega_T \stackrel{\text{def}}{=} \mathbf{I}_{L_T} \otimes \mathbf{1}_{L_R}^T$ and $\Omega_R \stackrel{\text{def}}{=} \mathbf{1}_{L_T}^T \otimes \mathbf{I}_{L_R}$. From (8), we observe that \mathbf{y} is the vectorized form of the transposed 4-mode unfolding of a 4-way tensor $\mathcal{Y} \in \mathbb{C}^{K_R \times K_T \times K_s^v \times K_s^h}$, i.e., $\mathbf{y} = [\mathcal{Y}]_{(4)}^T$ that admits a constrained CP decomposition as [16], [20]

$$\mathcal{Y} = \mathcal{I}_{4,L} \times_1 \bar{\mathbf{A}}_R \Omega_R \times_2 \bar{\mathbf{A}}_T \Omega_T \times_3 \bar{\mathbf{B}}_h \times_4 \bar{\mathbf{B}}_v + \mathcal{Z}, \quad (9)$$

where \mathcal{Z} is the noise tensor, $\mathcal{I}_{4,L} \in \mathbb{C}^{L \times L \times L \times L}$ is a super-diagonal tensor with ones on the super diagonal, and

$$\bar{\mathbf{A}}_R = \mathbf{W}^H \mathbf{A}_R = \mathbf{W}^H [\mathbf{v}_{1D}(\psi_{R,1}), \dots, \mathbf{v}_{1D}(\psi_{R,L_R})], \quad (10)$$

$$\bar{\mathbf{A}}_T = \mathbf{F}^T \mathbf{A}_T = \mathbf{F}^T [\mathbf{v}_{1D}(\psi_{T,1}), \dots, \mathbf{v}_{1D}(\psi_{T,L_T})], \quad (11)$$

$$\bar{\mathbf{B}}_h = \Phi_h^T \mathbf{B}_h = \Phi_h^T [\mathbf{v}_{1D}(\mu_1^h), \dots, \mathbf{v}_{1D}(\mu_L^h)], \quad (12)$$

$$\bar{\mathbf{B}}_v = \Phi_v^T \mathbf{B}_v \mathbf{G} = \Phi_v^T [\mathbf{v}_{1D}(\mu_1^v), \dots, \mathbf{v}_{1D}(\mu_L^v)] \mathbf{G}. \quad (13)$$

The n -mode unfoldings of tensor \mathcal{Y} , for $n \in \{1, 2, 3, 4\}$ can be expressed as

$$[\mathcal{Y}]_{(1)} = \bar{\mathbf{A}}_R \Omega_R (\bar{\mathbf{B}}_v \diamond \bar{\mathbf{B}}_h \diamond \bar{\mathbf{A}}_T \Omega_T)^T + [\mathcal{Z}]_{(1)} \quad (14)$$

$$[\mathcal{Y}]_{(2)} = \bar{\mathbf{A}}_T \Omega_T (\bar{\mathbf{B}}_v \diamond \bar{\mathbf{B}}_h \diamond \bar{\mathbf{A}}_R \Omega_R)^T + [\mathcal{Z}]_{(2)} \quad (15)$$

$$[\mathcal{Y}]_{(3)} = \bar{\mathbf{B}}_h (\bar{\mathbf{B}}_v \diamond \bar{\mathbf{A}}_T \Omega_T \diamond \bar{\mathbf{A}}_R \Omega_R)^T + [\mathcal{Z}]_{(3)} \quad (16)$$

$$[\mathcal{Y}]_{(4)} = \bar{\mathbf{B}}_v (\bar{\mathbf{B}}_h \diamond \bar{\mathbf{A}}_T \Omega_T \diamond \bar{\mathbf{A}}_R \Omega_R)^T + [\mathcal{Z}]_{(4)}. \quad (17)$$

Given the measurement tensor \mathcal{Y} , the CE task boils down to first estimating the tensor factor matrices. Several techniques have been proposed to achieve this end, e.g., in [21]–[23]. One of these techniques is the alternating least squares (ALS) [24], which minimizes the data fitting error with respect to one of the factor matrices, with the other three being fixed. For example, to estimate $\bar{\mathbf{A}}_R$, assuming that $\bar{\mathbf{A}}_T$, $\bar{\mathbf{B}}_h$, and $\bar{\mathbf{B}}_v$ are fixed, the problem can be formulated as

$$\hat{\bar{\mathbf{A}}}_R = \arg \min_{\bar{\mathbf{A}}_R} \left\| [\mathcal{Y}]_{(1)} - \bar{\mathbf{A}}_R \Omega_R (\bar{\mathbf{B}}_v \diamond \bar{\mathbf{B}}_h \diamond \bar{\mathbf{A}}_T \Omega_T)^T \right\|_F^2, \quad (18)$$

which is a convex problem and can be solved using the LS method. Using the same methodology, $\hat{\bar{\mathbf{A}}}_T$, $\hat{\bar{\mathbf{B}}}_h$, and $\hat{\bar{\mathbf{B}}}_v$ can be estimated similarly to (18). Therefore, an ALS-based method can be used to estimate the four factor matrices as summarized in Algorithm 1 (from step 3 to step 9), which is guaranteed to converge monotonically to a local optimum point [24].

Let $\hat{\bar{\mathbf{A}}}_R$, $\hat{\bar{\mathbf{A}}}_T$, $\hat{\bar{\mathbf{B}}}_h$, and $\hat{\bar{\mathbf{B}}}_v$ denote the estimated factor matrices at the convergence of the iterative steps of Algorithm 1. Then, the parameters associated with each factor matrix can be recovered, e.g., via a simple correlation-based scheme. For

Algorithm 1 Tensor-based RIS-aided CE (TenRICE)

- 1: Input: Measurement tensor $\mathcal{Y} \in \mathbb{C}^{K_R \times K_T \times K_S^h \times K_S^v}$ and I_{\max}
 - 2: Output: Estimated channels $\widehat{\mathbf{H}}_T$ and $\widehat{\mathbf{H}}_R$
 - 3: Initialization: $\widehat{\mathbf{B}}_v^{(0)}$, $\widehat{\mathbf{B}}_h^{(0)}$, and $\widehat{\mathbf{A}}_T^{(0)}$, e.g., randomly
 - 4: **while** not converged or $i < I_{\max}$ **do**
 - 5: $\widehat{\mathbf{A}}_R^{(i)} = [\mathcal{Y}]_{(1)} \left[\Omega_R (\widehat{\mathbf{B}}_v^{(i-1)} \diamond \widehat{\mathbf{B}}_h^{(i-1)} \diamond \widehat{\mathbf{A}}_T^{(i-1)} \Omega_T)^\top \right]^+$
 - 6: $\widehat{\mathbf{A}}_T^{(i)} = [\mathcal{Y}]_{(2)} \left[\Omega_T (\widehat{\mathbf{B}}_v^{(i-1)} \diamond \widehat{\mathbf{B}}_h^{(i-1)} \diamond \widehat{\mathbf{A}}_R^{(i)} \Omega_R)^\top \right]^+$
 - 7: $\widehat{\mathbf{B}}_h^{(i)} = [\mathcal{Y}]_{(3)} \left[(\widehat{\mathbf{B}}_v^{(i-1)} \diamond \widehat{\mathbf{A}}_T^{(i)} \Omega_T \diamond \widehat{\mathbf{A}}_R^{(i)} \Omega_R)^\top \right]^+$
 - 8: $\widehat{\mathbf{B}}_v^{(i)} = [\mathcal{Y}]_{(4)} \left[(\widehat{\mathbf{B}}_h^{(i)} \diamond \widehat{\mathbf{A}}_T^{(i)} \Omega_T \diamond \widehat{\mathbf{A}}_R^{(i)} \Omega_R)^\top \right]^+$
 - 9: **end while**
 - 10: Recover $\widehat{\psi}_R$, $\widehat{\psi}_T$, $\widehat{\mu}^h$, $\widehat{\mu}^v$ using, e.g., (19) or NOMP [25]
 - 11: Compute $\widehat{\mathbf{g}} = [\Phi_v^\top \widehat{\mathbf{B}}_v \diamond \Phi_h^\top \widehat{\mathbf{B}}_h \diamond \mathbf{F}^\top \widehat{\mathbf{A}}_T \Omega_T \diamond \mathbf{W}^H \widehat{\mathbf{A}}_R \Omega_R]^\top$
 - 12: Reconstruct $\widehat{\mathbf{H}}_c = (\widehat{\mathbf{A}}_R \otimes \widehat{\mathbf{A}}_T) \text{diag}\{\widehat{\mathbf{g}}\} (\widehat{\mathbf{B}}_v \diamond \widehat{\mathbf{B}}_h)^\top$
 - 13: Estimate $\widehat{\mathbf{H}}_T$ and $\widehat{\mathbf{H}}_R$ from $\widehat{\mathbf{H}}$ using [17, Algorithm 1]
-

example, the k th entry of ψ_R , i.e., $\psi_{R,k}$ associated with the k th column vector of $\widehat{\mathbf{A}}_R$, i.e., $\widehat{\mathbf{a}}_{R,k}$ can be recovered as

$$\widehat{\psi}_{R,k} = \arg \max_{\psi \in [0, 2\pi]} \frac{|\widehat{\mathbf{a}}_{R,k}^H \mathbf{W}^H \mathbf{v}_{\text{ID}}(\psi)|}{\|\widehat{\mathbf{a}}_{R,k}\| \|\mathbf{W}^H \mathbf{v}_{\text{ID}}(\psi)\|}, \quad (19)$$

which can be efficiently implemented by first employing a coarse grid and then gradually refining it around the maximizing grid points. Alternatively, (19) can be interpreted as an off-grid sparse recovery problem, where efficient methods like, Newtonized OMP (NOMP) [25] can be readily applied to recover $\widehat{\psi}_{R,k}$ with high accuracy and low complexity. A similar approach can be used to recover the vectors ψ_T , μ^h , and μ^v from $\widehat{\mathbf{A}}_T$, $\widehat{\mathbf{B}}_h$, and $\widehat{\mathbf{B}}_v$, respectively.

Next, using the estimated vectors $\widehat{\psi}_R$, $\widehat{\psi}_T$, $\widehat{\mu}^h$, and $\widehat{\mu}^v$ in step 10, we reconstruct $\widehat{\mathbf{A}}_T$, $\widehat{\mathbf{A}}_R$, $\widehat{\mathbf{B}}_h$, and $\widehat{\mathbf{B}}_v$. Then, the path gain vector \mathbf{g} can be estimated from (8) (or $[\mathcal{Y}]_{(4)}$) using a LS method as shown by step 11. Finally, the cascaded channel matrix $\widehat{\mathbf{H}}_c$ can be reconstructed as in step 12, which can be used to estimate $\widehat{\mathbf{H}}_T$ and $\widehat{\mathbf{H}}_R$, up to trivial scaling factors, using the LS Khatri-Rao factorization (LSKRF) method [17].

Uniqueness and identifiability conditions: It is well known that the CP decomposition is unique up to scaling and permutation ambiguities under mild conditions [24], [26]–[29]. In general, the uniqueness of a CP decomposition is guaranteed by Kruskal’s condition [27], which is also known as the k -rank. However, due to the definitions of Ω_R and Ω_T , the first two factor matrices, i.e., $\widehat{\mathbf{A}}_R \Omega_R = \widehat{\mathbf{A}}_R$ and $\widehat{\mathbf{A}}_T \Omega_T = \widehat{\mathbf{A}}_T$ contain repeated columns, where every column of $\widehat{\mathbf{A}}_R$ is repeated L_T times and every column of $\widehat{\mathbf{A}}_T$ is repeated L_R times. This implies that the k -rank of $\widehat{\mathbf{A}}_R$ and $\widehat{\mathbf{A}}_T$ is equal to one. Therefore, the sufficient condition of [27] fails [29]. As for Algorithm 1, which is an ALS-based algorithm, the identifiability in the LS sense requires that each of the following matrices: $\mathbf{C}_R = \Omega_R (\widehat{\mathbf{B}}_v \diamond \widehat{\mathbf{B}}_h \diamond \widehat{\mathbf{A}}_T \Omega_T)^\top \in \mathbb{C}^{L_R \times J_R}$, $\mathbf{C}_T = \Omega_T (\widehat{\mathbf{B}}_v \diamond \widehat{\mathbf{B}}_h \diamond \widehat{\mathbf{A}}_R \Omega_R)^\top \in \mathbb{C}^{L_T \times J_T}$, $\mathbf{C}_h = (\widehat{\mathbf{B}}_v \diamond \widehat{\mathbf{A}}_T \Omega_T \diamond \widehat{\mathbf{A}}_R \Omega_R)^\top \in \mathbb{C}^{L \times J_S^h}$, and $\mathbf{C}_v = (\widehat{\mathbf{B}}_h \diamond \widehat{\mathbf{A}}_T \Omega_T \diamond \widehat{\mathbf{A}}_R \Omega_R)^\top \in \mathbb{C}^{L \times J_S^v}$ to have a unique right Moore-Penrose pseudo-inverse,

i.e., full row-rank, where $J_R = K_T K_S$, $J_T = K_R K_S$, $J_S^h = K_R K_T K_S^v$, and $J_S^v = K_R K_T K_S^h$. This requires that $J_R \geq L_R$, $J_T \geq L_T$, $J_S^h \geq L$, and $J_S^v \geq L$, where $L = L_R \cdot L_T$. Since L_R and L_T are practically very small (i.e., $\max\{L_R, L_T\} \approx 3$ [18]), the above conditions are easily satisfied. For example, assuming that the TX is in line-of-sight with the RIS, we have that $L_T = 1$, as it has been assumed in [8].

Complexity analysis: Assuming that the complexity of calculating the Moore-Penrose pseudo-inverse of a $n \times m$ matrix is on the order of $\mathcal{O}(\min\{n, m\}^3)$. Then, the complexity of the ALS steps in Alg. 1 is on the order of $\mathcal{O}(I_{\max}(L_R^3 + L_T^3 + 2L^3))$. Moreover, assuming that the NOMP method from [25] is used in step 10, then the complexity of recovering the channel parameters is on the order of $\bar{L}(L_R + L_T + 2L)$, where \bar{L} denotes the number of grid points used by NOMP in the sparse-coding stage. In comparison, the complexity of TRICE-CS [13] is on the order of $\mathcal{O}(L(K_R K_T (\bar{L}^2 + L + L^2)) + 2L^3 + L K_S \bar{L}^2)$ and the Joint-CS method [14] is on the order of $\mathcal{O}(L(N_R K_T K_S (\bar{L}^4 + L + L^2)) + L^3)$. Clearly, TenRICE has a much lower complexity compared to both methods. The main reason is that TRICE and Joint-CS require multidimensional (xD) dictionaries (2D for TRICE and 4D for Joint-CS) compared to the 1D dictionary required by TenRICE. Moreover, in contrast to the TenRICE, TRICE and Joint-CS methods require a dictionary orthogonalization operation during the parameter recovery [30], which is very complex especially with large dictionaries.

IV. PHASE 2: THE PROPOSED RIS REFLECTION DESIGN METHOD (FROMAX)

In this section, given the estimated channels $\widehat{\mathbf{H}}_R$ and $\widehat{\mathbf{H}}_T$, we design the TX and the RX beamforming matrices and the RIS reflection coefficient vector as a solution to the following SE maximization problem:

$$\begin{aligned} \max_{\mathbf{Q}, \mathbf{P}, \boldsymbol{\omega}} \quad & \log_2 \det(\mathbf{I}_{N_s} + \mathbf{R}^{-1} \mathbf{Q}^H \widehat{\mathbf{H}}_c \mathbf{P} \mathbf{P}^H \widehat{\mathbf{H}}_c^H \mathbf{Q}) \\ \text{s.t.} \quad & \|\mathbf{P}\|_F^2 \leq P_{\max} \text{ and } |[\boldsymbol{\omega}]_m| = 1/\sqrt{M_S}, \forall m, \end{aligned} \quad (20)$$

where $\widehat{\mathbf{H}}_c \stackrel{\text{def}}{=} \widehat{\mathbf{H}}_R \text{diag}\{\boldsymbol{\omega}\} \widehat{\mathbf{H}}_T$ and P_{\max} is the transmit power at the TX. Note that (20) is non-convex, since the objective function is non-concave over $\boldsymbol{\omega}$ and the constant modulus constraints are non-convex functions. Moreover, \mathbf{P} , \mathbf{Q} , and $\boldsymbol{\omega}$ depend on each other, which makes (20) a difficult problem to solve. In the following, we propose a non-iterative solution to (20), which has a comparable performance to that of [2], but with a much lower complexity.

Initially, it is not hard to see that for any given $\boldsymbol{\omega}$, (20) reduces to a single-user multi-stream MIMO communication system. Let $\widehat{\mathbf{H}}_c = \mathbf{U}_{\widehat{\mathbf{H}}_c} \boldsymbol{\Sigma}_{\widehat{\mathbf{H}}_c} \mathbf{V}_{\widehat{\mathbf{H}}_c}^H$ be the singular value decomposition (SVD) of $\widehat{\mathbf{H}}_c$. Then, the optimal fully-digital³

³Here, we note that in mmWave-based communications, hybrid analog-digital (HAD) beamforming architectures [31]–[34] are generally assumed to reduce the power consumption. However, since in this section we focus on the RIS reflection coefficient design, we assume fully-digital beamforming architectures at the TX and the RX, to simplify the exposition.

solutions to \mathbf{Q} and \mathbf{P} , for fixed $\boldsymbol{\omega}$, are given as

$$\mathbf{Q} = \mathbf{U}_s \text{ and } \mathbf{P} = \mathbf{V}_s \text{diag}\{\sqrt{p_1}, \dots, \sqrt{p_{N_s}}\}, \quad (21)$$

where $\mathbf{U}_s = [\mathbf{U}_{\widehat{\mathbf{H}}_c}]_{[:,1:N_s]}$, $\mathbf{V}_s = [\mathbf{V}_{\widehat{\mathbf{H}}_c}]_{[:,1:N_s]}$, and $\{p_i\}_{i=1}^{N_s}$ are the power allocations found using the waterfilling method [35] such that $\sum_{i=1}^{N_s} p_i = P_{\max}$. Consequently, $\mathbf{Q}^H \mathbf{Q} = \mathbf{I}_{N_s}$, $\boldsymbol{\Sigma}_s = \mathbf{U}_s^H \widehat{\mathbf{H}}_R \text{diag}\{\boldsymbol{\omega}\} \widehat{\mathbf{H}}_T \mathbf{V}_s = \text{diag}\{\alpha_1, \dots, \alpha_{N_s}\}$, and the SE expression in (4) simplifies to

$$\text{SE} = \sum_{i=1}^{N_s} \log_2 \left(1 + \frac{1}{\sigma^2} \alpha_i^2 p_i \right), \quad (22)$$

where α_i is the i th dominant singular value in $\boldsymbol{\Sigma}_{\widehat{\mathbf{H}}_c}$. In the following, we turn our attention to the RIS reflection coefficient design and propose an efficient non-iterative solution to find $\boldsymbol{\omega}$ based on a FroMax design strategy.

FroMax-1: As a baseline method, the RIS reflection vector is found as a solution to

$$\begin{aligned} \boldsymbol{\omega} &= \arg \max_{\boldsymbol{\omega}} \|\widehat{\mathbf{H}}_R \text{diag}\{\boldsymbol{\omega}\} \widehat{\mathbf{H}}_T\|_F^2 = \arg \max_{\boldsymbol{\omega}} \|\mathbf{K}\boldsymbol{\omega}\|_2^2 \\ \text{s.t. } & |[\boldsymbol{\omega}]_m| = 1/\sqrt{M_S}, \forall m, \end{aligned} \quad (23)$$

where $\mathbf{K} \stackrel{\text{def}}{=} \widehat{\mathbf{H}}_T^H \diamond \widehat{\mathbf{H}}_R$ is obtained by applying Property 1. Note that (23) is non-convex due to the constant modulus constraints. Therefore, we first seek a solution to the following relaxed and convex version of (23) given as

$$\hat{\boldsymbol{\omega}} = \arg \max_{\hat{\boldsymbol{\omega}}} \|\mathbf{K}\hat{\boldsymbol{\omega}}\|_2^2, \quad \text{s.t. } \|\hat{\boldsymbol{\omega}}\|_2 = 1. \quad (24)$$

Let $\mathbf{K} = \mathbf{U}_K \boldsymbol{\Sigma}_K \mathbf{V}_K^H$ be the SVD of \mathbf{K} . Then, the optimal solution to (24) is given as $\hat{\boldsymbol{\omega}} = [\mathbf{V}_K]_{[:,1]}$. To satisfy the constant modulus constraints of (23), we use a simple projection function, where the m th entry of $\boldsymbol{\omega}$ is given as

$$[\boldsymbol{\omega}^{\text{FroMax-1}}]_m = \frac{1}{\sqrt{M_S}} \cdot \left(\frac{[\hat{\boldsymbol{\omega}}]_m}{|[\hat{\boldsymbol{\omega}}]_m|} \right). \quad (25)$$

However, using computer simulations, we have observed that FroMax-1 mainly maximizes the dominant singular value of $\widehat{\mathbf{H}}_c$, which makes it limited to single-stream scenarios.

FroMax-2: From (22), we can clearly see that $\boldsymbol{\omega}$ should be designed so that the singular values α_i are maximized. Thus, we propose to modify (23) as

$$\begin{aligned} \boldsymbol{\omega} &= \arg \max_{\boldsymbol{\omega}} \|\boldsymbol{\Sigma}_s\|_F^2 = \arg \max_{\boldsymbol{\omega}} \|\mathbf{D}\boldsymbol{\omega}\|_2^2 \\ \text{s.t. } & |[\boldsymbol{\omega}]_m| = 1/\sqrt{M_S}, \forall m, \end{aligned} \quad (26)$$

where \mathbf{D} , due to the diagonal structure of $\boldsymbol{\Sigma}_s$, is given as

$$\mathbf{D} \stackrel{\text{def}}{=} \begin{bmatrix} [\mathbf{V}_s]_{[:,1]}^T \widehat{\mathbf{H}}_T^H \diamond [\mathbf{U}_s]_{[:,1]}^H \widehat{\mathbf{H}}_R \\ \vdots \\ [\mathbf{V}_s]_{[:,N_s]}^T \widehat{\mathbf{H}}_T^H \diamond [\mathbf{U}_s]_{[:,N_s]}^H \widehat{\mathbf{H}}_R \end{bmatrix} \in \mathbb{C}^{N_s \times M_S}. \quad (27)$$

Similarly to (24), (26) can be relaxed to a convex form as

$$\bar{\boldsymbol{\omega}} = \arg \max_{\bar{\boldsymbol{\omega}}} \|\mathbf{D}\bar{\boldsymbol{\omega}}\|_2^2, \quad \text{s.t. } \|\bar{\boldsymbol{\omega}}\|_2 = 1. \quad (28)$$

However, differently from (24), we propose a solution that achieves a higher SE, where $\bar{\boldsymbol{\omega}}$ is obtained by taking the

Algorithm 2 FroMax-based methods for RIS reflection design.

- 1: Input: $\widehat{\mathbf{H}}_T$, $\widehat{\mathbf{H}}_R$, and P_{\max}
 - 2: **if** FroMax-1 based method **then**
 - 3: Construct \mathbf{K} as in (23) and get $\hat{\boldsymbol{\omega}}$ from \mathbf{V}_K
 - 4: Obtain $\boldsymbol{\omega}^* \leftarrow \boldsymbol{\omega}^{\text{FroMax-1}}$ using (25)
 - 5: **else if** FroMax-2 based method **then**
 - 6: Compute $\mathbf{U}_s = [\mathbf{U}_{\widehat{\mathbf{H}}_R}]_{[:,1:N_s]}$ and $\mathbf{V}_s = [\mathbf{V}_{\widehat{\mathbf{H}}_T}]_{[:,1:N_s]}$
 - 7: Construct \mathbf{D} as in (27) and get $\bar{\boldsymbol{\omega}}$ from \mathbf{V}_D
 - 8: Obtain $\boldsymbol{\omega}^* \leftarrow \boldsymbol{\omega}^{\text{FroMax-2}}$ using (29)
 - 9: **end if**
 - 10: For given $\boldsymbol{\omega}^*$, obtain \mathbf{Q} and \mathbf{P} as in (21)
-

contributions of the dominant N_s right singular vectors of \mathbf{D} . Specifically, let $\mathbf{D} = \mathbf{U}_D \boldsymbol{\Sigma}_D \mathbf{V}_D^H$ be the SVD of \mathbf{D} . Then, the proposed solution is given as $\bar{\boldsymbol{\omega}} = \frac{[\mathbf{V}_D]_{[:,1]} + \dots + [\mathbf{V}_D]_{[:,N_s]}}{\|[\mathbf{V}_D]_{[:,1]} + \dots + [\mathbf{V}_D]_{[:,N_s]}\|_2}$. Using $\bar{\boldsymbol{\omega}}$, the RIS reflection vector $\boldsymbol{\omega}$ is obtained as

$$[\boldsymbol{\omega}^{\text{FroMax-2}}]_m = \frac{1}{\sqrt{M_S}} \cdot \left(\frac{[\bar{\boldsymbol{\omega}}]_m}{|[\bar{\boldsymbol{\omega}}]_m|} \right), \forall m. \quad (29)$$

Remark 1: From (27), it is clear that the unitary matrices \mathbf{U}_s and \mathbf{V}_s are required to construct \mathbf{D} . However, since \mathbf{U}_s and \mathbf{V}_s depend on $\boldsymbol{\omega}$, an iterative two-step algorithm is required, where we update \mathbf{U}_s and \mathbf{V}_s in one step and $\boldsymbol{\omega}$ in the other step. However, we found that if \mathbf{U}_s and \mathbf{V}_s are appropriately initialized, then one iteration of such an algorithm is sufficient to have a comparable SE performance to that obtained by the iterative method of [2]. Here, we propose to initialize \mathbf{U}_s and \mathbf{V}_s as follows. Let $\widehat{\mathbf{H}}_R = \mathbf{U}_{\widehat{\mathbf{H}}_R} \boldsymbol{\Sigma}_{\widehat{\mathbf{H}}_R} \mathbf{V}_{\widehat{\mathbf{H}}_R}^H$ and $\widehat{\mathbf{H}}_T = \mathbf{U}_{\widehat{\mathbf{H}}_T} \boldsymbol{\Sigma}_{\widehat{\mathbf{H}}_T} \mathbf{V}_{\widehat{\mathbf{H}}_T}^H$ be the SVD of $\widehat{\mathbf{H}}_R$ and $\widehat{\mathbf{H}}_T$, respectively. Then, we assume that \mathbf{U}_s and \mathbf{V}_s in (26) are given as $\mathbf{U}_s = [\mathbf{U}_{\widehat{\mathbf{H}}_R}]_{[:,1:N_s]}$ and $\mathbf{V}_s = [\mathbf{V}_{\widehat{\mathbf{H}}_T}]_{[:,1:N_s]}$.

In summary, the proposed beamforming and RIS reflection coefficient design method is summarized in Algorithm 2.

Complexity analysis: Let the complexity of calculating the SVD⁴ of a $n \times m$ matrix on the order of $\mathcal{O}(nm^2)$. Then, the complexity of Algorithm 2 steps 3, 6, 7, and 10 is on the order of $\mathcal{O}(M_R M_T M_S^2)$, $\mathcal{O}(M_R M_S^2 + M_S M_T^2)$, $\mathcal{O}(N_s M_S^2)$, and $\mathcal{O}(M_R M_T^2)$ respectively. Accordingly, the complexity of FroMax-1 is on the order of $\mathcal{O}(M_R M_T M_S^2 + M_R M_T^2)$ and of FroMax-2 is on the order of $\mathcal{O}(M_R M_S^2 + M_S M_T^2 + N_s M_S^2 + M_R M_T^2)$. In comparison, the complexity of the alternation maximization (AltMax) method of [2] is on the order of $\mathcal{O}(J_{\max}(M_S(3M_R^3 + 2M_R^2 M_T + M_T^2) + M_R M_T^2))$, where J_{\max} is the maximum number of iterations.

V. NUMERICAL RESULTS

In this section, we show simulation results to evaluate the effectiveness of the proposed methods. In all simulation results, we assume that $M_T = 64$, $M_R = 16$, and $M_S^h = M_S^v = 16$, i.e., the RIS has $M_S = 256$ reflecting elements.

⁴Note that the complexity of calculating the SVD of $n \times m$ matrix can be reduced by using the *Power Iteration* method. However, to simplify the analysis, we assume that the SVD is calculated using the bidiagonalization and QR algorithm with a complexity on the order of $\mathcal{O}(nm^2)$.

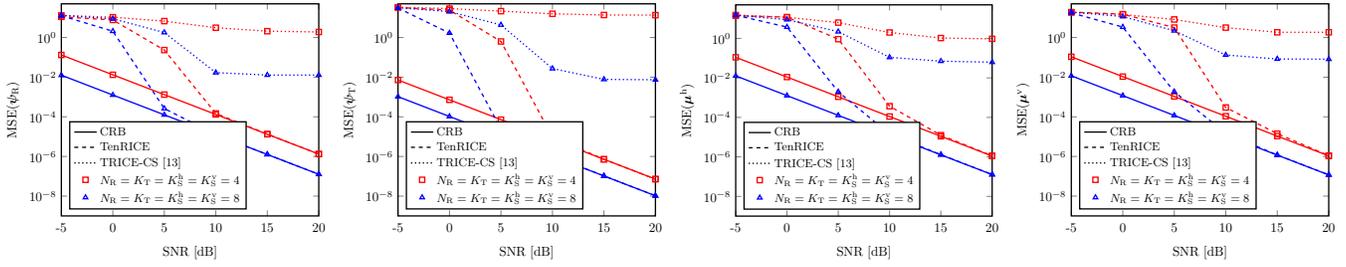


Fig. 3. MSE vs. SNR [$L_T = L_R = 2$].

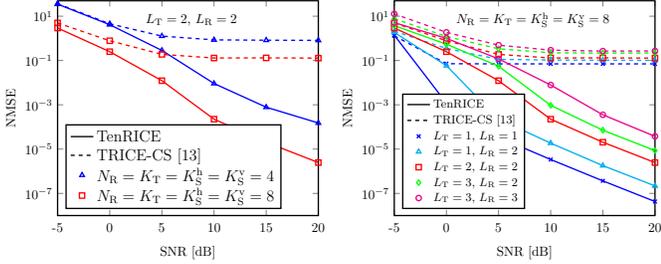


Fig. 4. NMSE vs. SNR.

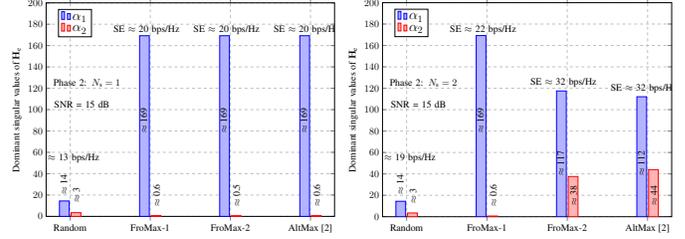


Fig. 6. Dominant singular values of a perfect effective channel $\mathbf{H}_e = \mathbf{H}_R \text{diag}\{\boldsymbol{\omega}\} \mathbf{H}_T$ for one channel realization [$L_T = L_R = 2$].

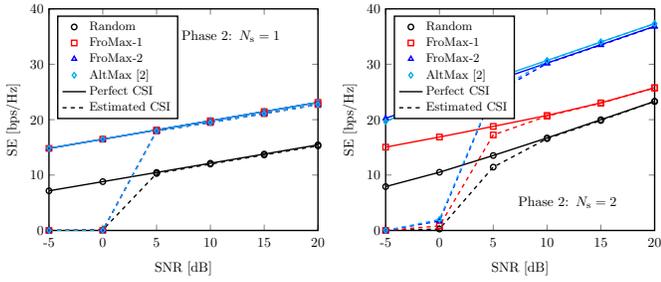


Fig. 5. SE vs. SNR. [$L_T = L_R = 2$. Phase 1: $K_R = K_T = K_S^h = K_S^v = 8$]

Phase 1 - CE: In the CE phase, we assume that the training matrices \mathbf{W} , \mathbf{F} , Φ^h , and Φ^v in (2) are randomly generated such that the (i, j) th entry of \mathbf{W} is given as $[\mathbf{W}]_{[i,j]} = \frac{1}{\sqrt{M_R}} e^{j\varphi_{i,j}}$, $\varphi_{i,j} \in [0, 2\pi]$, where \mathbf{F} , Φ^h , and Φ^v are similarly generated. We show results in terms of the mean-squared error (MSE) of ψ_R defined as $\text{MSE}(\psi_R) = \mathbb{E}\{\|\psi_R - \hat{\psi}_R\|_2^2\}$, where $\text{MSE}(\psi_T)$, $\text{MSE}(\boldsymbol{\mu}^h)$, and $\text{MSE}(\boldsymbol{\mu}^v)$ are similarly defined, and the normalized MSE (NMSE) of the cascaded channel is defined as $\text{NMSE} = \mathbb{E}\{\|\mathbf{H}_c - \hat{\mathbf{H}}_c\|_F^2\} / \mathbb{E}\{\|\mathbf{H}_c\|_F^2\}$. We define the signal-to-noise ratio (SNR) as $\text{SNR} = \mathbb{E}\{\|\mathbf{Y} - \mathbf{Z}\|_F^2\} / \mathbb{E}\{\|\mathbf{Z}\|_F^2\}$. For comparison, we include simulation results of the two-stage TRICE-CS framework [13], where the estimation is performed using the classical OMP technique [30] assuming a 2D dictionary of 128×128 grid points in both stages.

Figs. 3 and 4 show the MSE versus the SNR and the NMSE versus the SNR results, respectively, averaged over 1,000 channel realizations. From Fig. 3, we can see that TenRICE provides more accurate parameter estimates, compared to TRICE-CS, approaching the CRB⁵ as the SNR increases. The

⁵The CRB derivation to our 4-way CP tensor is a straightforward extension of the CRB derivation in [26] for a 3-way CP tensor. Therefore, it has been omitted here due to brevity.

main reason is that TenRICE not only exploits the low-rank nature of mmWave channels, but also the tensor structure of the received signals when estimating the channel parameters. Moreover, TenRICE employs a high-resolution parameter recovery method in NOMP, while TRICE-CS suffers from quantization errors, due to the on-grid assumption. These advantages lead to more accurate channel estimates, as can be seen from Fig. 4, with less training overhead and lower complexity.

Phase 2 - DT: Next, we show simulation results to illustrate the efficiency of the proposed RIS reflection design method, FroMax. For comparison, we include results when the RIS reflection coefficient vector $\boldsymbol{\omega}$ is designed according to the alternating maximization method in [2], termed AltMax, and Random, where the entries of $\boldsymbol{\omega}$ are randomly generated such that the m th entry is given as $[\boldsymbol{\omega}]_m = \frac{1}{\sqrt{M_S}} e^{j\omega_m}$, $\omega_m \in [0, 2\pi]$. We define the SNR as $\text{SNR} = P_{\max} / \sigma^2$.

Fig. 5 shows SE versus SNR results, averaged over 1,000 channel realizations. Clearly, we can see that FroMax-1 has an equal performance to that of FroMax-2 and AltMax when $N_s = 1$. However, FroMax-1 experiences a performance loss when $N_s = 2$, since it mainly maximizes the dominant singular value, as it can be seen from Fig. 6. Differently, the AltMax and FroMax-2 methods optimize the dominant N_s singular values of the effective channel such that it maximizes the system SE. Note that, in the low SNR regime, i.e., below 5 dB, all the simulated methods experience a very low SE performance, due to the CE errors. Therefore, a preprocessing *denoising* step will be required to improve the CE accuracy, which we leave for future work.

VI. CONCLUSIONS

In this work, we have considered the channel estimation and the RIS reflection coefficient design problems in point-to-point RIS-aided mmWave MIMO communication systems. We have proposed a CP tensor-based channel estimation method termed TenRICE, which estimates the transmitter to RIS and the RIS to receiver channels separately, up to a trivial scaling factor. We have shown that by jointly exploiting the low-rank nature of mmWave channels and the tensor structure of the received signals, not only the estimation accuracy can be improved, but also the training overhead and the complexity can be reduced. The proposed non-iterative RIS reflection design method based on a Frobenius-norm maximization (FroMax) design strategy has a comparable performance to a benchmark method but with significantly lower complexity.

REFERENCES

- [1] M. Di Renzo, A. Zappone, M. Debbah, M.-S. Alouini, C. Yuen, J. de Rosny, and S. Tretjakov, "Smart radio environments empowered by reconfigurable intelligent surfaces: How it works, state of research, and the road ahead," *IEEE J. Sel. Areas Commun.*, vol. 38, no. 11, pp. 2450–2525, 2020.
- [2] S. Zhang and R. Zhang, "Capacity characterization for intelligent reflecting surface aided MIMO communication," *IEEE J. Sel. Areas Commun.*, vol. 38, no. 8, pp. 1823–1838, Aug. 2020.
- [3] 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, 2019.
- [4] L. Dong and H.-M. Wang, "Enhancing secure MIMO transmission via intelligent reflecting surface," *IEEE Trans. Wireless Commun.*, vol. 19, no. 11, pp. 7543–7556, Nov. 2020.
- [5] Q.-U.-A. Nadeem, A. Kammoun, A. Chaaban, M. Debbah, and M.-S. Alouini, "Asymptotic max-min SINR analysis of reconfigurable intelligent surface assisted MISO systems," *IEEE Trans. Wireless Commun.*, vol. 19, no. 12, pp. 7748–7764, Dec. 2020.
- [6] D. Mishra and H. Johansson, "Channel estimation and low-complexity beamforming design for passive intelligent surface assisted MISO wireless energy transfer," in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 4659–4663.
- [7] T. L. Jensen and E. De Carvalho, "An optimal channel estimation scheme for intelligent reflecting surfaces based on a minimum variance unbiased estimator," in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 5000–5004.
- [8] Q. Nadeem, H. Alwazani, A. Kammoun, A. Chaaban, M. Debbah, and M. Alouini, "Intelligent reflecting surface-assisted multi-user MISO communication: Channel estimation and beamforming design," *IEEE Open J. Commun. Soc.*, vol. 1, pp. 661–680, 2020.
- [9] J. Zhang, C. Qi, P. Li, and P. Lu, "Channel estimation for reconfigurable intelligent surface aided massive MIMO system," in *Proc. IEEE 21st International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, May 2020, pp. 1–5.
- [10] D. L. Donoho, "Compressed sensing," *IEEE Trans. Inf. Theory*, vol. 52, no. 4, pp. 1289–1306, Apr. 2006.
- [11] K. Ardah, B. Sokal, A. L. F. de Almeida, and M. Haardt, "Compressed sensing based channel estimation and open-loop training design for hybrid analog-digital massive MIMO systems," in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, May 2020, pp. 4597–4601.
- [12] K. Ardah, A. L. F. de Almeida, and M. Haardt, "A gridless CS approach for channel estimation in hybrid massive MIMO systems," in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, May 2019, pp. 4160–4164.
- [13] K. Ardah, S. Gherekhloo, A. L. F. de Almeida, and M. Haardt, "TRICE: A channel estimation framework for RIS-aided millimeter-wave MIMO systems," *IEEE Signal Process. Lett.*, vol. 28, pp. 513–517, Feb. 2021.
- [14] P. Wang, J. Fang, H. Duan, and H. Li, "Compressed channel estimation for intelligent reflecting surface-assisted millimeter wave systems," *IEEE Signal Process. Lett.*, vol. 27, pp. 905–909, 2020.
- [15] Z. He and X. Yuan, "Cascaded channel estimation for large intelligent metasurface assisted massive MIMO," *IEEE Wireless Commun. Lett.*, vol. 9, no. 2, pp. 210–214, Feb. 2020.
- [16] G. Favier and A. L. de Almeida, "Overview of constrained PARAFAC models," *EURASIP Journal on Applied Signal Processing*, vol. 2014, p. 142, Dec. 2014.
- [17] G. T. de Araújo and A. L. F. de Almeida, "PARAFAC-based channel estimation for intelligent reflective surface assisted MIMO system," in *Proc. IEEE 11th Sensor Array and Multichannel Signal Processing Workshop (SAM)*, 2020, pp. 1–5.
- [18] T. S. Rappaport, Y. Xing, G. R. MacCartney, A. F. Molisch, E. Mellios, and J. Zhang, "Overview of millimeter wave communications for fifth-generation (5G) wireless networks—with a focus on propagation models," *IEEE Transactions on Antennas and Propagation*, vol. 65, no. 12, pp. 6213–6230, 2017.
- [19] J. Zhang and M. Haardt, "Channel estimation and training design for hybrid multi-carrier mmwave massive MIMO systems: The beamspace ESPRIT approach," in *Proc. 25th European Signal Processing Conference (EUSIPCO)*, 2017, pp. 385–389.
- [20] A. L. F. de Almeida, G. Favier, and J. C. M. Mota, "A constrained factor decomposition with application to MIMO antenna systems," *IEEE Trans. Signal Process.*, vol. 56, no. 6, pp. 2429–2442, 2008.
- [21] L. De Lathauwer, "A link between the canonical decomposition in multilinear algebra and simultaneous matrix diagonalization," *SIAM Journal on Matrix Analysis and Applications*, vol. 28, no. 3, pp. 642–666, Apr. 2006.
- [22] K. Ardah, A. L. F. de Almeida, and M. Haardt, "Low-complexity millimeter wave CSI estimation in MIMO-OFDM hybrid beamforming systems," in *Proc. 23rd International ITG Workshop on Smart Antennas (WSA)*, Apr. 2019, pp. 1–5.
- [23] F. Roemer and M. Haardt, "A semi-algebraic framework for approximate CP decompositions via simultaneous matrix diagonalizations (SECSI)," *Signal Processing*, vol. 93, no. 9, pp. 2722 – 2738, 2013.
- [24] P. Comon, X. Luciani, and A. L. F. de Almeida, "Tensor decompositions, alternating least squares and other tales," *Journal of Chemometrics*, vol. 23, no. 7-8, pp. 393–405, 2009.
- [25] B. Mamanidipoor, D. Ramasamy, and U. Madhoo, "Newtonized orthogonal matching pursuit: Frequency estimation over the continuum," *IEEE Trans. Signal Process.*, vol. 64, no. 19, pp. 5066–5081, Oct. 2016.
- [26] Z. Zhou, J. Fang, L. Yang, H. Li, Z. Chen, and R. S. Blum, "Low-rank tensor decomposition-aided channel estimation for millimeter wave MIMO-OFDM systems," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 7, pp. 1524–1538, Jul. 2017.
- [27] T. G. Kolda and B. W. Bader, "Tensor decompositions and applications," *SIAM Review*, vol. 51, no. 3, pp. 455–500, Sept. 2009.
- [28] A. L. F. de Almeida, G. Favier, and J. C. M. Mota, "Constrained tensor modeling approach to blind multiple-antenna CDMA schemes," *IEEE Trans. Signal Process.*, vol. 56, no. 6, pp. 2417–2428, Jun. 2008.
- [29] A. Stegeman and A. L. F. de Almeida, "Uniqueness conditions for constrained three-way factor decompositions with linearly dependent loadings," *SIAM Journal on Matrix Analysis and Applications*, vol. 31, no. 3, pp. 1469–1490, 2010.
- [30] B. L. Sturm and M. G. Christensen, "Comparison of orthogonal matching pursuit implementations," in *Proc. of the 20th European Signal Processing Conference (EUSIPCO)*, Aug. 2012, pp. 220–224.
- [31] R. W. Heath, N. González-Prelcic, S. Rangan, W. Roh, and A. M. Sayeed, "An overview of signal processing techniques for millimeter wave MIMO systems," *IEEE J. Sel. Topics Signal Process.*, vol. 10, no. 3, pp. 436–453, 2016.
- [32] S. Gherekhloo, K. Ardah, and M. Haardt, "Hybrid beamforming design for downlink MU-MIMO-OFDM millimeter-wave systems," in *Proc. IEEE 11th Sensor Array and Multichannel Signal Processing Workshop (SAM)*, Jun. 2020, pp. 1–5.
- [33] K. Ardah, G. Fodor, Y. C. B. Silva, W. C. Freitas, and F. R. P. Cavalcanti, "A unifying design of hybrid beamforming architectures employing phase shifters or switches," *IEEE Trans. Veh. Technol.*, vol. 67, no. 11, pp. 11 243–11 247, Nov. 2018.
- [34] K. Ardah, G. Fodor, Y. C. B. Silva, W. C. Freitas, and A. L. F. de Almeida, "Hybrid analog-digital beamforming design for SE and EE maximization in massive MIMO networks," *IEEE Trans. Veh. Technol.*, vol. 69, no. 1, pp. 377–389, Jan. 2020.
- [35] D. Palomar and J. Fonollosa, "Practical algorithms for a family of waterfilling solutions," *IEEE Trans. Signal Process.*, vol. 53, no. 2, pp. 686–695, 2005.