

Reducing the Control Overhead of Intelligent Reconfigurable Surfaces Via a Tensor-Based Low-Rank Factorization Approach

Bruno Sokal, *Student Member, IEEE*, Paulo R. B. Gomes, André L. F. de Almeida, *Senior Member, IEEE*, Behrooz Makki, *Senior Member, IEEE*, and Gabor Fodor, *Senior Member, IEEE*

Abstract—Passive intelligent reconfigurable surfaces (IRS) are becoming an attractive component of cellular networks due to their ability of shaping the propagation environment and thereby improving the coverage. While passive IRS nodes incorporate a great number of phase-shifting elements and a controller entity, the phase-shifts are typically determined by the cellular base station (BS) due to its computational capability. Since the fine granularity control of the large number of phase-shifters may become prohibitive in practice, it is important to reduce the control overhead between the BS and the IRS controller. To this end, in this paper we propose a low-rank approximation of the near-optimal phase-shifts, which would incur prohibitively high communication overhead on the BS-IRS controller links. The key idea is to represent the potentially large IRS phase-shift vector using a low-rank tensor model. This is achieved by factorizing a *tensorized* version of the IRS phase-shift vector, where each component is modeled as the Kronecker product of a predefined number of factors of smaller sizes, which can be obtained *via* tensor decomposition algorithms. We show that the proposed low-rank models drastically reduce the required feedback requirements associated with the BS-IRS control links.

Our simulation results indicate that the proposed method is especially attractive in scenarios with a strong line of sight component, in which case nearly the same spectral efficiency is reached as in the cases with near-optimal phase-shifts, but with a drastically reduced communication overhead.

Index Terms—Reconfigurable intelligent surface (RIS), feedback overhead, control signaling, low-rank approximation, tensor modeling, PARAFAC, Tucker.

I. INTRODUCTION

Intelligent reconfigurable surface (IRS) is a candidate technology for beyond fifth generation and sixth generation networks due to its ability to *control* the electromagnetic properties of the radio-frequency waves by performing an intelligent phase-shift to the desired direction [2]–[9]. Usually, IRS is defined as a planar (2-D) surface with a large number

of independent reflective elements, in which they can be fully passive or with some elements active [10]–[12]. IRS is connected to a smart controller that sets the desired phase-shift for each reflective element, by applying bias voltages at the elements e.g., PIN diodes. The main advantage of fully passive IRSs is its full-duplex nature, i.e., no noise amplification is observed since no signal processing is possible. However, the fully passive nature of the IRSs makes the channel state information (CSI) acquisition process difficult, since no pilots are processed, thus only the cascade channel can be estimated. Nevertheless, in the case of employing a few active elements in the IRS, this issue is suppressed and channel can be estimated using, for example, compressed sensing tools [10]. Another advantage of an IRS with fully passive elements is that the power consumption is concentrated at the controller. This makes the IRS a more attractive technology in terms of energy efficiency compared to alternative technologies, e.g., amplify-and-forward and decode-and-forward relays [13]–[15].

Several works have addressed the CSI acquisition problem in IRS-assisted networks, e.g., [16]–[23]. The work of [16] proposes a tensor-based method where the authors show the benefits of exploiting the multidimensional structure of the received signal by separating the cascade channel. The work of [18] proposes a compressed sensing approach in a multi-user uplink multiple-input multiple-output (MIMO) scenario. In [19], a two-timescale channel estimation framework is proposed to overcome the pilot overhead in a multi-user IRS-aided system. Also, [20] addresses the channel estimation problem in millimeter-wave MIMO systems. The work of [21] proposes a channel estimation framework for millimeter-wave (mmWave) IRS-assisted MIMO systems based on compressed sensing techniques. The authors of [22] propose a low-complexity framework for channel estimation and passive beamforming in MIMO IRS-assisted systems.

Although many works focus on channel estimation [16]–[23], achievable rate maximization [23]–[25], energy efficiency (EE) maximization [26]–[30], and interference mitigation problems [31]–[33], few works have addressed the problem of reducing the channel training or the feedback-overhead of IRS phase-shifts to the IRS controller. The work of [24] proposes a protocol design to maximize the transmission rate in IRS-assisted MIMO-OFDM systems. Also, [34] proposes a framework for overhead-aware feedback and resource allocation in IRS-assisted MIMO systems. The main idea of [34] is to optimize the network resource such as the bandwidth and the total power used for transmission and feedback. However, the number of phase-shifts to be conveyed

Bruno Sokal, Paulo R. B. Gomes, and André L. F. de Almeida are with the Wireless Telecom Research Group (GTel), Department of Teleinformatics Engineering, Federal University of Ceará, Fortaleza-CE, Brazil. E-mails: {brunosokal, andre}@gtel.ufc.br.

Behrooz Makki is with Ericsson Research, Göteborg, Sweden. E-mail: behrooz.makki@ericsson.com

Gabor Fodor is with Ericsson Research and KTH Royal Institute of Technology, Stockholm, Sweden. E-mail: gabor.fodor@ericsson.com

This work was supported by the Ericsson Research, Sweden, and Ericsson Innovation Center, Brazil, under UFC.48 Technical Cooperation Contract Ericsson/UFC. This study was financed in part by the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior - Brasil (CAPES)-Finance Code 001, and CAPES/PRINT Proc. 88887.311965/2018-00. André L. F. de Almeida acknowledges CNPq for its financial support under the grant 312491/2020-4. G. Fodor was partially supported by the Digital Futures project PERCy.

Part of this work has been submitted for possible presentation in IEEE GLOBECOM 2022 [1].

to the IRS controller can still be large and results in high feedback signaling overhead, especially for large IRS panels.

In this work, we propose an overhead-aware model for designing the IRS phase-shifts. Our idea is to represent the IRS phase-shift vector with a low-rank model. This is achieved by factorizing a *tensorized* version of the IRS phase-shift vector, where each component is modelled as the Kronecker product of a predefined number of factors. These factors are estimated using tensor decompositions such as the PARAllel FACTors (PARAFAC) [35] and Tucker [36]. After the estimation process, the phases of the factors are quantized and fed back to the IRS controller, which can reconstruct the IRS phase-shift vector based on the chosen low-rank tensor model. The main contributions of this work are the following:

- 1) Our proposed IRS phase-shift factorization allows to save network resources by reducing the total IRS phase-shift feedback overhead. This allows a more frequent IRS phase-shift feedback, for a fixed feedback load, which can significantly improve the end-to-end latency, crucial in fast varying channels, high mobility scenarios and/or the cases with moderate/large sizes of the IRS. Also, thanks to the significant reduction on the feedback overhead, the IRS-assisted network can decide to multiplex phase-shifts associated with a higher number of users in the same feedback channel.
- 2) The proposed IRS phase-shift factorization provides a flexible feedback design by controlling the parameters of the low-rank factorization model, such as the number of components, the number and the size of the factors, as well as their respective resolution. This is an important feature of our proposed feedback-aware model, since for limited feedback control links, the low-rank model and its factorization parameters can be efficiently adjusted to the available capacity of the feedback link, providing more degrees of freedom to system design.
- 3) Our tensor-based factorization approach relies on the optimum IRS phase-shift vector, which means that it can be implemented in every IRS-assisted network and in multiple communication links, i.e., downlink or uplink, in single-input single-output, multiple-input single-output, as well in MIMO systems.

Different from the works of [24] and [34], we aim to reduce the IRS phase-shifts feedback overhead by conveying to the IRS controller only the factors of our proposed low-rank model. Our approach is analytical and provides a systematic way of controlling the feedback overhead by adjusting the parameters of the low-rank IRS model, namely, its rank and the corresponding number of factors of each rank-one component. Our simulations show that the proposed low-rank model for the IRS phase-shifts can achieve the same spectral efficiency (SE) as the state-of-the-art in line of sight (LOS) scenarios, while the feedback payload (number of bits to be fed back) is dramatically reduced. For example, taking an IRS with $N = 1024$ elements, the feedback duration can be 50 times smaller than the state-of-the-art, depending on the low-rank model parameters. Also, when taking into account the total system SE and EE, i.e., both the IRS phase-shift feedback duration, and the channel estimation duration, our proposed model outperforms the state-of-the-art.

The rest of the paper is organized as follows. Section

II provides an introduction of the tensor notation and decompositions that are exploited in this paper. The system model is described in Section III. Section IV details our proposed feedback overhead-aware method and provides the details of the PARAFAC-IRS and Tucker-IRS models for IRS phase-shift vector factorization. Section V describes the phase-shift and weighting factors quantization procedure and the reconstruction of the IRS phase-shift vector at the IRS controller. The effects of the factorization parameters and the quantization process are also discussed in this section. Simulation results are provided in Section VI and the final conclusions and perspectives are discussed in Section VII.

A. Notation and Properties

Scalars are represented as non-bold lower-case letters a , column vectors as lower-case boldface letters \mathbf{a} , matrices as upper-case boldface letters \mathbf{A} , and tensors as calligraphic upper-case letters \mathcal{A} . The superscripts $\{\cdot\}^T$, $\{\cdot\}^*$, $\{\cdot\}^H$ and $\{\cdot\}^+$ stand for transpose, conjugate, conjugate transpose and pseudo-inverse operations, respectively. The operator $\|\cdot\|_F$ denotes the Frobenius norm of a matrix or tensor, and $\mathbb{E}\{\cdot\}$ is the expectation operator. The operator $\text{diag}(\mathbf{a})$ converts \mathbf{a} into a diagonal matrix, while $\text{diag}(\mathbf{A})$ returns a vector whose elements are the main diagonal of \mathbf{A} . Moreover, $\text{vec}(\mathbf{A})$ converts $\mathbf{A} \in \mathbb{C}^{I_1 \times R}$ to a column vector $\mathbf{a} \in \mathbb{C}^{I_1 R \times 1}$ by stacking its columns on top of each other, while the $\text{unvec}(\cdot)$ operator is the inverse of the vec operation. The symbol \circ denotes the outer product operator. Also, $\mathbf{a}_r \in \mathbb{C}^{I_1 \times 1}$ represents the r -th column of $\mathbf{A} \in \mathbb{C}^{I_1 \times R}$. Let us define two matrices $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_R] \in \mathbb{C}^{I_1 \times R}$ and $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_R] \in \mathbb{C}^{I_2 \times R}$. The Kronecker product between them is defined by

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{1,1}\mathbf{B} & \dots & a_{1,R}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{I_1,1}\mathbf{B} & & a_{I_1,R}\mathbf{B} \end{bmatrix} \in \mathbb{C}^{I_2 I_1 \times R R}.$$

The Khatri-Rao product, also known as the column wise Kronecker product, between two matrices, symbolized by \diamond , is defined as

$$\mathbf{A} \diamond \mathbf{B} = [\mathbf{a}_1 \otimes \mathbf{b}_1, \dots, \mathbf{a}_R \otimes \mathbf{b}_R] \in \mathbb{C}^{I_2 I_1 \times R R}.$$

We make use of the following properties

$$\text{vec}(\mathbf{ABC}) = (\mathbf{C}^T \otimes \mathbf{A}) \text{vec}(\mathbf{B}), \quad (1)$$

$$\text{vec}(\mathbf{A} \text{diag}(\mathbf{b}) \mathbf{C}) = (\mathbf{C}^T \diamond \mathbf{A}) \mathbf{b}, \quad (2)$$

$$\mathbf{a}^T \diamond \mathbf{B} = \mathbf{B} \text{diag}(\mathbf{a}), \quad (3)$$

$$\mathbf{a} \otimes \mathbf{b} = \text{vec}(\mathbf{b} \circ \mathbf{a}), \quad (4)$$

where the involved vectors and matrices have compatible dimensions in each case.

II. TENSOR PRE-REQUISITES

In this section, tensor preliminaries are provided by focusing on the main notation, operations and properties that will be useful in the rest of the paper.

Consider a set of matrices $\{\mathbf{X}_{i_3}\} \in \mathbb{C}^{I_1 \times I_2}$, for $i_3 = 1, \dots, I_3$. Concatenating all I_3 matrices, we form the third-order tensor $\mathcal{X} = [\mathbf{X}_1 \sqcup_3 \mathbf{X}_2 \sqcup_3 \dots \sqcup_3 \mathbf{X}_{I_3}] \in \mathbb{C}^{I_1 \times I_2 \times I_3}$, where \sqcup_3 indicates a concatenation in the third dimension. We can interpret \mathbf{X}_{i_3} as the i_3 -th frontal slice of \mathcal{X} , defined as

$\mathcal{X}_{..i_3} = \mathbf{X}_{i_3}$ where the “..” indicates that the dimensions I_1 and I_2 are fixed. The tensor \mathcal{X} can be *matricized* by letting one dimension vary along the rows and the remaining two dimensions along the columns. From \mathcal{X} , we can form three different matrices, referred to as the n -mode unfoldings (for $n = \{1, 2, 3\}$ in this case), which are respectively given by

$$[\mathcal{X}]_{(1)} = [\mathcal{X}_{..1}, \dots, \mathcal{X}_{..I_3}] \in \mathbb{C}^{I_1 \times I_2 I_3}, \quad (5)$$

$$[\mathcal{X}]_{(2)} = [\mathcal{X}_{1..}^T, \dots, \mathcal{X}_{I_3..}^T] \in \mathbb{C}^{I_2 \times I_1 I_3} \quad (6)$$

$$[\mathcal{X}]_{(3)} = [\text{vec}(\mathcal{X}_{..1}), \dots, \text{vec}(\mathcal{X}_{..I_3})]^T \in \mathbb{C}^{I_3 \times I_1 I_2}. \quad (7)$$

A. Tensorization

The tensorization operation consists of mapping the elements of a vector into a high-order tensor. Let us define the vector $\mathbf{y} \in \mathbb{C}^{N \times 1}$, in which $N = \prod_{p=1}^P N_p$, where N_p is the size of the p -th partition of this vector. By applying the tensorization operator, defined as $\mathcal{T}\{\cdot\}$, we can form a P -order tensor $\mathcal{Y} = \mathcal{T}\{\mathbf{y}\} \in \mathbb{C}^{N_1 \times N_2 \times \dots \times N_P}$. The mapping of elements from \mathbf{y} to \mathcal{Y} is defined as

$$\mathcal{Y}_{n_1, n_2, \dots, n_P} = \mathbf{y}_{n_1 + (n_2 - 1)N_1 + \dots + (n_P - 1)N_{P-1} \dots N_2 N_1}, \quad (8)$$

where $n_p = \{1, \dots, N_p\}$, for $p = \{1, \dots, P\}$. This operator plays a key role on the proposed method, and will be exploited to recast the IRS phase-shift vector as a tensor, from which the low-rank factorization schemes are proposed.

B. PARAFAC Decomposition

It is known that every matrix of rank R can be expressed as the summation of its rank-one components obtained by, e.g., singular value decomposition (SVD). In the case of tensors, a tensor of rank R is given by the summation of its rank-one tensor factors. This decomposition is called PARAFAC [35]. For a P order tensor $\mathcal{Y} \in \mathbb{C}^{I_1 \times I_2 \times \dots \times I_P}$, its PARAFAC decomposition is given as

$$\mathcal{Y} = \sum_{r=1}^R \mathbf{a}_r^{(1)} \circ \mathbf{a}_r^{(2)} \circ \dots \circ \mathbf{a}_r^{(P)} \in \mathbb{C}^{I_1 \times I_2 \times \dots \times I_P}, \quad (9)$$

where $\mathbf{a}_r^{(p)} \in \mathbb{C}^{I_p \times 1}$ is r -th column of the p -th factor matrix $\mathbf{A}^{(p)} \in \mathbb{C}^{I_p \times R}$, $p = \{1, \dots, P\}$. The p -th mode unfolding of \mathcal{Y} , defined as $[\mathcal{Y}]_{(p)} \in \mathbb{C}^{I_p \times I_1 \dots I_{p-1} I_{p+1} \dots I_P}$, is expressed as

$$[\mathcal{Y}]_{(p)} = \mathbf{A}^{(p)} \left(\mathbf{A}^{(1)} \diamond \dots \diamond \mathbf{A}^{(p+1)} \diamond \mathbf{A}^{(p-1)} \diamond \dots \diamond \mathbf{A}^{(1)} \right)^T. \quad (10)$$

Fig. 1 illustrates a PARAFAC tensor, for $P = 3$, as the summation of rank-one tensors. In this case, it can be shown that the three-mode unfoldings can be factorized as [37]

$$[\mathcal{Y}]_{(1)} = \mathbf{A}^{(1)} \left(\mathbf{A}^{(3)} \diamond \mathbf{A}^{(2)} \right)^T \in \mathbb{C}^{I_1 \times I_2 I_3}, \quad (11)$$

$$[\mathcal{Y}]_{(2)} = \mathbf{A}^{(2)} \left(\mathbf{A}^{(3)} \diamond \mathbf{A}^{(1)} \right)^T \in \mathbb{C}^{I_2 \times I_1 I_3}, \quad (12)$$

$$[\mathcal{Y}]_{(3)} = \mathbf{A}^{(3)} \left(\mathbf{A}^{(2)} \diamond \mathbf{A}^{(1)} \right)^T \in \mathbb{C}^{I_3 \times I_1 I_2}. \quad (13)$$

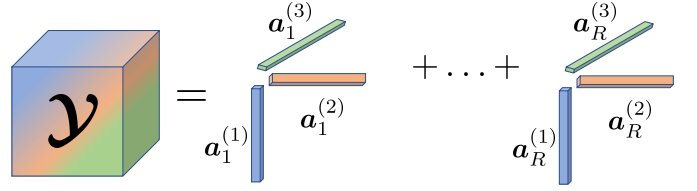


Fig. 1: Illustration of a third-order PARAFAC tensor as a sum of R rank-one tensors.

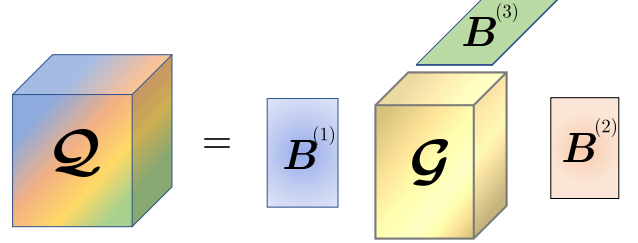


Fig. 2: Illustration of a third-order Tucker tensor and its factor matrices and core tensor.

C. Tucker Decomposition

The Tucker decomposition expresses a tensor as a set of factor matrices and a core tensor. A P -th order tensor $\mathcal{Q} \in \mathbb{C}^{I_1 \times \dots \times I_P}$ that admits a Tucker decomposition, can be written as

$$\mathcal{Q} = \mathcal{G} \times_1 \mathbf{B}^{(1)} \times_2 \dots \times_P \mathbf{B}^{(P)} \in \mathbb{C}^{I_1 \times \dots \times I_P}, \quad (14)$$

where $\mathbf{B}^{(p)} \in \mathbb{C}^{I_p \times R_p}$ is the p -th factor matrix, for $p = \{1, \dots, P\}$, and $\mathcal{G} \in \mathbb{C}^{R_1 \times \dots \times R_P}$ is the core tensor. The tensor \mathcal{Q} can also be represented as the outer product of its factors, given as

$$\mathcal{Q} = \sum_{r_1=1}^{R_1} \dots \sum_{r_P=1}^{R_P} \mathcal{G}_{r_1, \dots, r_P} \left(\mathbf{b}_{r_1}^{(1)} \circ \dots \circ \mathbf{b}_{r_P}^{(P)} \right),$$

where $\mathbf{b}_{r_p}^{(p)} \in \mathbb{C}^{I_p \times 1}$ is the r_p -th column of the p -th factor matrix $\mathbf{B}^{(p)} \in \mathbb{C}^{I_p \times R_p}$ for $p = \{1, \dots, P\}$ and $r_p = \{1, \dots, R_p\}$. The p -th mode unfolding matrix of \mathcal{Q} , defined as $[\mathcal{Q}]_{(p)} \in \mathbb{C}^{N_p \times N_1 \dots N_{p-1} N_{p+1} \dots N_P}$, is given by

$$[\mathcal{Q}]_{(p)} = \mathbf{B}^{(p)} [\mathcal{G}]_{(p)} \left(\mathbf{B}^{(1)} \otimes \dots \otimes \mathbf{B}^{(p+1)} \otimes \mathbf{B}^{(p-1)} \otimes \dots \otimes \mathbf{B}^{(1)} \right)^T. \quad (15)$$

For $P = 3$, Fig. 2, illustrates the decomposition. Its three mode unfoldings are given by

$$[\mathcal{Q}]_{(1)} = \mathbf{B}^{(1)} [\mathcal{G}]_{(1)} \left(\mathbf{B}^{(3)} \otimes \mathbf{B}^{(2)} \right)^T \in \mathbb{C}^{I_1 \times I_2 I_3}, \quad (16)$$

$$[\mathcal{Q}]_{(2)} = \mathbf{B}^{(2)} [\mathcal{G}]_{(2)} \left(\mathbf{B}^{(3)} \otimes \mathbf{B}^{(1)} \right)^T \in \mathbb{C}^{I_2 \times I_1 I_3}, \quad (17)$$

$$[\mathcal{Q}]_{(3)} = \mathbf{B}^{(3)} [\mathcal{G}]_{(3)} \left(\mathbf{B}^{(2)} \otimes \mathbf{B}^{(1)} \right)^T \in \mathbb{C}^{I_3 \times I_1 I_2}. \quad (18)$$

III. SYSTEM MODEL

We consider the system illustrated in Fig. 3, where the transmitter (TX) is equipped with a uniform linear array (ULA) with M_T antenna elements, the receiver (RX) is equipped with ULA with M_R antenna elements and the IRS

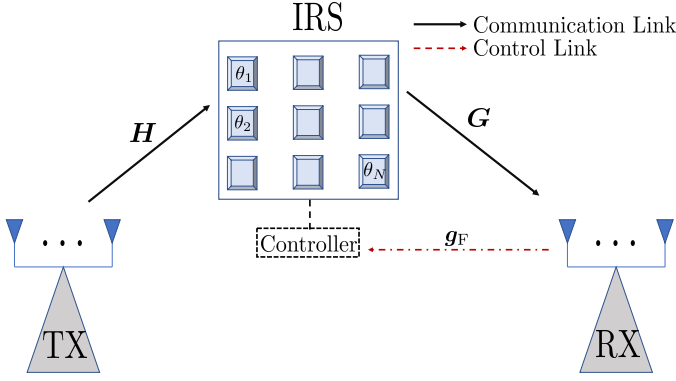


Fig. 3: System model illustration.

has N reflective elements. To simplify the discussion, let us consider a single stream transmission, and assume that there is no direct link between the TX and RX, e.g., base station (BS). First, the TX sends a pilot signal to the RX with the aid of the IRS. Since the IRS has no signal processing capabilities, the channel estimation and the IRS phase-shifts optimization are performed at the RX. The received signal after processing the pilots is given by

$$y = \mathbf{w}^H \mathbf{G} \mathbf{S} \mathbf{H} \mathbf{q} + \mathbf{w}^H \mathbf{b}, \quad (19)$$

where $\mathbf{b} \in \mathbb{C}^{M_R \times 1}$ is the additive noise at the receiver with $\mathbb{E}[\mathbf{b}\mathbf{b}^H] = \sigma_b^2 \mathbf{I}_{M_r}$, $\mathbf{w} \in \mathbb{C}^{M_R \times 1}$ and $\mathbf{q} \in \mathbb{C}^{M_T \times 1}$ are the receiver and transmitter combiner and precoder, respectively. $\mathbf{H} \in \mathbb{C}^{N \times M_T}$ and $\mathbf{G} \in \mathbb{C}^{M_R \times N}$ are the TX-IRS and IRS-RX involved channels, and $\mathbf{S} = \text{diag}(\mathbf{s}) \in \mathbb{C}^{N \times N}$ with $\mathbf{s} = [e^{j\theta_1}, \dots, e^{j\theta_N}] \in \mathbb{C}^{N \times 1}$ being the IRS phase-shift vector, and θ_n is the phase-shift applied to the n -th IRS element.

After the channel estimation step, the precoder and the combiner (active beamformers) vectors \mathbf{w} and \mathbf{q} , and the IRS phase-shift vector \mathbf{s} (passive beamformer) are optimized. Later, the RX needs to feedback to the IRS controller the designed phase-shifts so that the IRS controller tunes the phase-shift for each IRS element. Considering the fact that this feedback occurs in a limited capacity control channel and that the IRS may contain several hundreds to thousands of reflecting elements, the feedback of each phase-shift with a certain resolution imposes a signaling overhead. In this regard, the work [34] models the feedback duration as

$$T_F = \frac{N b_F}{B_F \log \left(1 + \frac{p_F |g_F|^2}{B_F N_0} \right)}, \quad (20)$$

where N is the total number of IRS phase-shifts to be fed back, B_F , p_F are the feedback bandwidth and power, respectively, g_F is the scalar control channel used, b_F is the resolution of each phase-shift, and N_0 is the noise power density. The authors of [34] focus on the problem of rate and EE maximization, where the rate is given by

$$\text{SE} = \left(1 - \frac{T_E + T_F}{T} \right) B \log \left(1 + \frac{p_{\text{TX}} |\mathbf{w}^H \mathbf{G} \mathbf{S} \mathbf{H} \mathbf{q}|^2}{B N_0} \right), \quad (21)$$

with T_E and T being the duration of the channel estimation phase and the total time interval, and B the transmission

bandwidth. The EE is given by $\text{EE} = \text{Rate}/P_{\text{tot}}$, and the total power consumption P_{tot} can be expressed as

$$P_{\text{tot}} = P_E + \frac{T - T_E - T_F}{T} \mu p + \frac{\mu_F p_F T_F}{T} + P_c, \quad (22)$$

where P_E is the power used for the channel estimation phase, $1/\mu$ is the efficiency of the transmitter power amplifier, p_F is the power used during T_F seconds, and μ_F is the efficiency of the transmit amplifier used for feedback. The work [34] maximizes (21) and (22) by optimizing the values of the p , p_F , B , B_F .

Based on the model provided by [34] in (20), we propose to reduce the feedback overhead by factorizing the IRS phase-shift vector into smaller factors, as explained in the following section.

IV. PROPOSED FEEDBACK-AWARE METHOD

In this section, we describe the proposed tensor low-rank approximation based feedback-aware methods that focus on reducing the feedback duration T_F , given in (20). First, we assume that the RX has an estimate of the involved channels \mathbf{H} and \mathbf{G} . The N phase-shifts of the IRS can be determined based on different state-of-the-art algorithms (see, e.g., [6], [34]), and are represented in a vector format as $\mathbf{s} = [e^{j\theta_1}, e^{j\theta_2}, \dots, e^{j\theta_N}] \in \mathbb{C}^{N \times 1}$. Our initial idea consists of factorizing \mathbf{s} as the Kronecker product of P factors, i.e.,

$$\mathbf{s} = \mathbf{s}^{(P)} \otimes \dots \otimes \mathbf{s}^{(1)} \in \mathbb{C}^{N_{P \dots N_1 \times 1}}, \quad (23)$$

where $\mathbf{s}^{(p)} \in \mathbb{C}^{N_p \times 1}$ and $N = \prod_{p=1}^P N_p$.

Example: To get a first insight into the impact of this factorization on the IRS phase-shift feedback overhead, let us consider a simple scenario with $N = 1024$ phase-shifts, and we apply our factorization method by choosing $P = 3$ factors. Consider, as one example, the following factors $\mathbf{s}^{(1)} = [e^{j\theta_1^{(1)}}, \dots, e^{j\theta_{32}^{(1)}}] \in \mathbb{C}^{32 \times 1}$, $\mathbf{s}^{(2)} = [e^{j\theta_1^{(2)}}, \dots, e^{j\theta_8^{(2)}}] \in \mathbb{C}^{8 \times 1}$ and $\mathbf{s}^{(3)} = [e^{j\theta_1^{(3)}}, \dots, e^{j\theta_4^{(3)}}] \in \mathbb{C}^{4 \times 1}$, i.e., $N_1 = 32$, $N_2 = 8$ and $N_3 = 4$. Note that, N_1 , N_2 , N_3 can have every size as long $N_1 \times N_2 \times N_3 = N = 1024$. In this scenario, instead of conveying to the IRS controller 1024 phase-shifts, we only need to convey the phase-shifts of the factors, i.e., $32 + 8 + 4 = 44$, reducing drastically the total amount of phase-shift overhead. Physically, the Kronecker product in (23) represents a summation of the factors phase-shifts. It is clear that, in a general model for a large N , and based on the choice of P , we have that $\sum_{p=1}^P N_p \ll N = \prod_{p=1}^P N_p$. The discussed example is illustrated in Fig. 4. ■

In a general view, the proposed factorization consists of three steps, illustrated in Fig. 5 (for the PARAFAC-IRS model):

- 1) **Rearrangement of elements:** In this step, the optimum phase-shift vector $\mathbf{s} \in \mathbb{C}^{N \times 1}$ is rearranged into a P -th order tensor $\mathcal{S} \in \mathbb{C}^{N_1 \times N_2 \times \dots \times N_P}$, with $N = \prod_{p=1}^P N_p$. This is accomplished by mapping the elements of the IRS phase-shift vector \mathbf{s} into the tensor \mathcal{S} , using the tensorization operator, given in (8).
- 2) **Low-rank approximation (LRA):** In this step, the RX selects an LRA model for \mathcal{S} based on its unfoldings

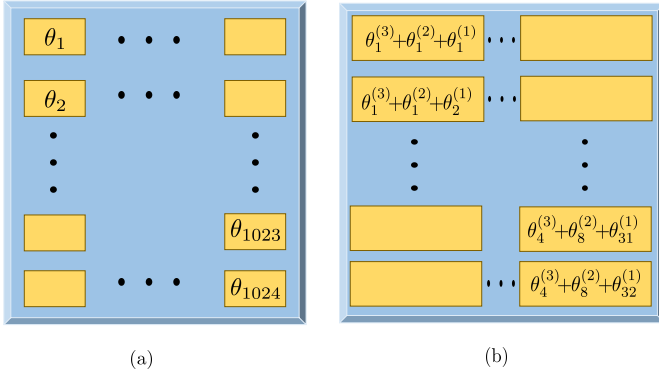


Fig. 4: (a) IRS with $N = 1024$ elements without factorization, (b) IRS with $N = 1024$ elements factorized into $P = 3$ factors.

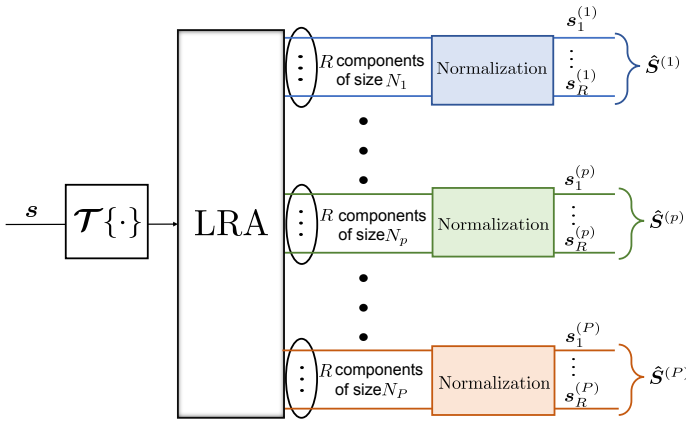


Fig. 5: Proposed method for the IRS phase-shift vector factorization based on the PARAFAC-IRS model.

matrices. For example, the RX can approximate the tensor \mathcal{S} as a PARAFAC or a Tucker model, and makes use of classical tensor algorithms, such as the alternating least squares (ALS) [37] and high order singular value decomposition (HOSVD) [38], to estimate the factors matrices (and the core tensor, in the case of the Tucker model).

- 3) **Normalization:** The factors outputs of the LRA are normalized due to the unitary modulus constraint of the phase-shift vector. In other words, the RX convey to the IRS controller only the angles of the computed factors.

In the next section, we generalize the concept of (23) by factorizing the IRS phase-shift vector $\mathbf{s} \in \mathbb{C}^{N \times 1}$ based on the PARAFAC and Tucker LRA models.

A. PARAFAC-IRS Low-Rank Approximation

In the tensorization step, a rearranging of elements from the IRS phase-shift vector \mathbf{s} to a tensor is performed, i.e., $\mathcal{S} = \mathcal{T}\{\mathbf{s}\}$. Then, the RX will approximate the optimum phase-shift tensor \mathcal{S} using a PARAFAC model, i.e.,

$$\mathcal{S} \approx \sum_{r=1}^R \mathbf{s}_r^{(1)} \circ \dots \circ \mathbf{s}_r^{(P)} \in \mathbb{C}^{N_1 \times \dots \times N_P}, \quad (24)$$

Algorithm 1 PARAFAC-IRS ALS

- 1: **Inputs:** Tensor $\mathcal{S} \in \mathbb{C}^{N_1 \times \dots \times N_P}$, the number of components R , and maximum number of iterations I .
- 2: Randomly initialize the factors $\hat{\mathbf{S}}_0^{(2)}, \dots, \hat{\mathbf{S}}_0^{(P)}$. Iteration $i = 0$.
- 3: Define a maximum number of iteration I .
- 4: **for** $i = 1 : I$ **do**
- 5: **for** $p = 1 : P$ **do**
- 6: Compute an estimate of the p -th factor $\mathbf{S}_i^{(p)}$ as

$$\hat{\mathbf{S}}_i^{(p)} = [\mathcal{S}]_{(p)} \left(\left(\hat{\mathbf{S}}_{i-1}^{(P)} \diamond \dots \diamond \hat{\mathbf{S}}_{i-1}^{(p+1)} \diamond \hat{\mathbf{S}}_{i-1}^{(p-1)} \diamond \dots \diamond \hat{\mathbf{S}}_{i-1}^{(1)} \right)^{\dagger} \right)^+$$
- 7: **for** $r = 1 : R$ **do**
- 8: Normalize the r -th column of $\hat{\mathbf{S}}_i^{(p)}$, defined as $\mathbf{s}_{r,(i)}^{(p)}$, and store its norm as the r -th element of the vector $\boldsymbol{\lambda}^{(p)} \in \mathbb{R}^{R \times 1}$

$$\lambda_r^{(p)} = \left\| \hat{\mathbf{s}}_{r,(i)}^{(p)} \right\|_2, \quad \hat{\mathbf{s}}_{r,(i)}^{(p)} = \frac{\hat{\mathbf{s}}_{r,(i)}^{(p)}}{\lambda_r^{(p)}}.$$
- 9: **end for**
- 10: **end for**
- 11: Define the weighting vector $\boldsymbol{\lambda} = \boldsymbol{\lambda}^{(1)} \odot \dots \odot \boldsymbol{\lambda}^{(P)} \in \mathbb{C}^{R \times 1}$.
- 12: $i = i + 1$
- 13: **end for**
- 14: **Return** $\hat{\mathbf{S}}^{(1)}, \dots, \hat{\mathbf{S}}^{(P)}$ and $\boldsymbol{\lambda}$.

where R is the number of components and $\mathbf{s}_r^{(p)} \in \mathbb{C}^{N_p \times r}$ is the r -th column of the p -th factor matrix $\mathbf{S}^{(p)} = [\mathbf{s}_1^{(p)}, \dots, \mathbf{s}_R^{(p)}] \in \mathbb{C}^{N_p \times R}$, for $p = \{1, \dots, P\}$.

Note that, applying (4) into (23), (24) is a straight-forward generalization where we have R components, and the approximation comes from the fact that an independent phase-shift is fitted as a combination of $P \times R$ sets of phase-shifts, thus an approximation error is expected. However, as it will be explained in Section VI, for scenarios with moderate/strong LOS components (approximated rank-one channels) the effect of the fitting error on the SE performance is negligible.

The RX estimates the factor components solving the following problem

$$\left[\hat{\mathbf{s}}_r^{(1)}, \dots, \hat{\mathbf{s}}_r^{(P)} \right] = \underset{\mathbf{s}_r^{(1)}, \dots, \mathbf{s}_r^{(P)}}{\operatorname{argmin}} \left\| \mathcal{S} - \sum_{r=1}^R \mathbf{s}_r^{(1)} \circ \dots \circ \mathbf{s}_r^{(P)} \right\|_{\text{F}}^2, \quad (25)$$

where $\mathbf{s}_r^{(p)} \in \mathbb{C}^{N_p \times 1}$ is the p -th factor component. Let us define $\mathbf{S}^{(p)} = [\mathbf{s}_1^{(p)}, \dots, \mathbf{s}_R^{(p)}] \in \mathbb{C}^{N_p \times R}$ as the p -th factor matrix, for $p = \{1, \dots, P\}$. From (10), the p -mode unfolding of \mathcal{S} , defined as $[\mathcal{S}]_{(p)} \in \mathbb{C}^{N_p \times N_1 \dots N_{p-1} N_{p+1} \dots N_P}$, is given as

$$[\mathcal{S}]_{(p)} \approx \mathbf{S}^{(p)} \left(\mathbf{S}^{(P)} \diamond \dots \diamond \mathbf{S}^{(p+1)} \diamond \mathbf{S}^{(p-1)} \diamond \dots \diamond \mathbf{S}^{(1)} \right)^{\text{T}}. \quad (26)$$

To solve the problem in (25), the RX can use the ALS algorithm [37], described in Algorithm 1. Basically, the ALS

algorithm contains I iterations, where, in each iteration, P LS problems are solved. The p -th LS problem is defined as

$$\hat{\mathbf{S}}^{(p)} = \underset{\mathbf{S}^{(p)}}{\operatorname{argmin}} \left\| \begin{array}{c} \hat{\mathbf{S}}^{(p)} - \mathbf{S}^{(p)} (\mathbf{S}^{(P)} \diamond \dots \diamond \mathbf{S}^{(p+1)} \diamond \dots \diamond \mathbf{S}^{(1)})^{\top} \\ \mathbf{S}^{(p-1)} \diamond \dots \diamond \mathbf{S}^{(1)} \end{array} \right\|_{\text{F}}^2, \quad (27)$$

where its solution is given by

$$\hat{\mathbf{S}}^{(p)} = [\mathbf{S}]_{(p)} \left(\left(\mathbf{S}^{(P)} \diamond \dots \diamond \mathbf{S}^{(p+1)} \diamond \mathbf{S}^{(p-1)} \diamond \dots \diamond \mathbf{S}^{(1)} \right)^{\top} \right)^+ \quad (28)$$

In the first iteration, the first step is to estimate $[\mathbf{S}]_{(1)}$ based on (28), for $p = 1$. Then, its R columns are normalized to unit norm and stored in as elements of the vector $\boldsymbol{\lambda}^{(1)} \in \mathbb{R}^{R \times 1}$. After the normalization, the estimated $[\hat{\mathbf{S}}]_{(1)}$ is plugged in the LS solution (28) for $p = 2$. Likewise, the columns of the estimated factor $[\hat{\mathbf{S}}]_{(2)}$ are normalized and stored in a vector defined $\boldsymbol{\lambda}^{(2)} \in \mathbb{R}^{R \times 1}$, and then, the normalized estimations $[\hat{\mathbf{S}}]_{(1)}$ and $[\hat{\mathbf{S}}]_{(2)}$ are plugged into the LS solution (28) for $p = 3$. This process continues for the $P - 3$ remaining LS problems. Then, we compute the weighting vector $\boldsymbol{\lambda} \in \mathbb{R}^{R \times 1}$ as the Hadamard product of all P factors norms, i.e., $\boldsymbol{\lambda} = \boldsymbol{\lambda}^{(1)} \odot \boldsymbol{\lambda}^{(2)} \odot \dots \odot \boldsymbol{\lambda}^{(P)}$, finalizing the first iteration of the ALS. Then, the process repeats for all I iterations or until reaching a convergence threshold by checking the normalized mean square error (NMSE) of the reconstructed tensor in a window of consecutive iterations. The NMSE at the i -th iteration is given as

$$e_{(i)} = \frac{\left\| [\mathbf{S}]_{(1),(i)} - [\hat{\mathbf{S}}]_{(1),(i)} \right\|_{\text{F}}^2}{\left\| [\mathbf{S}]_{(1),(i)} \right\|_{\text{F}}^2},$$

where $[\hat{\mathbf{S}}]_{(1),(i)}$ is the reconstructed 1-mode unfolding at the i -th ALS iteration, given by

$$[\hat{\mathbf{S}}]_{(1),(i)} = \hat{\mathbf{S}}^{(1)} \operatorname{diag}(\boldsymbol{\lambda}) \left(\hat{\mathbf{S}}^{(P)} \diamond \dots \diamond \hat{\mathbf{S}}^{(2)} \right)^{\top}. \quad (29)$$

If $|e_{(i)} - e_{(i-1)}| \leq \epsilon$, where ϵ is a pre-defined threshold, the algorithm stops [37]. In this paper, we consider $\epsilon = 10^{-6}$.

After the ALS algorithm, the phase-shifts of each factor and the weighting vector $\boldsymbol{\lambda}$ are quantized to be conveyed to the IRS controller. In this case, the feedback duration is given by

$$T_{\text{F}}^{(\text{PARAFAC})} = \frac{T_{\text{PR}} + R \sum_{p=1}^P N_p \cdot b_{\text{F}}^{(p)} + (R-1) \cdot b_{\text{F}}^{(w)}}{B_{\text{F}} \log \left(1 + \frac{p_{\text{F}} |g_{\text{F}}|^2}{B_{\text{F}} N_0} \right)}, \quad (30)$$

where T_{PR} is the number of bits required for a preamble of the frame, in order to inform the IRS controller the factorization parameters, such as P , R and the quantization bits $b_{\text{F}}^{(p)}$ and $b_{\text{F}}^{(w)}$, where the $b_{\text{F}}^{(p)}$ is the number of bits used for quantize the phase-shifts of the p -th factor, while $b_{\text{F}}^{(w)}$ is the number of bits for quantizing the elements of the weighting vector $\boldsymbol{\lambda}$.

As one example, Fig. 6 illustrates the ratio between state-of-the-art approach, where the N IRS phase-shifts are

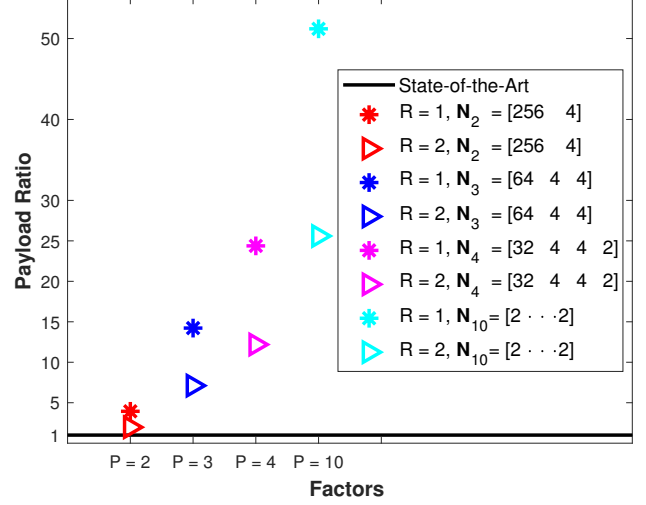


Fig. 6: Feedback payload ratio for $N = 1024$.

fed back, with the proposed PARAFAC-IRS approach, i.e., $N / (R \sum_{p=1}^P N_p)$, not taking into account the preamble T_{PR} and the resolution. Let us define a vector $N_p = [N_1 \dots N_P] \in \mathbb{R}^{P \times 1}$ that contains the factor's size for a certain P . We can observe that, for $R = 1$, the feedback duration of the proposed approach for the case of $P = 2$, $N_1 = 256$, $N_2 = 4$, is almost five times smaller than the state-of-the-art, while, when we increase the number of factors P , the size of the factors can be reduced, thus decreasing the feedback duration, such that, when we have $P = 10$ and $N_p = 2$, for $p = \{1, \dots, P\}$, the feedback overhead of the proposed approach is approximately fifty times smaller than the state-of-the-art [34]. As noticed, with increasing P , the feedback duration of our proposed approach decreases, but also, as it will be discussed in Section VI, the SE in non-line of sight (NLOS) scenarios. Thus, to overcome this loss the RX can increase the number of components R at the cost of a higher feedback overhead. In this way, the proposed overhead-aware method shows off a trade-off between SE and feedback overhead.

B. Tucker-IRS Low-Rank Approximation

Let us consider that the RX node opts to fit the phase-shift tensor \mathcal{S} as a Tucker model, i.e.,

$$\mathcal{S} \approx \sum_{r_1=1}^{R_1} \dots \sum_{r_P=1}^{R_P} \mathcal{G}_{r_1, \dots, r_P} \left(\mathbf{s}_{r_1}^{(1)} \circ \dots \circ \mathbf{s}_{r_P}^{(P)} \right) \in \mathbb{C}^{N_1 \times \dots \times N_P}. \quad (31)$$

where $\mathcal{G} \in \mathbb{C}^{R_1 \times \dots \times R_P}$ is the P -th order core tensor and $\mathbf{s}_{r_p}^{(p)} \in \mathbb{C}^{N_p \times 1}$ is the r_p -th column of the p -th factor matrix $\mathbf{S}^{(p)} \in \mathbb{C}^{N_p \times R_p}$, for $p = \{1, \dots, P\}$ and $r_p = \{1, \dots, R_p\}$. According to (32), the p -th mode unfolding of \mathcal{S} is given by

$$[\mathbf{S}]_{(p)} = \mathbf{S}^{(p)} [\mathcal{G}]_{(p)} \left(\mathbf{S}^{(P)} \otimes \dots \otimes \mathbf{S}^{(p+1)} \otimes \mathbf{S}^{(p-1)} \otimes \dots \otimes \mathbf{S}^{(1)} \right)^{\top} \quad (32)$$

Based on the Tucker-IRS model, the RX estimates each factor matrix $\mathbf{S}^{(p)} \in \mathbb{C}^{N_p \times R}$, for $p = \{1, \dots, P\}$, and the core tensor \mathcal{G} . This estimation procedure can be performed using,

e.g., the HOSVD algorithm [38] given in Algorithm 2, which, in this case, consists of the RX estimating the factors matrices by computing the SVD of all P -mode unfolding matrices of \mathcal{S} independently. Defining the SVD of $[\mathcal{S}]_{(p)}$ as $\mathbf{U}^{(p)}\mathbf{\Sigma}^{(p)}\mathbf{V}^H$, an estimate of $\mathcal{S}^{(p)}$ is given by

$$\hat{\mathcal{S}}^{(p)} = \mathbf{U}_{:1:R_p}^{(p)} \in \mathbb{C}^{N_p \times R_p}, \quad (33)$$

which is the truncation of the left singular matrix $\mathbf{U}^{(p)}$ to its first R_p columns, $p = \{1, \dots, P\}$. The diagonal of the truncated singular matrix $\mathbf{\Sigma}^{(p)}$, defined as $\boldsymbol{\sigma}^{(p)} = \text{diag}(\mathbf{\Sigma}_{1:R_p, 1:R_p}^{(p)}) \in \mathbb{C}^{R_p \times 1}$, is stored to provide the weights to the R_p components in the quantization procedure. Once the P factor matrices are estimated, the RX obtains an estimate of the core tensor \mathcal{G} as

$$\hat{\mathbf{g}} = \left(\hat{\mathcal{S}}^{(P)} \otimes \dots \otimes \hat{\mathcal{S}}^{(1)} \right)^H \mathbf{s}, \mathbf{s} \in \mathbb{C}^{R_1 \dots R_P \times 1}, \quad (34)$$

where $\hat{\mathbf{g}} = \text{vec}(\hat{\mathcal{G}})$ and $\mathbf{s} = \text{vec}(\mathcal{S})$ are the vectorization of the core tensor and the IRS phase-shift tensor, respectively.

The feedback duration of the Tucker-IRS model is given as

$$T_F^{(\text{Tucker})} = \frac{T_{\text{PR}} + \left(\sum_{p=1}^P R_p N_p b_F^{(p)} \right) + \prod_{p=1}^P R_p + b_F^{(w)} \prod_{p=1}^P (R_p - 1)}{B_F \log \left(1 + \frac{P F |g_F|^2}{B_F N_0} \right)}, \quad (35)$$

where T_{PR} is the preamble duration that informs to the IRS controller the chosen LRA model, the number of factors P , and the number of components R_p , for $p = \{1, \dots, P\}$. The term $\sum_{p=1}^P R_p N_p b_F^{(p)}$ represents the cost, in bits, of the conveyed phase-shifts, $\prod_{p=1}^P R_p$ is the cost of the phase-shifts of the core tensor, and $b_F^{(w)} \prod_{p=1}^P (R_p - 1)$ is the term related to the cost of the weighting factors.

V. DISCUSSION ON QUANTIZATION, RECONSTRUCTION AND PARAMETER CHOICES

A. Phase-shift Quantization

After estimating the factors in Algorithms 1 or 2, the RX quantizes the phase-shifts of each factor with $b_F^{(p)}$ bits. Let us define $\tilde{\mathbf{a}} = \mathcal{Q}\{\mathbf{a}, b\}$ as the quantization operation, which quantizes a phase-shift vector \mathbf{a} with b bits. For the PARAFAC-IRS model, we have the following quantized factors $\tilde{\mathbf{s}}_r = \mathcal{Q}\{\hat{\mathbf{s}}_r^{(p)}, b_F^{(p)}\}$ for $p = \{1, \dots, P\}$ and $r = \{1, \dots, R\}$. In addition, for the Tucker-IRS model, we have the following quantized factors and core tensor $\tilde{\mathbf{s}}_{r_p} = \mathcal{Q}\{\hat{\mathbf{s}}_{r_p}^{(p)}, b_F^{(p)}\}$ and $\tilde{\mathcal{G}}_{r_1, \dots, r_P} = \mathcal{Q}\{\hat{\mathcal{G}}_{r_1, \dots, r_P}, b_F^{(p)}\}$, for $p = \{1, \dots, P\}$ and $r_p = \{1, \dots, R_p\}$. For the phase-shift quantization of the p -th factor, we use the following codebook

$$\mathcal{C}_\phi^{(p)} = \left\{ -\pi + \frac{2\pi}{2^{b_F^{(p)}}}, -\pi + \frac{4\pi}{2^{b_F^{(p)}}}, \dots, \pi \right\}.$$

B. Weighting Factor Quantization

For the PARAFAC-IRS model, let us define λ_{\max} as the largest element of $\boldsymbol{\lambda}$. Then, we define a new weighting vector $\boldsymbol{\lambda}' = \boldsymbol{\lambda} / \lambda_{\max} \in \mathbb{R}^{R \times 1}$. Since the largest element of $\boldsymbol{\lambda}'$ is one, we do not need to quantize this element. Hence, we define a new vector $\tilde{\boldsymbol{\lambda}} \in \mathbb{R}^{R-1 \times 1}$ that contains all elements of $\boldsymbol{\lambda}'$, with the exception of the largest one. Then, we quantize the weighting vector by defining $\tilde{\tilde{\boldsymbol{\lambda}}} = \mathcal{Q}\{\tilde{\boldsymbol{\lambda}}, b_F^{(w)}\}$. Finally, we define $\tilde{\boldsymbol{\lambda}} \in \mathbb{R}^{R \times 1}$ as the quantized weighting vector by inserting in the correct position the largest element of $\boldsymbol{\lambda}'$ (one) in $\tilde{\tilde{\boldsymbol{\lambda}}} \in \mathbb{R}^{R-1 \times 1}$. At the end, the weighting vector quantization cost is $(R-1)b_F^{(w)}$ bits.

For the Tucker model, a similar approach is made, with the difference that there are P weighting vectors sorted by their largest value due to the SVD procedure. Considering the p -th weighting vector $\boldsymbol{\sigma}^{(p)} \in \mathbb{R}^{R_p \times 1}$, we normalize it by the first element, yielding $\boldsymbol{\sigma}^{(p)'} = \boldsymbol{\sigma}^{(p)} / \sigma_1^{(p)}$. For the quantization, we define a vector $\tilde{\boldsymbol{\sigma}}^{(p)} \in \mathbb{R}^{R_p-1 \times 1}$ that contains all elements of $\boldsymbol{\sigma}^{(p)'}$ with exception of the first one. Then, we define the quantized p -th weighting factor as $\tilde{\tilde{\boldsymbol{\sigma}}}^{(p)} = \mathcal{Q}\{\tilde{\boldsymbol{\sigma}}^{(p)}, b_F^{(w)}\} \in \mathbb{R}^{R_p-1 \times 1}$. Finally, for the p -th quantized vector, we define the quantized vector $\tilde{\boldsymbol{\sigma}}^{(p)} = [1, \tilde{\tilde{\boldsymbol{\sigma}}}^{(p)'}] \in \mathbb{R}^{R_p \times 1}$. At the end, the quantization of the P weighting factors for the Tucker model costs $b_F^{(w)} \cdot \prod_{p=1}^P (R_p - 1)$ bits.

For both PARAFAC and Tucker models, we define the following amplitude codebook

$$\mathcal{C}_w = \{0.01, 0.01 + l, 0.01 + 2l, \dots, 1\}, \quad (36)$$

where $l = \frac{1-0.01}{2^{b_F^{(w)}} - 1}$ is the pre-defined step. For simplicity, the values of the amplitudes in (36) are rounded to the second decimal point.

C. IRS Phase-shift Vector Reconstruction

After quantization, the RX conveys the factors to the IRS controller. Then, the phase-shift vector is reconstructed as

$$\mathbf{s} = e^{j\angle \hat{\mathbf{s}}} \in \mathbb{C}^{N \times 1}, \quad (37)$$

where $\hat{\mathbf{s}}$ is given by

$$\hat{\mathbf{s}} = \sum_{r=1}^R \tilde{\boldsymbol{\lambda}}_r \left(\tilde{\mathbf{s}}_r^{(P)} \otimes \dots \otimes \tilde{\mathbf{s}}_r^{(1)} \right), \quad (38)$$

for the PARAFAC-IRS model, while for the Tucker-IRS model, $\hat{\mathbf{s}}$ is factorized as

$$\hat{\mathbf{s}} = \sum_{r_1=1}^{R_1} \dots \sum_{r_P=1}^{R_P} \tilde{\mathcal{G}}_{r_1, \dots, r_P} \left(\tilde{\boldsymbol{\sigma}}_{r_P}^{(P)} \tilde{\mathbf{s}}_{r_P}^{(P)} \right) \otimes \dots \otimes \left(\tilde{\boldsymbol{\sigma}}_{r_1}^{(1)} \tilde{\mathbf{s}}_{r_1}^{(1)} \right). \quad (39)$$

D. On the Effect of the Factorization Parameters

In this section, we discuss the choice of the factorization parameters and the system performance implications.

- **Number of factors P :** This parameter defines the total number of factors used in the LRA. Its minimum value for the proposed factorization is $P = 2$, i.e., $P = 1$ means that no factorization is employed, its maximum value is

Algorithm 2 Tucker-IRS HOSVD

1: **Inputs:** Tensor \mathcal{S} , the number of components R_p , for $p = \{1, \dots, P\}$.

2: **for** $p = 1 : P$ **do**

3: Compute the SVD of the p -mode unfolding of \mathcal{S} as

$$[\mathcal{S}]_{(p)} = \mathbf{U}^{(p)} \boldsymbol{\Sigma}^{(p)} \mathbf{V}^{(p)H}.$$

4: Store the diagonal of the truncated singular matrix defined as $\boldsymbol{\sigma}^{(p)} = \text{diag} \left(\boldsymbol{\Sigma}_{1:R_p, 1:R_p}^{(p)} \right) \in \mathbb{C}^{R_p \times 1}$.

5: Set an estimation of $\mathcal{S}^{(p)}$ by truncating the left singular matrix to its first R columns

$$\hat{\mathcal{S}}^{(p)} = \mathbf{U}_{.:1:R_p}^{(p)}.$$

6: **end for**

7: Compute an estimate of the core tensor $\mathbf{g} = \text{vec}(\mathcal{G})$ as

$$\text{vec}(\hat{\mathbf{g}}) = \left(\hat{\mathcal{S}}^{(P)H} \otimes \dots \otimes \hat{\mathcal{S}}^{(1)H} \right) \text{vec}(\mathcal{S}).$$

8: Define $\hat{\mathcal{G}} = \mathcal{T}\{\text{vec}(\hat{\mathbf{g}})\}$.

9: Return $\hat{\mathcal{S}}^{(1)}, \dots, \hat{\mathcal{S}}^{(P)}$ and $\hat{\mathcal{G}}$.

$\log_2(N)$, for the case where all the factors have size $N_p = 2$ for $p = \{1, \dots, P\}$. By increasing the value of P , the number of factors increases, allowing to reduce the size of the factor components N_p . Consequently, increasing P reduces the phase-shift feedback overhead. Nevertheless, by selecting the minimum value of P , the size of each factor component increases, which increases the spectral efficiency at the cost of a higher feedback overhead.

- **Number of components:** For the PARAFAC model, we have R components, while for the Tucker model we have $P \cdot \sum_{p=1}^P R_p$ components. For both models, the number of components is a performance indicator since when increases, the approximation error of the LRA in (24) and (31) decreases. The RX selects its value based on the estimated channels. For example, if the channels have low-rank, or in the presence of a moderate/strong LOS component, the RX may choose $R = 1$. Also, $R = 1$ (PARAFAC model) or $R_p = 1$ (for the Tucker model) are the choices that minimizes the feedback overhead. On the other hand, by increasing R (or R_p), the SE increases at the cost of a higher feedback load.
- **Size of factor components N_p :** The size of the factor components indicates the total number of independent phase-shifts in the proposed solution, which it also affects the performance. For example, for $N = 256$, $P = 2$ and $R = 1$ for the PARAFAC-IRS model, two possible configurations are $(N_1 = 128, N_2 = 2)$ and $(N_1 = N_2 = 16)$. For the first choice, the system has more independent phase-shifts (130), thus a higher SE. However, its feedback overhead is higher than that of the second configuration that requires only 32 phase-shifts to be reported in the feedback channel.

E. On the Effect of the Phase-shift and Weighting Factor Quantization

After the factorization step, the phase-shifts of the factor matrices $\mathcal{S}^{(p)}$ are quantized before being conveyed to the IRS controller. From the fact that the proposed method factorizes the IRS phase-shift vector into P smaller factors, we can select different numbers of bits for the quantization of each factor, unlike the conventional IRS-assisted systems, where the RX (or TX) conveys the N phase-shifts with the same quantization resolution of b_F bits. The proposed method allows the system to adapt the phase-shift resolution of the factors to the available control link capacity, i.e., for each factor we may have a different resolution of $b_F^{(p)}$ in bits, $p = \{1, \dots, P\}$, providing more flexibility to the system design. Regarding the weighting factors, they play a more important role when the number of components $R > 1$ (in the PARAFAC model), or $R_p > 1$, $p = \{1, \dots, P\}$ (in the Tucker model), since they control the importance of the rank-one components in each model.

VI. SIMULATION RESULTS

In this section, we evaluate the performance of the proposed IRS phase-shift overhead-aware feedback model in terms of feedback duration, achievable data rate, SE and EE. The channels in (19) are modeled as

$$\mathbf{H} = \sqrt{\alpha_H \frac{K_H}{K_H + 1}} \mathbf{H}_{\text{LOS}} + \sqrt{\alpha_H \frac{1}{K_H + 1}} \mathbf{H}_{\text{NLOS}}, \quad (40)$$

$$\mathbf{G} = \sqrt{\alpha_G \frac{K_G}{K_G + 1}} \mathbf{G}_{\text{LOS}} + \sqrt{\alpha_G \frac{1}{K_G + 1}} \mathbf{G}_{\text{NLOS}}, \quad (41)$$

where α_H and α_G are the path-loss components of the TX-IRS and IRS-RX links, respectively. The matrices K_H and K_G are the Rician factors associated with \mathbf{H} and \mathbf{G} , respectively. \mathbf{H}_{LOS} , \mathbf{G}_{LOS} follow a geometric-based channel model, while the entries of \mathbf{H}_{NLOS} , \mathbf{G}_{NLOS} are modeled as circularly symmetric complex Gaussian random variables, with zero mean and unit variance, i.e., $\mathbf{H}_{\text{NLOS}} \sim \mathcal{CN}(0, \mathbf{I}_{M_T})$ and $\mathbf{G}_{\text{NLOS}} \sim \mathcal{CN}(0, \mathbf{I}_{M_R})$. More details of (40) and (41) are given in Appendix A.

For a fair comparison between the state-of-the-art [34] and the proposed PARAFAC-IRS and Tucker-IRS models, we optimize the precoder (\mathbf{q}), combiner (\mathbf{w}), and the IRS phase-shifts (\mathbf{s}) using the upper-bound solution of [34]. In this case, they are given as

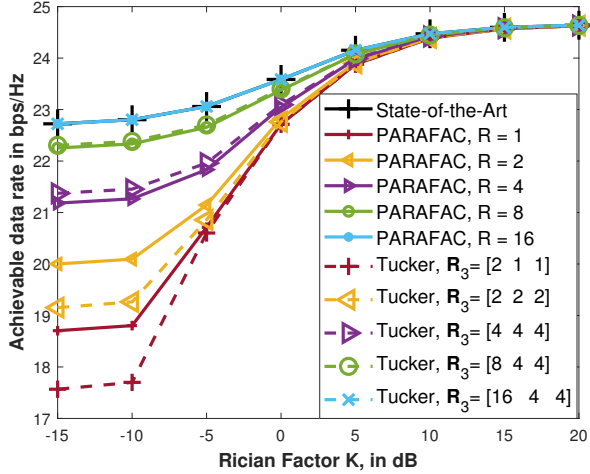
$$\mathbf{w} = \mathbf{U}_{.:1}^{(G)}, \quad \mathbf{q} = \mathbf{V}_{.:1}^{(H)}, \quad \mathbf{s}_n^{(\text{opt})} = e^{-\angle(\mathbf{v}_{n,1}^{(G)} \cdot \mathbf{u}_{n,1}^{(H)})},$$

with $n = \{1, \dots, N\}$, and $\mathbf{U}_{.:1}^{(G)} \in \mathbb{C}^{M_R \times 1}$, $\mathbf{V}_{.:1}^{(G)} \in \mathbb{C}^{N \times 1}$ are the dominant left and right singular vectors of \mathbf{G} , while $\mathbf{U}_{.:1}^{(H)} \in \mathbb{C}^{N \times 1}$, $\mathbf{V}_{.:1}^{(H)} \in \mathbb{C}^{M_T \times 1}$ are the dominant left and right singular vectors of \mathbf{H} .

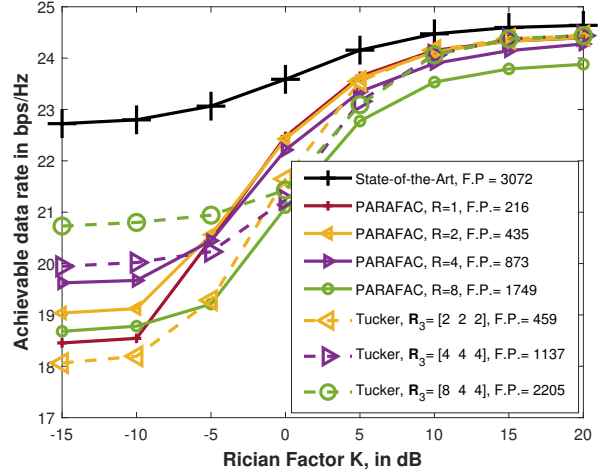
In Figs. 7-9, we set $\alpha_H = \alpha_G = 1$, and consider $K_H = K_G = K$ to simplify the presentation of the figures. However, we have tested the results for a broad range of channel models and parameter settings and observed the same qualitative conclusions as those presented.

A. PARAFAC-IRS vs Tucker-IRS

As a first experiment, we compare, in terms of achievable data rate, the two proposed strategies, PARAFAC-IRS and



(a) No quantization and feedback.



(b) Comparison of the feedback payload (F.P.) in bits.

Fig. 7: Comparison between the state-of-the-art [34], PARAFAC-IRS and Tucker-IRS models with different numbers of components. $N = 1024$, $P = 3$, with $N_1 = 64$, $N_2 = 4$.

Tucker-IRS, with the state-of-the-art method [34], where the IRS phase-shift vector is not factorized. The achievable data rate is given by

$$\log_2 \left(1 + \frac{|w^H G S H q|^2}{\sigma_b^2} \right), \text{ in bits/s/Hz}, \quad (42)$$

where $S = \text{diag}(s^{(\text{opt})}) \in \mathbb{C}^{N \times N}$ is the diagonal matrix containing the optimum IRS phase-shifts, which are given in (38) for the PARAFAC-IRS model and in (39) for the Tucker model.

In Fig. 7, we assume $P = 3$ for the proposed IRS factorization models, with $N_1 = 64$, and $N_2 = N_3 = 4$. As expected, in this scenario the state-of-the-art solution [34], provides the performance upper bound, since no factorization is applied.

In Fig. 7 (a), we compare the models in an ideal scenario with continuous phase-shift and continuous values for the weighting factors. For simplicity, let us define for the Tucker model, the vector $\mathbf{R}_P = [R_1, R_2, \dots, R_P] \in \mathbb{R}^{P \times 1}$ that contains the number of components for each factor for a certain P (with $P = 3$ in this case). As expected, when the number of components R or $\mathbf{R}_{(3)}$ increases, the achievable data rate also increases, and we can observe that for the PARAFAC-IRS model with $R = 16$ and for the Tucker-IRS model with $\mathbf{R}_3 = [16, 4, 4]$, the proposed models achieves the optimum performance of the benchmark method [34].

In practice, both the phase-shift and the weighting factors have to be quantized, as illustrated in Fig. 7 (b), there is an optimal point for the PARAFAC ($R = 4$, for this case) since when $R > 4$ the performance degrades due to overfitting. For the Tucker model, when the number of components of \mathbf{R}_3 increases, the performance in the NLOS region ($K < -5$ dB) also improves at the cost of a higher feedback overhead. Note that, for the moderate/strong LOS scenario ($K > 5$ dB), the number of components for both models does not give a noticeable performance enhancement. In this way, a proper model for a NLOS scenario would be the Tucker one, while PARAFAC is preferable in moderate/strong LOS

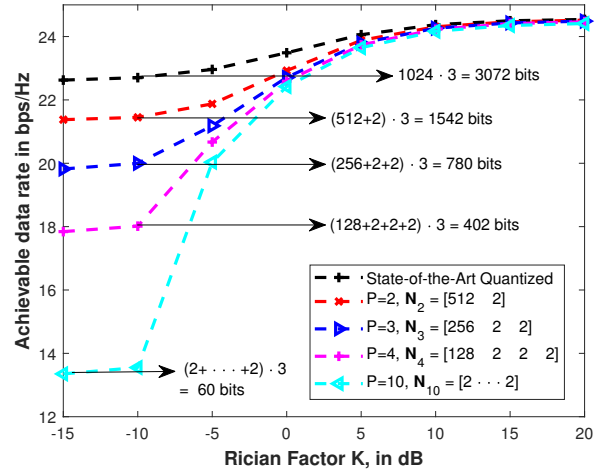


Fig. 8: For an IRS with $N = 1024$, TX and RX with $M_R = M_T = 2$ and $b_F^{(p)} = b_F = 3$ bits, for the IRS phase-shift quantization resolution, for $p = \{1, \dots, P\}$.

cases, since it leads to the best performance with the lowest feedback cost, which can be explained by the fact the channel matrices have low rank and the contributions of the components compared to the largest one are negligible.

In the following, we consider the PARAFAC-IRS model, due its simplicity and lower phase-shift and weight feedback cost. However, we have tested the results for the cases with Tucker method and observed the same qualitative conclusions as those presented.

B. On the Effect of the Number of Factors P

In Fig. 8, we compare the achievable data rate of the PARAFAC-IRS model with $R = 1$, by varying the number of factors. We can observe that, for the NLOS region ($K < -5$ dB), increasing P leads to a degradation on

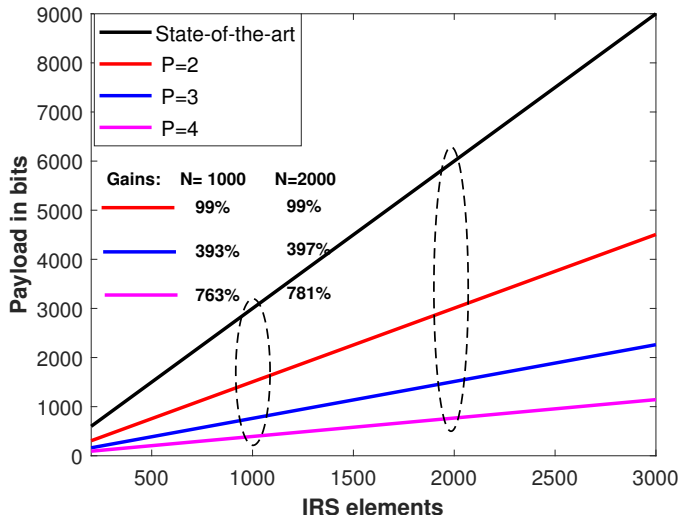


Fig. 9: Feedback payload for the PARAFAC-IRS model with $R = 1$, varying the number of IRS elements.

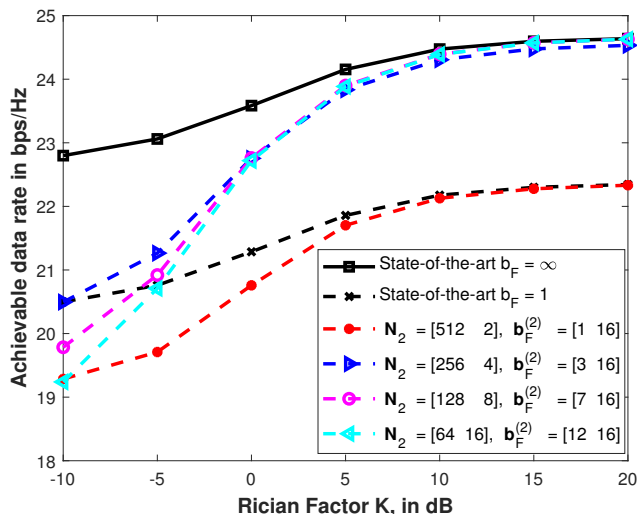


Fig. 10: Performance of the PARAFAC-IRS method by varying the resolution $b_F^{(p)}$ per factor, for fixed control link of 1024 bits.

the performance. This is due to the fact that, for a larger P , we have less independent phase-shifts. For example, for $P = 10$, the phase-shifts of the IRS elements are given by the sum of 10 factorized phase-shifts. However, when the Rician factor K increases, the performance gap between our proposed model and the state-of-the-art [34] reduces. This is explained by the fact that, the IRS phase-shift optimization is based on the channel estimation, thus when K increases, the LOS components become stronger, and we have a better approximation of the PARAFAC-IRS model for $R = 1$. In terms of feedback overhead, when $P = 2$ and for the $K < -5$ dB region, our proposed method has a data rate loss of approximately 1 bps/Hz. However, the feedback overhead is 50% less than that of the benchmark solution [34]. On the other hand, when the scenario changes to $K > 5$ dB, the proper parameter choice is $P = 10$, since this

configuration has a negligible performance loss compared to the state-of-the-art one, while having a lower feedback cost compared to the other proposed configurations ($P = 2, 3, 4$).

Physically, the results illustrated on Fig. 8 can be interpreted as a performance adaptation in the NLOS case, i.e., the RX can properly choose the factorization parameters to meet a required data rate performance or feedback saving. For instance, in this example, by choosing $P = 10$, the RX can afford more often feedback than configurations with smaller values of P .

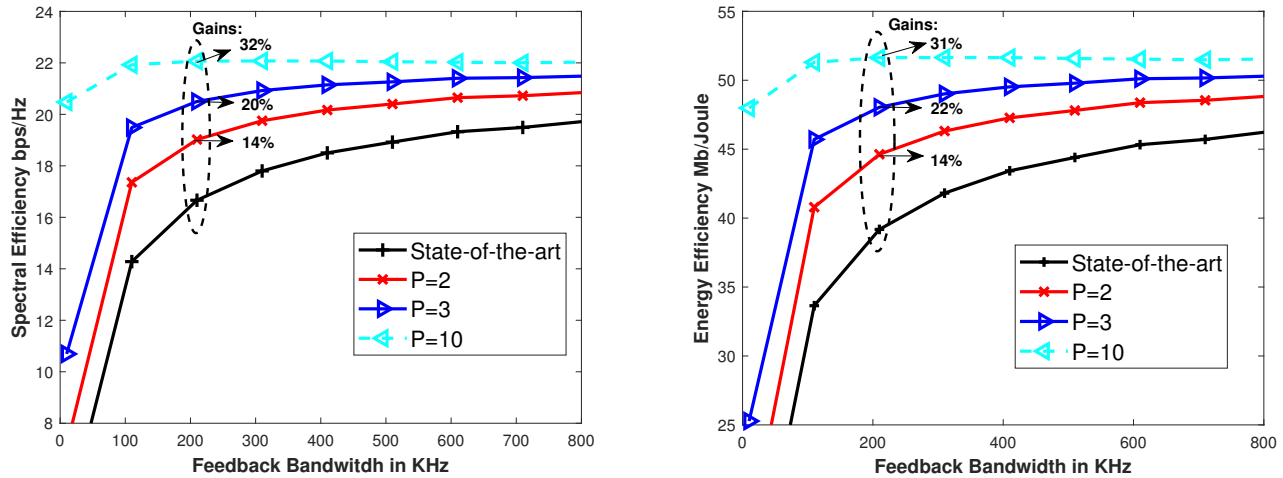
For a better understanding of the merits of the proposed method, Fig. 9 shows the feedback payload in bits by varying the number of IRS elements. As shown, for different methods the payload increases linearly with the number of IRS elements. For a given P , we may have different sets of factor sizes defined by $\mathbf{N}_P = [N_1, \dots, N_P]^T \in \mathbb{R}^{P \times 1}$, where the values of P are set to $P = 2, 3, 4$. We select the size configuration that leads to the better performance, which is the one that has the maximum possible number of independent phase-shifts. For example, assuming $N = 1000$, the size configuration is $\mathbf{N}_2 = [500, 2]$ for $P = 2$, $\mathbf{N}_3 = [250, 2, 2]$ for $P = 3$, and $\mathbf{N}_4 = [125, 2, 2, 2]$ for $P = 4$. Thus, it becomes clear that increasing P drastically reduces the feedback overhead.

C. On the Effect of the Factor Quantization

Here, we evaluate the performance of the proposed method in a limited feedback channel, i.e., we assume that the feedback control link has a maximum capacity of 1024 bits. In this case, traditional quantization applied to the unconstrained IRS phase shift vector (without factorization) is limited to a one bit resolution. We assume this challenging scenario to observe the performance impact of the proposed method when the resolution of the individual factors are adapted. To this end, we assume $N \cdot b_F \geq \mathbf{N}_P \cdot \mathbf{b}_F^{(p)T}$. In Fig. 10, different sets of size configurations for $P = 2$ are evaluated, with different resolutions per factor. The configuration $\mathbf{N}_2 = [512, 2]$ has the worst performance due to the fact that the first factor (512 elements) can only be quantized with 1 bit. However, the size of the factors is reduced, the resolution per factor can be increased accordingly to meet the limited control link capacity limit. For instance, when $\mathbf{N}_2 = [256, 4]$ and $\mathbf{b}_F^{(p)} = [3, 16]$, the total number of bits is $256 \cdot 3 + 4 \cdot 16 = 832$. We can observe that, by increasing the resolution of the factors, the performance gets closer to that of the state-of-the-art phase shift quantization (solid curve). In particular, note that for $K > -5$ dB, our approach provides the best results. Thus, the proposed method can not only reduce the feedback overhead, as illustrated in Figs. 8 and 9, but also it effectively provides higher data rates than traditional quantization over the unconstrained IRS phase shifts, approaching the continuous phase-shift case.

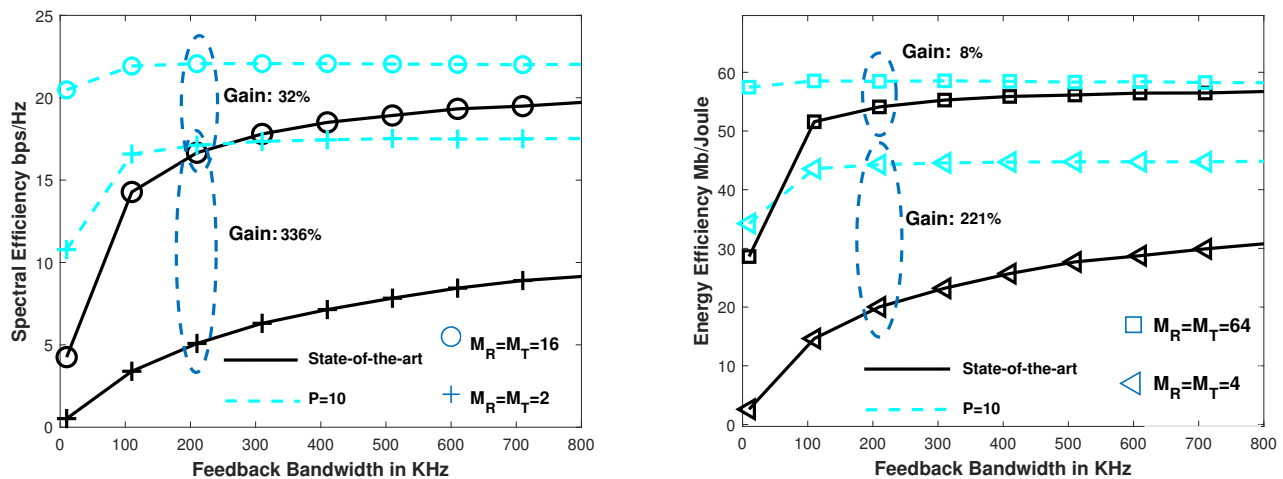
D. Total System SE and EE Evaluation

In this section, we evaluate the performance, in terms of SE and EE, of the proposed method by considering the total system rate, i.e., taking into account the channel estimation procedure duration and the IRS phase-shift feedback duration. To this end, we make use of the expressions given in (21) and (22). The channel estimation period, in (21), is given as $T_E = (M_T N + 1)T_0$, where $T_0 = 0.8\mu$ seconds denotes



(a) SE performance of the state-of-the-art and the proposed method.. (b) EE performance of the state-of-the-art and the proposed method.

Fig. 11: SE and EE performance of the proposed method varying the feedback bandwidth, with $N = 1024$, $M_R = M_T = 16$, $b_F^{(p)} = b_F = 3$ bits, for $p = 2, 3, 10$, for a Rician factor $K = 10$ dB.



(a) SE performance of the state-of-the-art and the proposed method.. (b) EE performance of the state-of-the-art and the proposed method.

Fig. 12: SE and EE performance of the proposed method varying the feedback bandwidth, with $N = 1024$, $b_F^{(p)} = b_F = 4$ bits, for $p = 2, 3, 10$, for a Rician factor $K = 10$ dB.

the duration of the pilot tones [34]. The frame duration is given by $T = T_{PD} + T_F$, where $T_{PD} = T_E + T_D$, is divided into 30% for pilot transmissions (T_E) and 70% for data transmission T_D . Regarding the power parameters of (22), we have $P_E = P_0(1 + NM_T)T_0$, where $P_0 = 0.8$ mW is the pilot tone power. Other parameter definitions can be found in Table I. The feedback channel g_F is generated from a circular symmetric complex Gaussian distribution, normalized by $\sqrt{\beta_F} = \sqrt{\alpha_H} = \sqrt{\alpha_G}$ to account for the effects of pathloss and shadowing, as given in Table I. In our next experiments, we assume $K = 10$ dB, $N = 1024$. For the proposed method, we consider the PARAFAC-IRS model with $R = 1$. As for the number of factors, we study three configurations, with $P = 2$, ($N_2 = [512, 2]$), $P = 3$ ($N_3 = [256, 2, 2]$) and $P = 10$ ($N_{10} = [2, \dots, 2] \in \mathbb{R}^{10 \times 1}$).

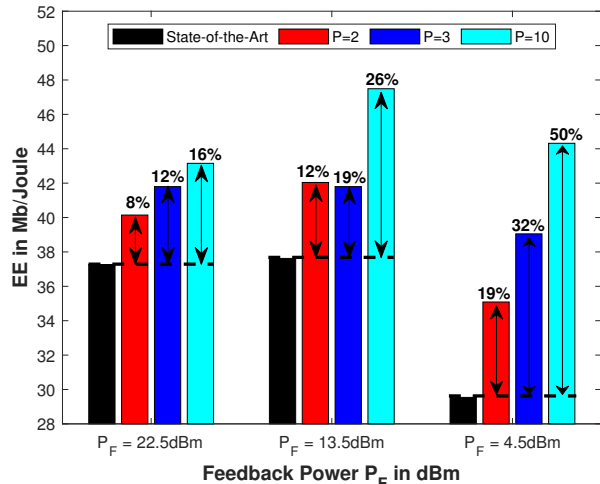
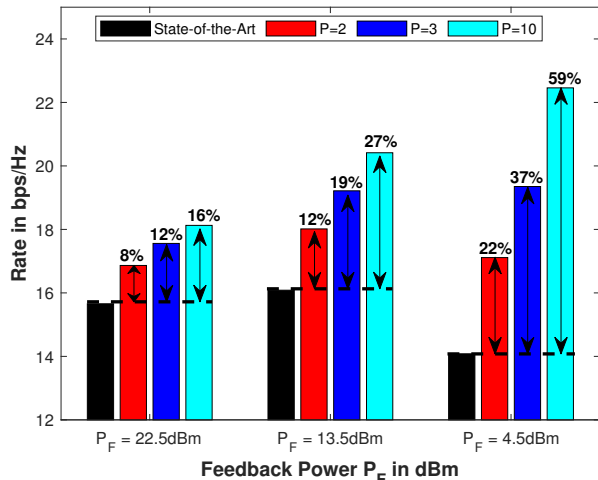
Figs.11 and 12, we analyze the total SE and EE of the

$P_{\max}/P_{c,0}/P_{c,n}$	B_{\max}	N_0	α_H/α_G	μ/μ_F
45/45/10 dBm	100 MHz	-174 dBm/Hz	110/110 dB	1/1

TABLE I

proposed method with the state-of-the-art [34], by varying the feedback bandwidth $B_F = B_{\max} - B$, where B_{\max} is the total available bandwidth given in Table I. As shown, in Figs. 11a and 11b, when the number of factors increases the feedback duration reduction pays off in the total system SE and EE. The proposed method achieves a gain in the SE of 32% for $P = 10$, 20% for $P = 3$, and 14% for $P = 2$, over the state-of-the-art, considering the $B_F = 200$ kHz, with a similar gain in the EE.

In Figs. 12a and 12b we compare the proposed PARAFAC-IRS model under $R = 1$ and $P = 10$ with the



(a) SE performance of the state-of-the-art and the proposed method. (b) EE performance of the state-of-the-art and the proposed method.

Fig. 13: EE and SE performance of the proposed method varying the feedback power, with $N = 1024$, $M_R = M_T = 2$, $b_F^{(p)} = b_F = 3$ bits, for $p = 2, 3, 10$, Rician factor $K = 10$ dB.

state-of-the-art by varying the number of antennas. In this case, we observe that, for a feedback bandwidth $B_F \leq 200$ kHz, the proposed factorization with a 2×2 setup outperforms (in terms of SE and EE) the state-of-the-art one under the 4×4 setup, while presenting the same performance than the state-of-the-art one under the 16×16 setup. Finally, Figs. 13a and 13b show the SE and EE performances of the proposed method as a function of the feedback power p_F , with $p_F = P_{\max} - p_{TX}$. We notice that the proposed configurations provide the best results in all scenarios.

To summarize the results illustrated in Figs.11-13, we conclude that the proposed tensor-based LRA IRS phase-shift factorization models allows to reduce the number of phase-shifts to be conveyed to the IRS-controller, which significantly reduces the feedback overhead, resulting in SE and EE performance enhancements. In addition, our approach reaches similar performance to the non-factorized IRS, especially in moderate/strong LOS scenarios, as it can be seen in Figs. 7-10. From a system-level viewpoint, the network can resort to the proposed overhead-aware IRS model to increase the feedback periodicity, i.e., by providing more frequent feedback, which is crucial in fast time-varying channels, where the IRS should be reconfigured more frequently to follow the environment changes. Moreover, the proposed IRS factorization methods allow the network to multiplex more IRS phase-shifts in the same feedback channel, which is useful to accommodate multi-user IRS-assisted communications.

VII. CONCLUSIONS AND PERSPECTIVES

In this paper, we proposed two IRS phase-shift feedback overhead-aware methods based on tensor signal processing, namely, PARAFAC-IRS and Tucker-IRS. We showed that the proposed methods significantly reduce the IRS phase-shift feedback overhead, compared to the state-of-the-art approach, where the IRS phase shifts are not factorized. The PARAFAC-IRS method is preferable in the case of moderate/strong LOS scenarios, achieving a spectral efficiency

that is close to that of the state-of-the-art, while providing a feedback overhead reduction. Moreover, in NLOS scenarios, the Tucker-IRS model achieves a higher data rate than the PARAFAC-IRS model at the expense of a higher feedback overhead. By controlling the factorization parameters, we showed how to trade off data rate for feedback-overhead, allowing the network controller to adapt the IRS factorization parameters to meet a determined quality of service.

APPENDIX A CHANNEL MODEL

We provide details on the channel models for \mathbf{H} and \mathbf{G} , given in (40) and (41), respectively. As mentioned, the NLOS components of \mathbf{H} and \mathbf{G} are modeled as random channels with $\mathbb{E}[\mathbf{H}_{\text{NLOS}}^H \mathbf{H}_{\text{NLOS}}] = \mathbf{I}_{M_T}$ and $\mathbb{E}[\mathbf{G}_{\text{NLOS}} \mathbf{G}_{\text{NLOS}}^H] = \mathbf{I}_{M_R}$. Nonetheless, the LOS components are given as

$$\begin{aligned} \mathbf{H}_{\text{LOS}} &= \alpha_H \mathbf{b}_{\text{IRS}} \cdot \mathbf{a}_{\text{TX}}^H \in \mathbb{C}^{N \times M_T}, \\ \mathbf{G}_{\text{LOS}} &= \alpha_G \mathbf{b}_{\text{RX}} \cdot \mathbf{a}_{\text{IRS}}^H \in \mathbb{C}^{M_R \times N}, \end{aligned}$$

where α_H and α_G are the path-loss components of the TX-IRS and IRS-RX links, respectively. Assuming that the TX and the RX are equipped with ULAs with half-wavelength inter-element spacing, their steering vectors can be written as

$$\mathbf{a}_{\text{TX}} = \left[1, e^{j\pi \sin \theta_{\text{TX}}}, \dots, e^{j\pi(M_T-1)\sin \theta_{\text{TX}}} \right]^T \in \mathbb{C}^{M_T \times 1}, \quad (43)$$

$$\mathbf{b}_{\text{RX}} = \left[1, e^{j\pi \sin \theta_{\text{RX}}}, \dots, e^{j\pi(M_R-1)\sin \theta_{\text{RX}}} \right]^T \in \mathbb{C}^{M_R \times 1}, \quad (44)$$

where θ_{TX} and θ_{RX} are the TX and RX angle of departure (AOD) and angle of arrival (AOA), respectively, which are generated from a uniform random distribution with $\{\theta_{\text{TX}}, \theta_{\text{RX}}\} \in [-\pi, \pi]$. Since the IRS is a 2-D panel, the steering vectors associated with arrival and departure angles can be factorized as the Kronecker product of horizontal and vertical component vectors, respectively, as follows:

$$\mathbf{b}_{\text{IRS}} = \mathbf{b}_{\text{IRS}}^{(v)} \otimes \mathbf{b}_{\text{IRS}}^{(h)} \in \mathbb{C}^{N_h N_v \times 1}, \quad (45)$$

$$\mathbf{a}_{\text{IRS}} = \mathbf{a}_{\text{IRS}}^{(v)} \otimes \mathbf{a}_{\text{IRS}}^{(h)} \in \mathbb{C}^{N_h N_v \times 1}, \quad (46)$$

where $N = N_h N_v$, $\mathbf{b}_{\text{IRS}}^{(h)} \in \mathbb{C}^{N_h \times 1}$ and $\mathbf{b}_{\text{IRS}}^{(v)} \in \mathbb{C}^{N_v \times 1}$ are the AOA steering vectors in the azimuth and elevation directions, respectively. Likewise, $\mathbf{a}_{\text{IRS}}^{(h)} \in \mathbb{C}^{N_h \times 1}$ and $\mathbf{a}_{\text{IRS}}^{(v)} \in \mathbb{C}^{N_v \times 1}$ are the AOD steering vectors in the azimuth and elevation directions, respectively.

$$\begin{aligned}\mathbf{b}_{\text{IRS}}^{(h)} &= [1, e^{j\pi \sin \psi_{\text{IRS}}^{\text{AOA}} \cos \phi_{\text{IRS}}^{\text{AOA}}}, \dots, e^{j\pi(N_h-1) \sin \psi_{\text{IRS}}^{\text{AOA}} \cos \phi_{\text{IRS}}^{\text{AOA}}}, \\ \mathbf{b}_{\text{IRS}}^{(v)} &= [1, e^{j\pi \cos \phi_{\text{IRS}}^{\text{AOA}}}, \dots, e^{j\pi(N_h-1) \cos \phi_{\text{IRS}}^{\text{AOA}}}, \\ \mathbf{a}_{\text{IRS}}^{(h)} &= [1, e^{j\pi \sin \psi_{\text{IRS}}^{\text{AOD}} \cos \phi_{\text{IRS}}^{\text{AOD}}}, \dots, e^{j\pi(N_h-1) \sin \psi_{\text{IRS}}^{\text{AOD}} \cos \phi_{\text{IRS}}^{\text{AOD}}}, \\ \mathbf{a}_{\text{IRS}}^{(v)} &= [1, e^{j\pi \cos \phi_{\text{IRS}}^{\text{AOD}}}, \dots, e^{j\pi(N_h-1) \cos \phi_{\text{IRS}}^{\text{AOD}}},\end{aligned}$$

where $\phi_{\text{IRS}}^{\text{AOA}}$ and $\phi_{\text{IRS}}^{\text{AOD}}$ are the elevation angles of arrival and departure, while $\psi_{\text{IRS}}^{\text{AOA}}$ and $\psi_{\text{IRS}}^{\text{AOD}}$ are the azimuth angles of arrival and departure. The azimuth angles $\psi_{\text{IRS}}^{\text{AOA}}$ and $\psi_{\text{IRS}}^{\text{AOD}}$ are generated from a uniform random distribution with $\{\psi_{\text{IRS}}^{\text{AOA}}, \psi_{\text{IRS}}^{\text{AOD}}\} \in [-\pi, \pi]$, while the elevation angles $\phi_{\text{IRS}}^{\text{AOA}}$ and $\phi_{\text{IRS}}^{\text{AOD}}$ are generated from a uniform random distribution with $\{\phi_{\text{IRS}}^{\text{AOA}}, \phi_{\text{IRS}}^{\text{AOD}}\} \in [0, \pi/2]$.

REFERENCES

- [1] B. Sokal, P. R. B. Gomes, A. L. F. de Almeida, B. Makki, and G. Fodor, "IRS phase-shift feedback overhead-aware model based on rank-one tensor approximation," 2022, [Online]. Available: <https://arxiv.org/pdf/2205.12024v1.pdf>.
- [2] Q. Wu, S. Zhang, B. Zheng, C. You, and R. Zhang, "Intelligent reflecting surface aided wireless communications: A tutorial," *IEEE Trans. Commun.*, pp. 1–1, May 2021.
- [3] S. Gong, X. Lu, D. T. Hoang, D. Niyato, L. Shu, D. I. Kim, and Y.-C. Liang, "Toward smart wireless communications via intelligent reflecting surfaces: A contemporary survey," *IEEE Commun. Surveys Tuts.*, vol. 22, no. 4, pp. 2283–2314, June 2020.
- [4] M. Jian, G. C. Alexandropoulos, E. Basar, C. Huang, R. Liu, Y. Liu, and C. Yuen, "Reconfigurable intelligent surfaces for wireless communications: Overview of hardware designs, channel models, and estimation techniques," *arXiv preprint arXiv:2203.03176*, 2022.
- [5] M. Di Renzo, A. Zappone, M. Debbah, M.-S. Alouini, C. Yuen, J. De Rosny, and S. Tretyakov, "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, July 2020.
- [6] E. Basar, M. Di Renzo, J. De Rosny, M. Debbah, M.-S. Alouini, and R. Zhang, "Wireless communications through reconfigurable intelligent surfaces," *IEEE Access*, vol. 7, pp. 116753–116773, Aug. 2019.
- [7] E. Basar, "Reconfigurable intelligent surface-based index modulation: A new beyond MIMO paradigm for 6G," *IEEE Trans. Commun.*, vol. 68, no. 5, pp. 3187–3196, Feb. 2020.
- [8] Ö. Özdoğan, E. Björnson, and E. G. Larsson, "Intelligent reflecting surfaces: Physics, propagation, and pathloss modeling," *IEEE Wireless Commun. Lett.*, vol. 9, no. 5, pp. 581–585, Dec. 2019.
- [9] N. Rajatheva, I. Atzeni, S. Bicaïs, E. Björnson, A. Bourdoux, S. Buzzi, C. D'Andrea, J.-B. Dore, S. Erkucuk, M. Fuentes, et al., "Scoring the terabit/s goal: Broadband connectivity in 6G," *arXiv preprint arXiv:2008.07220*, 2020.
- [10] A. Taha, M. Alrabeiah, and A. Alkhateeb, "Enabling large intelligent surfaces with compressive sensing and deep learning," pp. 44304–44321, March 2021.
- [11] M. H. Khoshafa, T. M. Ngatched, M. H. Ahmed, and A. R. Ndjiongue, "Active reconfigurable intelligent surfaces-aided wireless communication system," *IEEE Wireless Commun. Lett.*, vol. 25, no. 11, pp. 3699–3703, Sept. 2021.
- [12] G. C. Alexandropoulos and E. Vlachos, "A hardware architecture for reconfigurable intelligent surfaces with minimal active elements for explicit channel estimation," in *Proc. in ICASSP 2020*. Barcelona, Spain: IEEE, May 2020, pp. 9175–9179.
- [13] 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, p. 4157–4170, June 2019.
- [14] E. Björnson, Ö. Özdoğan, and E. G. Larsson, "Intelligent reflecting surface versus decode-and-forward: How large surfaces are needed to beat relaying?" *IEEE Wireless Commun. Lett.*, vol. 9, no. 2, pp. 244–248, Oct. 2019.
- [15] H. Guo, B. Makki, M. Åström, M.-S. Alouini, and T. Svensson, "Dynamic blockage pre-avoidance using reconfigurable intelligent surfaces," *arXiv preprint arXiv:2201.06659*, 2022.
- [16] G. T. de Araújo, A. L. F. de Almeida, and R. Boyer, "Channel estimation for intelligent reflecting surface assisted MIMO systems: A tensor modeling approach," *IEEE J. Sel. Topics Signal Process.*, vol. 15, no. 3, pp. 789–802, Feb. 2021.
- [17] G. Tavares de Araújo, P. R. Brboza Gomes, A. Lima Férrer de Almeida, G. Fodor, and B. Makki, "Semi-blind joint channel and symbol estimation in ired-assisted multi-user MIMO networks," *arXiv e-prints*, pp. arXiv:2202.2022, 2022.
- [18] J. Chen, Y.-C. Liang, H. V. Cheng, and W. Yu, "Channel estimation for reconfigurable intelligent surface aided multi-user MIMO systems," 2019, [Online]. Available: <https://arxiv.org/pdf/1912.03619.pdf>.
- [19] C. Hu, L. Dai, S. Han, and X. Wang, "Two-timescale channel estimation for reconfigurable intelligent surface aided wireless communications," *IEEE Trans. Commun.*, vol. 69, no. 11, pp. 7736–7747, April 2021.
- [20] B. Li, Z. Zhang, Z. Hu, and Y. Chen, "Joint array diagnosis and channel estimation for RIS-aided mmwave MIMO system," *IEEE Access*, vol. 8, pp. 193992–194006, 2020.
- [21] K. Ardah, S. Gherekhloo, A. L. 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.
- [22] J. An, C. Xu, L. Gan, and L. Hanzo, "Low-complexity channel estimation and passive beamforming for RIS-assisted MIMO systems relying on discrete phase shifts," *IEEE Trans. Commun.*, Nov. 2021.
- [23] L. Wei, C. Huang, G. C. Alexandropoulos, C. Yuen, Z. Zhang, and M. Debbah, "Channel estimation for RIS-empowered multi-user MISO wireless communications," *IEEE Trans. Commun.*, pp. 1–1, March 2021.
- [24] Y. Yang, B. Zheng, S. Zhang, and R. Zhang, "Intelligent reflecting surface meets OFDM: Protocol design and rate maximization," *IEEE Trans. Commun.*, March 2020.
- [25] N. S. Perović, L.-N. Tran, M. Di Renzo, and M. F. Flanagan, "Achievable rate optimization for MIMO systems with reconfigurable intelligent surfaces," *IEEE Trans. Wireless Commun.*, pp. 1–1, Feb. 2021.
- [26] Z. Gao, Y. Xu, Q. Wang, Q. Wu, and D. Li, "Outage-constrained energy efficiency maximization for RIS-assisted WPCNs," *IEEE Trans. Wireless Commun.*, vol. 25, no. 10, pp. 3370–3374, July 2021.
- [27] L. Du, W. Zhang, J. Ma, and Y. Tang, "Reconfigurable intelligent surfaces for energy efficiency in multicast transmissions," *IEEE Trans. Veh. Technol.*, vol. 70, no. 6, pp. 6266–6271, May 2021.
- [28] L. You, J. Xiong, D. W. K. Ng, C. Yuen, W. Wang, and X. Gao, "Energy efficiency and spectral efficiency trade off in RIS-aided multiuser MIMO uplink transmission," *IEEE Trans. Signal Process.*, vol. 69, pp. 1407–1421, Dec. 2020.
- [29] J. Chen, Y. Xie, X. Mu, J. Jia, Y. Liu, and X. Wang, "Energy efficient resource allocation for IRS assisted CoMP systems," *IEEE Trans. Commun.*, Jan. 2022.
- [30] M. Z. Siddiqi, R. Mackenzie, M. Hao, and T. Mir, "On energy efficiency of wideband RIS-aided cell-free network," *IEEE Access*, vol. 10, pp. 19742–19752, Feb. 2022.
- [31] A. Khaleel and E. Basar, "A novel NOMA solution with RIS partitioning," *IEEE J. Sel. Topics Signal Process*, Nov. 2021.
- [32] K. Ntougias and I. Krikidis, "Interference-constrained IRS-aided SWIPT," in *Proc. IEEE SPAWC 2021*. Lucca, Italy: IEEE, Sept. 2021, pp. 116–120.
- [33] M. A. ElMossallamy, K. G. Seddik, W. Chen, L. Wang, G. Y. Li, and Z. Han, "RIS optimization on the complex circle manifold for interference mitigation in interference channels," *IEEE Trans. Veh. Technol.*, vol. 70, no. 6, pp. 6184–6189, April 2021.
- [34] A. Zappone, M. Di Renzo, F. Shams, X. Qian, and M. Debbah, "Overhead-aware design of reconfigurable intelligent surfaces in smart radio environments," *IEEE Trans. Wireless Commun.*, vol. 20, no. 1, pp. 126–141, Sept. 2020.
- [35] R. A. Harshman et al., "Foundations of the parafac procedure: Models and conditions for an explanatory multimodal factor analysis," *UCLA working paper in Phonetics*, vol. 16, pp. 1–84, 1970.
- [36] L. R. Tucker, "Some Mathematical Notes on Three-mode Factor Analysis," *Psychometrika*, vol. 31, no. 3, pp. 279–311, Sept. 1966.
- [37] T. G. Kolda and B. W. Bader, "Tensor decompositions and applications," *SIAM review*, vol. 51, no. 3, pp. 455–500, Aug. 2009.
- [38] L. De Lathauwer, B. De Moor, and J. Vandewalle, "On the best rank-1 and rank-(r1, r2, ..., rn) approximation of higher-order tensors," *SIAM journal on Matrix Analysis and Applications*, vol. 21, no. 4, pp. 1324–1342, March 2000.