

Reducing Channel Estimation and Feedback Overhead in IRS-Aided Downlink System: A Quantize-then-Estimate Approach

Rui Wang, Zhaorui Wang, Liang Liu, Shuowen Zhang, and Shi Jin

Abstract—Channel state information (CSI) acquisition is essential for the base station (BS) to fully reap the beamforming gain in intelligent reflecting surface (IRS)-aided downlink communication systems. Recently, [1] revealed a strong correlation in different users' cascaded channels stemming from their common BS-IRS channel component, and leveraged such a correlation to significantly reduce the pilot transmission overhead in IRS-aided uplink communication. In this paper, we aim to exploit the above channel property to reduce the overhead for both pilot and feedback transmission in IRS-aided downlink communication. Note that in the downlink, the distributed users merely receive the pilot signals containing their own CSI and cannot leverage the correlation in different users' channels, which is in sharp contrast to the uplink counterpart considered in [1]. To tackle this challenge, this paper proposes a novel "quantize-then-estimate" protocol in frequency division duplex (FDD) IRS-aided downlink communication. Specifically, the users quantize and feed back their received pilot signals, instead of the estimated channels, to the BS. After de-quantizing the pilot signals received by all the users, the BS estimates all the cascaded channels by leveraging their correlation, similar to the uplink scenario. Under this protocol, we manage to propose efficient user-side quantization and BS-side channel estimation methods. Moreover, we analytically quantify the pilot and feedback transmission overhead to reveal the significant performance gain of our proposed scheme over the conventional "estimate-then-quantize" scheme.

Index Terms—Intelligent reflecting surface (IRS), channel estimation, channel feedback, distributed source coding.

I. INTRODUCTION

A. Motivation

Intelligent reflecting surface (IRS) has been recognized as a promising technique to enhance the capacity and coverage

Manuscript received March 16, 2024; revised July 15, 2024 and October 19, 2024; accepted November 10, 2024. This work was supported in part by the National Key Research and Development Project of China under Grant 2022YFB2902800; in part by the National Natural Science Foundation of China under Grant 62101474; in part by the Research Grants Council, Hong Kong, China, under Grant 15203222 and 15230022; in part by the Basic Research Project under Grant HZQB-KCZY-2021067 of Hetao Shenzhen-Hong Kong Science and Technology Cooperation Zone. (corresponding author: Liang Liu)

R. Wang, L. Liu, and S. Zhang are with the Department of Electrical and Electronic Engineering, The Hong Kong Polytechnic University, Hong Kong SAR, China (e-mails: rui-eie.wang@connect.polyu.hk, {liang-eie.liu, shuowen.zhang}@polyu.edu.hk).

Z. Wang is with the FNii and SSE, The Chinese University of Hong Kong, Shenzhen, Shenzhen 518172, China (e-mail: wang-wang2020channel@cuhk.edu.cn).

S. Jin is with the National Mobile Communications Research Laboratory, Southeast University, Nanjing 210096, China (e-mail: jinshi@seu.edu.cn).

The materials in this paper have been presented in part at the IEEE Global Communications Conference, December 2023 [2].

of the future 6G cellular networks, thanks to its ability to tune the channels to be favorable for communication. To design the best propagation conditions via the IRS, channel state information (CSI) acquisition is of paramount importance. However, such a task is challenging due to the vast number of channel coefficients associated with the IRS [3]–[6].

This paper considers IRS-assisted downlink communication in a frequency division duplex (FDD) system, where a multi-antenna base station (BS) needs to know the BS-IRS-user cascaded channels of all the users for designing its own and the IRS's beamforming vectors. In time division duplex (TDD) systems, the downlink CSI can be acquired by estimating the uplink CSI thanks to channel reciprocity. In our considered FDD systems, the channel reciprocity does not hold and the users have to feed back some useful information from their received pilot signals to let the BS acquire the CSI. Under the conventional systems without the IRS, the "estimate-then-quantize" scheme [7]–[10] was widely used for downlink channel estimation and feedback, where each user first estimates its downlink channels with the BS based on its received pilot signals and then feeds back the estimated channels to the BS. However, the overall overhead of this classic protocol is unaffordable in IRS-aided systems due to the following reasons. First, to enable each user to estimate a huge number of coefficients in its BS-IRS-user channel, the BS has to send a long pilot sequence, leading to high channel estimation overhead. Second, after the channel estimation phase, each user has to feed back a huge number of quantization bits for transmitting the estimated channel coefficients to the BS, leading to high feedback overhead. This calls for some innovative protocol to replace the conventional "estimate-then-quantize" protocol in IRS-assisted downlink systems for low-overhead communication.

B. Prior Work

Previous research has been conducted to reduce channel estimation overhead in IRS-assisted uplink communication systems, by utilizing multi-user channel correlation [1], [11], two-timescale property [12], beamspace channel sparsity [11], [13]–[17], IRS elements grouping [18], etc. In TDD downlink systems, the above methods can be used for the BS to obtain the CSI based on channel reciprocity. However, in our considered FDD downlink systems, dedicated methods should be proposed for low-overhead channel estimation and feedback.

Under FDD IRS-assisted downlink communication systems, most prior works are under the “estimate-then-quantize” scheme. In particular, they are interested in reducing the overhead in the feedback phase, assuming the channels are already estimated by the users [19]–[24]. For the single-user system, [19] proposed a novel cascaded codebook for the BS-IRS and IRS-user subchannel, respectively, that is synthesized by two sub-codebooks whose codewords are cascaded by line-of-sight and non-line-of-sight components. It is demonstrated that the cascaded codebook outperforms the naive random vector quantization codebook with lower feedback bits requirement. [20] reduced feedback overhead by selecting several dominant BS-IRS-user cascaded channel paths based on their contributions to spectral efficiency, instead of feeding back CSI of all the cascaded paths. However, these two works only focused on channel path gain information feedback, without considering path angle information. [21] proposed to customize a cascaded channel with a reduced number of paths for the sub-6 GHz rich scattering environment by path selection and multi-IRS phase shifter design, thus reducing the number of feedback parameters. [22] leveraged the two-timescale property of the cascaded BS-IRS-user channels [12] to build a neural network consisting of large and small timescale feedback. The BS-IRS channel, which is assumed to be unchanged for a long time, only needs to be fed back in each large timescale, and the IRS-user channel, which changes frequently but with low dimension, is fed back in each small timescale. For the multi-user system, [23] exploited the single-structured sparsity of the cascaded BS-IRS-user channel, i.e., in the hybrid spatial-angular domain channel matrix, all users share the same indices of non-zero columns, due to the common sparse BS-IRS channel. These user-independent non-zero columns’ indices are fed back by only one user, and the user-specific non-zero column vectors are fed back by different users, thus reducing feedback overhead. In addition to the common indices of non-zero columns in the beamspace channel, [24] found that the non-zero values in different non-zero columns have the same location offsets and amplitude ratios. These offsets and amplitudes are shared among all users and can be exploited to further reduce the number of feedback parameters.

Downlink CSI acquisition consists of both the downlink pilot transmission phase and the uplink feedback transmission phase. On one hand, the above works all assumed that each user knows its cascaded channels perfectly. However, due to the huge number of IRS reflecting elements, the overhead for the users to estimate their channels is huge, which is not considered in the above works. On the other hand, from feedback overhead perspective, there is a potential to significantly improve the performance of the schemes proposed in [19]–[24]. Recently, [1] revealed a unique property of the BS-IRS-user cascaded channels among different users — due to the common BS-IRS channel to all the users, each user’s cascaded channel vector is a scaled version of another user’s cascaded channel vector. However, [19]–[24] did not consider how to exploit this property to reduce the feedback overhead.

C. Main Contributions

This paper aims to significantly reduce the overhead of both the pilot transmission phase and the feedback transmission phase for CSI acquisition in IRS-assisted downlink communication. Actually, the CSI acquisition problem is essentially a distributed source coding (DSC) problem [25], where the core question is what information should be fed back by the users after they receive the pilot signals from the BS, such that the BS can acquire the CSI with the shortest pilot signals transmitted in the downlink and the minimum quantization bits transmitted in the uplink. Due to the correlation in cascaded channels among different users revealed in [1], the conventional “estimate-then-quantize” scheme, under which the users independently estimate the correlated channels, is not optimal. In this paper, our goal is to leverage the correlation among different users’ channels to reduce the overhead for CSI acquisition in IRS-assisted downlink communication. The main contributions of this paper are summarized as follows:

- We consider an FDD IRS-assisted downlink communication system with multiple single-antenna users. Moreover, we assume that the channels are quasi-static block fading channel model and the direct channels between the BS and users are blocked. In the FDD IRS-assisted downlink communication system, the conventional “estimate-then-quantize” scheme cannot utilize the channel correlation revealed in [1] because although the received signals of all the users are correlated, each user independently estimates its channels based on its own received signals. To overcome this issue, in this paper, we propose a novel “quantize-then-estimate” protocol. In sharp contrast to the “estimate-then-quantize” counterpart, our strategy makes each user first quantize its received pilot signals and then transmit the quantization bits to the BS. After de-quantization, the BS knows the pilot signals received by all the users, which contain the global CSI. Therefore, the BS can exploit the channel correlation revealed in [1] to jointly estimate the cascaded channels of all the users. The benefits of the proposed protocol for overhead reduction are two-fold. First, the BS is able to estimate the channels based on shorter pilot signals as shown in [1], i.e., the number of time samples for BS’s pilot transmission can be reduced. Second, because each user receives fewer pilot symbols, the number of quantization bits for feedback transmission is also reduced.
- Under our proposed protocol, we first design the codebook to quantize the received pilot signals via the Lloyd algorithm. Next, we design an efficient method such that the BS can estimate all the users’ cascaded channels based on its received feedback from the users. Specifically, we select several reference users, and the BS should estimate the ratios among the power of channels of the non-reference users and that of the reference users based on the feedback received in the first a few time samples, and estimate the channels of the reference users based on the feedback received in the remaining time samples. The linear minimum mean-squared error (LMMSE) estimation technique is proposed to estimate the channel ratios

and the channels under the above scheme. At last, we characterize the minimum overhead for transmitting the pilot signals in the downlink and the feedback signals in the uplink under our proposed “quantize-then-estimate” protocol, and the significant overhead reduction compared to the conventional “estimate-then-quantize” protocol is analytically demonstrated.

- To further improve the accuracy of CSI acquisition under our proposed protocol, we consider quantization bit allocation when users feed back their received pilot signals to the BS. Based on [26], when the number of IRS sub-surfaces is large, the received pilot signals tend to be Gaussian distributed. Then, rely on the the approximate Gaussian test channel [27], we characterize the rate-distortion trade-off for feedback transmission and propose an efficient approach to design the quantization bit allocation of each user to optimize the rate-distortion trade-off.

D. Organization

The rest of this paper is organized as follows. Section II describes the system model for our considered IRS-assisted multi-user downlink communication system. Section III reviews the traditional “estimate-then-quantize” CSI acquisition protocol and introduces our proposed “quantize-then-estimate” protocol. Section IV describes how to implement the “quantize-then-estimate” protocol in practice and characterize its minimum overhead. Section V designs a quantization bit allocation method for each user based on an approximated Gaussian test channel model. Section VI provides numerical examples to demonstrate the effectiveness of our proposed “quantize-then-estimate” scheme. Section VII concludes the paper.

Notation: \mathbf{I} and \mathbf{O} denote an identify matrix and an all-zero matrix, with appropriate dimensions. For a square full-rank matrix \mathbf{A} , \mathbf{A}^{-1} denotes its inverse. For a matrix \mathbf{B} , \mathbf{B}^T , \mathbf{B}^H , and \mathbf{B}^\dagger denotes its transpose, conjugate transpose, and pseudo-inverse matrix, respectively. $\lceil \cdot \rceil$ denotes the ceiling function. \otimes denotes the Kronecker product.

II. SYSTEM MODEL

We study the downlink communication in an FDD system which consists of a BS with M antennas, K single-antenna users, and an IRS with N passive reflecting elements, as shown in Fig. 1. In practice, the number of BS antennas is usually much larger than the number of users. Therefore, in this paper, we assume that $M > K$. In our considered IRS-assisted systems, the overhead to acquire CSI in IRS-assisted communication systems is high due to the large number of IRS elements N . To tackle this challenge, we adopt an IRS element grouping strategy as in [18], [28]. Specifically, the IRS elements are divided into D sub-surfaces to reduce the number of channels to be estimated, and the IRS elements within each sub-surface share a common reflection coefficient. Let $\phi_{d,i} = e^{j\theta_{d,i}} \in \mathbb{C}$ denote the reflection coefficient of the d -th IRS sub-surface at time sample i , where $\theta_{d,i} \in [0, 2\pi)$ denotes the phase shift of $\phi_{d,i}$.

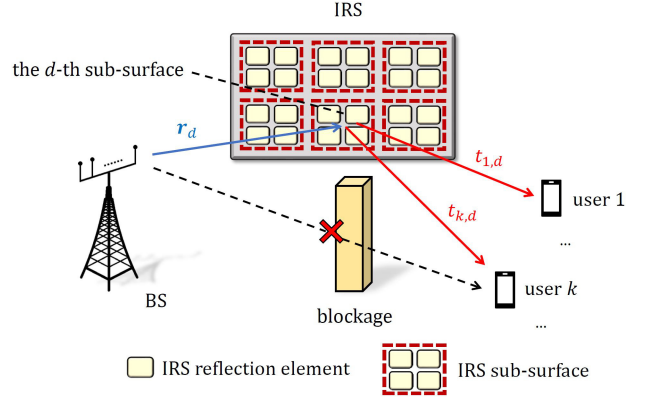


Fig. 1. An IRS-aided downlink communication system.

We assume a quasi-static block fading channel model, where the channels remain approximately constant in each coherence block. In addition, as shown in Fig. 1, we assume that the direct channels between the BS and users are blocked, and the signals can only be transmitted through the IRS reflecting channels to the users. The baseband equivalent channels from the BS to the d -th IRS sub-surface, and from the d -th IRS sub-surface to user k are denoted by $\mathbf{r}_d \in \mathbb{C}^{M \times 1}$ and $t_{k,d} \in \mathbb{C}$, $k = 1, \dots, K$, $d = 1, \dots, D$, respectively. For convenience, define $\mathbf{R} = [\mathbf{r}_1, \dots, \mathbf{r}_D]$ as the overall channels from the BS to the IRS, and $\mathbf{t}_k = [t_{k,1}, \dots, t_{k,D}]^T$ as the channels from the IRS to user k . Then, the cascaded reflecting channels from the BS to user k through the d -th IRS sub-surface is expressed as

$$\mathbf{g}_{k,d} = t_{k,d} \mathbf{r}_d \in \mathbb{C}^{M \times 1}, \quad \forall d, k. \quad (1)$$

In downlink communication, the pilot signals transmitted from the BS at time sample i are denoted by $\mathbf{x}_i = [x_{i,1}, \dots, x_{i,M}]^T \in \mathbb{C}^{M \times 1}$, $i = 1, \dots, T$, where $x_{i,m}$ is the i -th pilot sample transmitted by antenna m , and T is the number of time samples to transmit the pilots. Then, the received signal of user k at time sample i is expressed as

$$y_{k,i} = \sum_{d=1}^D \phi_{d,i} \mathbf{g}_{k,d}^T \sqrt{p_i} \mathbf{x}_i + z_{k,i}, \quad \forall k, i, \quad (2)$$

where p_i denotes the transmit power at the BS at time sample i , and $z_{k,i} \sim \mathcal{CN}(0, \sigma_z^2)$ denotes the additive white Gaussian noise (AWGN) of user k at time sample i .

III. A NOVEL QUANTIZE-THEN-ESTIMATE PROTOCOL

In this paper, we aim to propose a low-overhead channel estimation and feedback scheme such that the BS can efficiently acquire the cascaded channels $\mathbf{g}_{k,d}$'s, $\forall k, d$, to design the BS beamforming vectors and the IRS reflecting coefficients. In this section, we will first overview the conventional method for channel estimation and feedback in our considered IRS-assisted FDD systems, and show its disadvantages. Then, we will propose a novel strategy that can exploit the unique channel property in IRS-assisted communication to reduce the channel estimation and feedback overhead.

A. Traditional Estimate-then-Quantize Scheme

Without IRS, channel estimation and feedback have been widely studied for downlink FDD systems [7], [8], [29]. In the IRS-assisted network, we may adopt similar philosophy and apply the following “estimate-then-quantize” strategy for channel estimation and feedback.

- **Phase I (Estimation):** In the first phase, each user k first receives the pilot signals from the BS as shown in (2) and then estimates its cascaded channel $\mathbf{g}_{k,d}$'s, $\forall d$, based on the existing method proposed in [18], [28], [30], [31].
- **Phase II (Feedback):** In the second phase, each user quantizes its estimated channels and feeds back the quantization bits to the BS. The BS de-quantizes the quantization bits to recover the cascaded channels of all the users.

Although the above “estimate-then-quantize” strategy works well in conventional systems without IRS, its overall overhead for pilot and feedback transmission is significant in our considered IRS-assisted systems. Specifically, under the above scheme, the minimum number of pilot samples for user k to estimate $\mathbf{g}_{k,d}$'s, $d = 1, \dots, D$, in Phase I is $T_{\min} = MD$ [1]. Moreover, each user k needs to feed back MD channel coefficients in $\mathbf{g}_{k,d}$'s, $d = 1, \dots, D$, to the BS in Phase II. Recently, it has been revealed in [1] that there is a lot of redundancy in users' cascaded channels and the number of independent unknown variables in $\mathbf{g}_{k,d}$'s, $k = 1, \dots, K, d = 1, \dots, D$, is much smaller than KMD . Specifically, if we focus on a particular IRS sub-surface d , then the channel between the BS and the d -th sub-surface of the IRS, i.e., r_d , is common among the cascaded channels $\mathbf{g}_{k,d}$'s of all the users. As a result, according to (1), we have

$$\mathbf{g}_{k,d} = \lambda_{k,d} \mathbf{g}_{k_d,d}, \quad \forall k \neq k_d, \quad d = 1, \dots, D, \quad (3)$$

where k_d is the index of the reference user selected for IRS sub-surface d , and the channel ratio between user k and reference user k_d is given by

$$\lambda_{k,d} = \frac{t_{k,d}}{t_{k_d,d}}. \quad (4)$$

(3) and (4) indicate that for each IRS sub-surface d , if the cascaded channel of the reference user k_d , i.e., $\mathbf{g}_{k_d,d}$, is known, a scalar $\lambda_{k,d}$ is sufficient for the BS to know the cascaded channel vector of user $k \neq k_d$, i.e., $\mathbf{g}_{k,d}$. In other words, the BS just needs to know $MD + (K - 1)D$ channel coefficients in $\mathbf{g}_{k_d,d}$'s, $d = 1, \dots, D$, and $\lambda_{k,d}$'s, $\forall k \neq k_d, d = 1, \dots, D$. Therefore, the number of time samples for pilot transmission in Phase I and the number of quantization bits in Phase II can be hugely reduced in IRS-assisted downlink communication, if the channel property shown in (3) and (4) can be properly utilized. However, the conventional “estimate-then-quantize” scheme does not take advantage of (3) and (4) for reducing the overhead.

B. Proposed Quantize-then-Estimate Scheme

In this sub-section, we propose a novel strategy that can leverage (3) and (4) to significantly reduce the overhead for channel estimation and feedback in IRS-assisted downlink

communication. Note that (3) and (4) reveal the correlation among different users' cascaded channels, while the users are distributed and cannot cooperate with each other to leverage such correlation for channel estimation. To overcome this issue, we propose a “quantize-then-estimate” protocol, where the channels of all the users are estimated at the BS side by utilizing (3) and (4), rather than at the distributed user side. The proposed scheme is detailed as below.

- **Phase I (Feedback):** In the first phase, all the users receive the pilot signals from the BS as shown in (2). However, instead of estimating its own channels, each user k quantizes its received pilot signals over T time samples, i.e., $\mathbf{y}_k = [y_{k,1}, \dots, y_{k,T}]^T$, and feeds back the quantization bits to the BS.
- **Phase II (Estimation):** In the second phase, the BS de-quantizes the received quantization bits for recovering the pilot signals received by all the users. Then, with the above global information about users' received pilots, the BS is able to leverage the correlation among different users' channels shown in (3) and (4) to estimate $\mathbf{g}_{k,d}$'s more efficiently.

Note that the key difference between our proposed “quantize-then-estimate” scheme and the conventional “estimate-then-quantize” scheme shown in Section III-A lies in what is fed back from the users to the BS and who performs channel estimation. Under our scheme, users feed back their received pilot signals such that the BS can perform a joint estimation of different users' channels by leveraging (3) and (4). The benefits of the above joint estimation are two-fold. First, as shown in [1], the minimum number of pilot samples for channel estimation, i.e., T_{\min} , can be reduced from MD to $\max\{M - 1, D + \lceil (M - 1)D/K \rceil\}$ by leveraging (3) and (4). Second, because signals from fewer time samples are quantized, the feedback overhead is significantly reduced. Specifically, without utilizing (3) and (4), all the users need to quantize KMD samples in $\mathbf{g}_{k,d}$'s; while under our proposed scheme, as will be shown later in Section IV, all the users only need to quantize $KD + K \cdot \max\{M - 1, D + \lceil (M - 1)D/K \rceil\}$ samples in \mathbf{y}_k 's.

IV. QUANTIZATION AND ESTIMATION DESIGN

In this section, we will introduce how to quantize $y_{k,i}$'s at the user side and how to estimate the channels via leveraging (3) and (4) at the BS side, in Phase I and Phase II of our proposed “quantize-then-estimate” protocol, respectively. We will also analytically characterize the overhead for pilot and feedback transmission under our proposed scheme and show the significant overhead reduction over the conventional “estimate-then-quantize” protocol.

A. Quantization at User Side

In this sub-section, we introduce how the users quantize the received pilot signals. For each user k , we assume independent signal quantization at different time samples. Specifically, each user k quantizes $y_{k,1}, \dots, y_{k,T}$ subsequently via scalar

quantization. At time sample i , the codebook for quantizing $y_{k,i}$ is denoted by

$$\mathcal{C}_{k,i} = \{c_{k,i,1}, \dots, c_{k,i,L_{k,i}}\}, \quad \forall k, i, \quad (5)$$

which consists of $L_{k,i}$ codewords and is shared by user k and the BS. We will introduce how to design $L_{k,i}$'s, $\forall k$ and i , in Section V. Based on the probability density function of $y_{k,i}$, the codebook given any $L_{k,i}$ can be designed via the Lloyd algorithm [32]. Given the codebook $\mathcal{C}_{k,i}$, the codeword index and the corresponding codeword to quantize $y_{k,i}$ are given by

$$\begin{aligned} l_{k,i}^* &= \arg \min_{l_{k,i}=1, \dots, L_{k,i}} \mathcal{D}(y_{k,i}, c_{k,i,l}), \\ \tilde{y}_{k,i} &= c_{k,i,l_{k,i}^*}, \quad \forall k, i, \end{aligned} \quad (6)$$

where $\mathcal{D}(y_{k,i}, c_{k,i,l}) = \|y_{k,i} - c_{k,i,l}\|_2^2$ denotes the distortion function between $y_{k,i}$ and the codeword $c_{k,i,l}$. The quantized signal of $y_{k,i}$ can be expressed as

$$\tilde{y}_{k,i} = y_{k,i} + e_{k,i}, \quad \forall k, i, \quad (7)$$

where $e_{k,i}$ denotes the error to quantize the signal $y_{k,i}$ with zero mean and variance $q_{k,i}$. Note that $e_{k,i}$'s are independent over i due to scalar quantization at each time sample.

For user k , each codeword index $l_{k,i}^*$ is represented by $\lceil \log_2 L_{k,i} \rceil$ quantization bits. Since independent quantization is performed at different time samples, the total number of bits for user k to quantize all its received pilot signals is the sum of the number of bits over T time samples, i.e.,

$$B_k = \sum_{i=1}^T \lceil \log_2 L_{k,i} \rceil, \quad \forall k. \quad (8)$$

Then, each user modulates its quantization bits onto quadrature amplitude modulation (QAM) symbols and sends these symbols to the BS via a feedback channel, which is assumed to be error-free [33]. Denote the modulation rate of user k as μ_k bits/sample, $\forall k$. Thus, the number of time samples for user k to feed back the QAM symbols to the BS is

$$T_{\text{fb},k} = \frac{B_k}{\mu_k} = \frac{1}{\mu_k} \sum_{i=1}^T \lceil \log_2 L_{k,i} \rceil, \quad \forall k. \quad (9)$$

As show in [33], in our considered setup where the number of BS antennas is larger than that of users, i.e., $M > K$, the BS can mitigate inter-user interference via zero-forcing beamforming design. Therefore, all the users can transmit the feedback symbols simultaneously to the BS without inter-user interference. As a result, the number of time samples required for all the users to finish feedback transmission is determined by the user that needs the largest number of time samples, i.e.,

$$T_{\text{fb}} = \max_{1 \leq k \leq K} T_{\text{fb},k}. \quad (10)$$

B. Channel Estimation at BS Side

In this sub-section, we introduce how the BS can estimate the channels. To begin with, the BS collects the QAM symbols from the users and de-modulates the symbols to quantization bits. Similar to [33], we assume that the feedback channels

from the users to the BS are error-free such that the BS can perfectly decode the quantization bits and then recover $\tilde{y}_{k,i}$'s.

Subsequently, based on the de-quantized signals $\tilde{y}_{k,i}$'s, the BS adopts a two-step channel estimation method, where in the first step with a duration of $\tau_1 < T$ samples, the BS estimates $\lambda_{k,d}$'s, $\forall k \neq k_d$ and $\forall d$, based on $\tilde{y}_{k,1}, \dots, \tilde{y}_{k,\tau_1}$ for non-reference users; while in the second step with a duration of $\tau_2 = T - \tau_1$ samples, the BS estimates $\mathbf{g}_{k_d,d}$'s based on $\tilde{y}_{k,\tau_1+1}, \dots, \tilde{y}_{k,\tau_1+\tau_2}$, $\forall k, d$, for reference users. Last, user k 's cascaded channels can be recovered based on (3), $\forall k$. In the following, we introduce how to estimate $\lambda_{k,d}$'s in Step 1 and $\mathbf{g}_{k_d,d}$'s in Step 2.

Step 1 (Estimation of channel ratios $\lambda_{k,d}$'s): In the first step, the de-quantized received signals $\tilde{y}_{k,i}$'s given in (7) can be re-written as

$$\tilde{y}_{k,i} = \sum_{d=1}^D \phi_{d,i} \alpha_{k,d,i} + z_{k,i} + e_{k,i}, \quad \forall k, i = 1, \dots, \tau_1, \quad (11)$$

where

$$\alpha_{k,d,i} = \begin{cases} \sqrt{p_i} \mathbf{g}_{k_d,d}^T \mathbf{x}_i, & \text{if } k = k_d, \\ \sqrt{p_i} \lambda_{k,d} \mathbf{g}_{k_d,d}^T \mathbf{x}_i, & \text{if } k \neq k_d, \end{cases} \quad \forall d, i = 1, \dots, \tau_1. \quad (12)$$

Note that if we can perfectly recover $\alpha_{k,d,i}$'s in (12), then the channel ratio can be obtained by

$$\lambda_{k,d} = \frac{\alpha_{k,d,i}}{\alpha_{k_d,d,i}}, \quad \forall k \neq k_d, \forall d, i = 1, \dots, \tau_1. \quad (13)$$

The main challenge to estimate $\alpha_{k,d,i}$'s in (12) is that for each k , the BS has to estimate $\tau_1 D$ unknown variables, i.e., $\alpha_{k,d,i}$'s, $d = 1, \dots, D, i = 1, \dots, \tau_1$, using $\tau_1 < \tau_1 D$ observations, i.e., $\tilde{y}_{k,i}$'s, $i = 1, \dots, \tau_1$. To solve this challenge, we propose to set identical pilot signals over all time samples, i.e.,

$$\sqrt{p_i} \mathbf{x}_i = \sqrt{p} \mathbf{x}, \quad \forall i. \quad (14)$$

In this case, the received signal in (11) reduces to

$$\tilde{y}_{k,i} = \sum_{d=1}^D \phi_{d,i} \alpha_{k,d} + z_{k,i} + e_{k,i}, \quad \forall k, i = 1, \dots, \tau_1, \quad (15)$$

where

$$\alpha_{k,d} = \begin{cases} \sqrt{p} \mathbf{g}_{k_d,d}^T \mathbf{x}, & \text{if } k = k_d, \\ \sqrt{p} \lambda_{k,d} \mathbf{g}_{k_d,d}^T \mathbf{x}, & \text{if } k \neq k_d, \end{cases} \quad \forall d. \quad (16)$$

Note that in (15), for each k , there are D unknown variables, i.e., $\alpha_{k,d}$, $d = 1, \dots, D$, rather than $\tau_1 D$ variables as in (12). Moreover, if $\alpha_{k,d}$'s can be perfectly estimated, $\lambda_{k,d}$'s can be estimated as

$$\lambda_{k,d} = \frac{\alpha_{k,d}}{\alpha_{k_d,d}}, \quad \forall k \neq k_d, \forall d. \quad (17)$$

In the following, we show how to estimate $\alpha_{k,d}$'s based on (15). Let $\tilde{\mathbf{y}}_k^{(1)} = [\tilde{y}_{k,1}, \dots, \tilde{y}_{k,\tau_1}]^T$ denote the overall quantized received signals of user k over τ_1 time samples. Then, (15) is equivalent to

$$\tilde{\mathbf{y}}_k^{(1)} = \Phi_1 \alpha_k + \mathbf{z}_k^{(1)} + \mathbf{e}_k^{(1)}, \quad \forall k, \quad (18)$$

where

$$\Phi_1 = \begin{bmatrix} \phi_{1,1} & \cdots & \phi_{D,1} \\ \vdots & \ddots & \vdots \\ \phi_{1,\tau_1} & \cdots & \phi_{D,\tau_1} \end{bmatrix}, \quad (19)$$

$\alpha_k = [\alpha_{k,1}, \dots, \alpha_{k,D}]^T$, $\mathbf{z}_k^{(1)} = [z_{k,1}, \dots, z_{k,\tau_1}]^T$, and $\mathbf{e}_k^{(1)} = [e_{k,1}, \dots, e_{k,\tau_1}]^T$. If there is no noise and quantization error, i.e., $\mathbf{z}_k^{(1)} = \mathbf{0}$ and $\mathbf{e}_k^{(1)} = \mathbf{0}$, then we can set Φ_1 as a discrete Fourier transform (DFT) matrix and then perfectly estimate α_k using

$$\bar{\tau}_1 = D \quad (20)$$

time samples. In the practical case with noise and quantization error, we can set $\tau_1 \geq D$ and Φ_1 as the first D columns of a $\tau_1 \times \tau_1$ DFT matrix. Then, we can apply the LMMSE technique to estimate α_k as

$$\begin{aligned} \hat{\alpha}_k &= [\hat{\alpha}_{k,1}, \dots, \hat{\alpha}_{k,D}]^T \\ &= \Lambda_k \Phi_1^H \left(\Phi_1 \Lambda_k \Phi_1^H + \sigma_z^2 \mathbf{I}_D + \mathbf{E}_k^{(1)} \right)^{-1} \tilde{\mathbf{y}}_k^{(1)}, \end{aligned} \quad (21)$$

where Λ_k denotes the covariance matrix of α_k , and $\mathbf{E}_k^{(1)}$ denotes the covariance matrix of $\mathbf{e}_k^{(1)}$. Then, the estimated $\lambda_{k,d}$ is expressed as

$$\hat{\lambda}_{k,d} = \frac{\hat{\alpha}_{k,d}}{\hat{\alpha}_{k,d}}, \quad \forall k \neq k_d, \forall d. \quad (22)$$

Given any reference user selection strategy k_d 's, $d = 1, \dots, D$, we can estimate $\lambda_{k,d}$'s, $k \neq k_d$, using the above method. However, the selection of the reference user for each IRS sub-surface d can significantly affect the accuracy for estimating $\lambda_{k,d}$'s, $\forall k \neq k_d, \forall d$. This is because if user k with a very weak value of $|\hat{\alpha}_{k,d}|$ is selected as the reference user, then a very small error for estimating $\alpha_{k,d}$, i.e., $\hat{\alpha}_{k,d} - \alpha_{k,d}$, can cause a significant error for estimating $\lambda_{k,d}$ in (22). Therefore, for each IRS sub-surface d , we select the reference user as follows

$$k_d = \arg \max_{k=1, \dots, K} |\hat{\alpha}_{k,d}|^2, \quad \forall d. \quad (23)$$

Step 2 (Estimation of reference users' channels $\mathbf{g}_{k_d,d}$'s): In the second step, we estimate the channels of reference users, i.e., $\mathbf{g}_{k_d,d}$'s, $d = 1, \dots, D$. Before introducing Step 2, we want to emphasize that after $\hat{\alpha}_{k,d}$'s are estimated in Step 1, we already have some useful information about $\mathbf{g}_{k_d,d}$'s according to (16):

$$\hat{\alpha}_{k_d,d} = \sqrt{p} \mathbf{x}^T \mathbf{g}_{k_d,d} + \beta_{k_d,d}, \quad d = 1, \dots, D, \quad (24)$$

where $\beta_{k_d,d}$ denotes the error for estimating $\alpha_{k_d,d}$. Define that $\hat{\alpha} = [\hat{\alpha}_{k_1,1}, \dots, \hat{\alpha}_{k_D,D}]^T$ and $\mathbf{g} = [\mathbf{g}_{k_1,1}^T, \dots, \mathbf{g}_{k_D,D}^T]^T$. Then, we have

$$\hat{\alpha} = \mathbf{F} \mathbf{g} + \beta, \quad (25)$$

where $\mathbf{F} = \sqrt{p} \mathbf{I}_D \otimes \mathbf{x}^T \in \mathbb{C}^{D \times MD}$, and $\beta = [\beta_{k_1,1}, \dots, \beta_{k_D,D}]^T$. This information should be used in Step 2 to estimate \mathbf{g} .

In Step 2, the received pilots at time sample i is given as

$$\begin{aligned} \tilde{\mathbf{y}}_i^{(2)} &= [\tilde{y}_{1,i}, \dots, \tilde{y}_{K,i}]^T = \sqrt{p} \sum_{d=1}^D \phi_{d,i} \mathbf{x}_i^T \mathbf{g}_{k_d,d} \lambda_d + \mathbf{e}_i^{(2)} \\ &= \mathbf{F}_i \mathbf{g} + \mathbf{e}_i^{(2)}, \quad i = \tau_1 + 1, \dots, \tau_1 + \tau_2, \end{aligned} \quad (26)$$

where $\lambda_d = [\lambda_{1,d}, \dots, \lambda_{K,d}]^T$ with $\lambda_{k_d,d} = 1, \forall d$, $\mathbf{e}_i^{(2)} = [z_{1,i} + e_{1,i}, \dots, z_{K,i} + e_{K,i}]^T$, and

$$\mathbf{F}_i = \sqrt{p} [\phi_{1,i} \mathbf{x}_i^T \otimes \lambda_1, \dots, \phi_{D,i} \mathbf{x}_i^T \otimes \lambda_D] \in \mathbb{C}^{K \times MD}. \quad (27)$$

Since both $\hat{\alpha}$ in Step 1 and $\tilde{\mathbf{y}}_i^{(2)}$'s in Step 2 contain information about \mathbf{g} , we define

$$\tilde{\mathbf{y}}^{(2)} = [(\tilde{\mathbf{y}}_{\tau_1+1}^{(2)})^T, \dots, (\tilde{\mathbf{y}}_{\tau_1+\tau_2}^{(2)})^T, \hat{\alpha}^T]^T. \quad (28)$$

According to (25) and (26), we have

$$\begin{aligned} \tilde{\mathbf{y}}^{(2)} &= \Theta \mathbf{g} + \mathbf{e}^{(2)} \\ &= \hat{\Theta} \mathbf{g} + (\Theta - \hat{\Theta}) \mathbf{g} + \mathbf{e}^{(2)}, \end{aligned} \quad (29)$$

where $\Theta = [\mathbf{F}_{\tau_1+1}^T, \dots, \mathbf{F}_{\tau_1+\tau_2}^T, \mathbf{F}^T]^T$, $\hat{\Theta}$ is an estimation of Θ with $\lambda_{k,d}$'s (as shown in (27), \mathbf{F}_i 's are functions of $\lambda_{k,d}$'s) replaced by their estimations $\hat{\lambda}_{k,d}$'s given in (22), and $\mathbf{e}^{(2)} = [(\mathbf{e}_{\tau_1+1}^{(2)})^T, \dots, (\mathbf{e}_{\tau_1+\tau_2}^{(2)})^T, \beta^T]^T$. Note that in (29), $\hat{\Theta}$ is known by the BS, and $\Theta - \hat{\Theta}$ is the unknown estimation error. If there is no noise and quantization error in both Step 1 (such that $\hat{\lambda}_{k,d} = \lambda_{k,d}, \forall k, d$, and $\Theta = \hat{\Theta}$) and Step 2, i.e., $\mathbf{e}^{(2)} = \mathbf{0}$, then (29) reduces to

$$\tilde{\mathbf{y}}^{(2)} = \Theta \mathbf{g}. \quad (30)$$

Theorem 1: In the ideal case that there is no noise in users' received pilots given in (2) and quantization error in BS's received signals given in (29), i.e., $\Theta = \hat{\Theta}$ and $\mathbf{e}^{(2)} = \mathbf{0}$, the minimum number of time samples for the BS to perfectly estimate \mathbf{g} based on (30) is

$$\bar{\tau}_2 = \max \left\{ M - 1, \left\lceil \frac{(M-1)D}{K} \right\rceil \right\}. \quad (31)$$

Proof: Please refer to Appendix. \blacksquare

To summarize, the minimum number of time samples to transmit pilot signals from the BS to the users is

$$T_{\min}^* = \bar{\tau}_1 + \bar{\tau}_2 = D + \max \left\{ M - 1, \left\lceil \frac{(M-1)D}{K} \right\rceil \right\}. \quad (32)$$

In the practical case with noise and quantization error, we can use $\tau_2 \geq \bar{\tau}_2$ time samples for pilot transmission. However, the unknown error propagated from Step 1, i.e., $\Theta - \hat{\Theta}$ in (29), makes it hard to obtain the LMMSE estimator of \mathbf{g} . Actually, we can increase the number of time samples for pilot transmission in Step 1 such that $\Theta - \hat{\Theta}$ is sufficiently small. In this case, we assume that $\Theta - \hat{\Theta} \approx \mathbf{0}$ as in [34]. Then, (29) reduces to

$$\tilde{\mathbf{y}}^{(2)} \approx \hat{\Theta} \mathbf{g} + \mathbf{e}^{(2)}. \quad (33)$$

Based on (33), the LMMSE estimator of \mathbf{g} can be designed. Specifically, we can set pilot signals \mathbf{x}_i 's and IRS reflection

coefficients $\phi_{d,i}$'s, $d = 1, \dots, D$, $i = \tau_1 + 1, \dots, \tau_1 + \tau_2$, according to the orthogonal transmission and reflection strategy in Section V-C in [1]. In this case, based on (33), the LMMSE estimator of \mathbf{g} is given as

$$\hat{\mathbf{g}} = \mathbf{G}\hat{\Theta}^H \left(\hat{\Theta}\mathbf{G}\hat{\Theta}^H + \mathbf{E}^{(2)} \right)^{-1} \tilde{\mathbf{y}}^{(2)}, \quad (34)$$

where \mathbf{G} denotes the covariance matrix of \mathbf{g} , and $\mathbf{E}^{(2)}$ denotes the covariance matrix of $\mathbf{e}^{(2)}$. With the estimations of $\lambda_{k,d}$'s given in (22) and those of $\mathbf{g}_{k,d}$'s given in (34), the cascaded channels can be estimated as

$$\hat{\mathbf{g}}_{k,d} = \hat{\lambda}_{k,d} \hat{\mathbf{g}}_{k,d}, \quad k = 1, \dots, K, d = 1, \dots, D. \quad (35)$$

Remark 1: Until now, we have shown that under our proposed quantize-and-estimate protocol, the minimum numbers of time samples for pilot transmission and feedback transmission are (32) and (10), respectively. Therefore, the minimum number of time samples for pilot and feedback transmission is

$$\begin{aligned} T_{\text{tot}} &= T_{\text{min}}^* + T_{\text{fb}} \\ &= T_{\text{min}}^* + \max_{1 \leq k \leq K} \left\{ \frac{1}{\mu_k} \sum_{i=1}^{T_{\text{min}}^*} \lceil \log_2 L_{k,i} \rceil \right\}. \end{aligned} \quad (36)$$

Note that under the conventional ‘‘estimate-then-quantize’’ strategy, the minimum number of time samples for pilot transmission is [1]

$$T_{\text{cov}}^* = MD. \quad (37)$$

Then, the minimum number of time samples for pilot and feedback transmission is

$$T_{\text{tot,cov}} = T_{\text{cov}}^* + \max_{1 \leq k \leq K} \left\{ \frac{MD \lceil \log_2 L_k \rceil}{\mu_k} \right\}, \quad (38)$$

where L_k is the codebook size for each user k to quantize each estimated channel coefficient. Because fewer pilot samples are transmitted in Phase I and quantized in Phase II thanks to the utilization of channel correlation shown in (3), the overhead of our proposed ‘‘quantize-then-estimate’’ strategy characterized in (36) is significantly reduced compared with that of the conventional ‘‘estimate-then-quantize’’ strategy characterized in (38).

At last, given the minimum number of time samples for pilot and feedback transmission, we characterize the computational complexity of our proposed scheme and the conventional ‘‘estimate-then-quantize’’ scheme. For the feedback phase, we count the total time of exhaustively searching the codebook for selecting the optimal codeword as the measure of complexity. If B bits are used to quantize each complex symbol, then the complexity in the feedback phase under our proposed scheme is $\mathcal{O}((D + \lceil (M-1)D/K \rceil)2^B)$, and that under the ‘‘estimate-then-quantize’’ scheme is $\mathcal{O}(MD2^B)$, respectively. For the estimation phase, we count the number of complex multiplication (CM) [11] as the measure of complexity. The main complexity of our proposed scheme comes from computing (21) and (34). It is straightforward to see that the two operations require $\Delta_1 = D^2(1 + \beta + 3D)$ CMs and $\Delta_2 = (K \lceil (M-1)D/K \rceil + D)(MD)^2 + (K \lceil (M-$

$1)D/K \rceil + D)^2(2MD + \beta + 1)$ CMs, respectively, where β is a scaling factor depending on the specific algorithm for the matrix inversion. As a result, the total computational complexity can be expressed as $\mathcal{O}(K\Delta_1 + \Delta_2)$. The complexity of the ‘‘estimate-then-quantize’’ scheme can be expressed as $\mathcal{O}((MD)^2(3MD + \beta + 1))$. Therefore, the complexity in the estimation phase of the two schemes are both nearly $\mathcal{O}((MD)^3)$, while the complexity in the feedback phase under our proposed scheme is much lower than that under the ‘‘estimate-then-quantize’’ scheme.

V. QUANTIZATION BIT ALLOCATION

In section IV, we have proposed efficient methods to quantize pilot signals at the user side and estimate the channels at the BS side. A remaining issue that is not tackled is how to determine the quantization bit allocation policy for each user to achieve the best rate-distortion trade-off under our proposed ‘‘quantize-then-estimate’’ scheme. The quantized version of user k 's received pilot signal at time sample i is given in (7). However, the distribution of the quantization errors under the Lloyd algorithm is usually non-trivial, and it is thus difficult to analyze the rate-distortion trade-off. In this section, we will apply the Gaussian test channel model to approximate the rate-distortion trade-off achieved by the Lloyd algorithm, based on which we are able to optimize the quantization bit allocation, i.e., $L_{k,i}$'s. As will be shown later in this section, the Gaussian test channel model yields an analytical expression for the rate-distortion performance. Moreover, [27] has shown analytically and numerically that the rate-distortion trade-off obtained under the Gaussian test channel model is a very tight approximation to that achieved by the practical Lloyd algorithm. We will also numerically verify the tightness of the Gaussian test channel model in terms of rate-distortion approximation later in this section.

A. Gaussian Test Channel and Rate-Distortion Trade-off

The Gaussian test channel model to approximate (7) is given as

$$\tilde{y}_{k,i} = y_{k,i} + \tilde{e}_{k,i}, \quad \forall k, i, \quad (39)$$

where $\tilde{e}_{k,i}$ has the same variance as $e_{k,i}$ in (7) and is assumed to be Gaussian distributed, i.e., $\tilde{e}_{k,i} \sim \mathcal{CN}(0, q_{k,i})$, and is independent with $y_{k,i}$, $\forall k, i$. Moreover, because each user k applies scalar quantization on $y_{k,i}$'s, $\forall i$, $\tilde{e}_{k,i}$'s are independent over i .

As shown in Lemma 1 of [26], when the number of IRS sub-surfaces D is large, $y_{k,i}$, $\forall k, i$, tends to be Gaussian distributed in general, i.e., $y_{k,i} \sim \mathcal{CN}(0, u_{k,i})$, where $u_{k,i} = \mathbb{E}[|y_{k,i}|^2]$ denotes the variance of $y_{k,i}$, $\forall k, i$. Then, under the Gaussian test channel model given in (39), the BS can adopt an MMSE estimator to recover $y_{k,i}$ based on $\tilde{y}_{k,i}$, and the corresponding estimation MSE is [35]

$$\gamma_{k,i} = \frac{u_{k,i}q_{k,i}}{u_{k,i} + q_{k,i}} \quad \forall k, i. \quad (40)$$

Therefore, given $\mathbf{q}_k = [q_{k,1}, \dots, q_{k,T}]^T$, the overall MSE to estimate $\mathbf{y}_k = [y_{k,1}, \dots, y_{k,T}]^T$ is given as

$$\Gamma_k(\mathbf{q}_k) = \sum_{i=1}^T \gamma_{k,i} = \sum_{i=1}^T \frac{u_{k,i} q_{k,i}}{u_{k,i} + q_{k,i}}, \quad \forall k. \quad (41)$$

Last, because different users independently quantize their received signals, the overall MSE to estimate $\mathbf{y}_1, \dots, \mathbf{y}_K$ is given as

$$\Gamma_{\text{sum}}(\mathbf{q}) = \sum_{k=1}^K \Gamma_k(\mathbf{q}_k) = \sum_{k=1}^K \sum_{i=1}^T \frac{u_{k,i} q_{k,i}}{u_{k,i} + q_{k,i}}, \quad (42)$$

where $\mathbf{q} = [\mathbf{q}_1^T, \dots, \mathbf{q}_K^T]^T$. Moreover, according to (39), the number of bits to quantize the sample $y_{k,i}$ is [36]

$$\begin{aligned} \mathcal{I}(\tilde{y}_{k,i}, y_{k,i}) &= \mathcal{H}(\tilde{y}_{k,i}) - \mathcal{H}(\tilde{y}_{k,i} | y_{k,i}) = \mathcal{H}(\tilde{y}_{k,i}) - \mathcal{H}(\tilde{e}_{k,i}) \\ &= \log_2 [\pi e (u_{k,i} + q_{k,i})] - \log_2 (\pi e q_{k,i}) \\ &= \log_2 \left(1 + \frac{u_{k,i}}{q_{k,i}} \right), \quad \forall k, i, \end{aligned} \quad (43)$$

where $\mathcal{I}(\tilde{y}_{k,i}, y_{k,i})$ is the mutual information between $\tilde{y}_{k,i}$ and $y_{k,i}$, $\mathcal{H}(\cdot)$ denotes the differential entropy, and the second equality is due to independence of $y_{k,i}$ and $\tilde{e}_{k,i}$. Then, the overall number of bits for user k to quantize \mathbf{y}_k over T time samples is

$$B_k^{(G)}(\mathbf{q}_k) = \sum_{i=1}^T \log_2 \left(1 + \frac{u_{k,i}}{q_{k,i}} \right), \quad \forall k. \quad (44)$$

Thus, similar to (10), the feedback transmission time (in terms of samples) under the Gaussian test channel model (39) is

$$T_{\text{fb}}^{(G)} = \max_{1 \leq k \leq K} T_{\text{fb},k}^{(G)}, \quad (45)$$

where

$$T_{\text{fb},k}^{(G)} = \frac{B_k^{(G)}(\mathbf{q}_k)}{\mu_k}, \quad \forall k. \quad (46)$$

According to (45) and (42), the rate-distortion trade-off under the Gaussian test channel model can be characterized by the following optimization problem

$$\begin{aligned} \text{(P0):} \quad & \text{Minimize}_{\mathbf{q}} \Gamma_{\text{sum}}(\mathbf{q}) \\ \text{Subject to} \quad & \frac{B_k^{(G)}(\mathbf{q}_k)}{\mu_k} \leq \bar{T}_{\text{fb}}, \quad \forall k. \end{aligned} \quad (47)$$

where \bar{T}_{fb} is the given time constraint for quantization bit transmission. In other words, we aim to characterize given the feedback time constraint, what is the minimum quantization MSE.

B. Quantization Bit Allocation

In this section, we aim to design the quantization bit allocation solution of each user by solving problem (P0). According to (42), which is due to the fact that users independently quantize their received signals, all the users can obtain their quantization bit allocation solutions in parallel. Specifically,

the quantization bit allocation policy of user k can be obtained by solving the following sub-problem:

$$\begin{aligned} \text{(P1-}k\text{):} \quad & \text{Minimize}_{\{q_{k,i}\}} \Gamma_k(\mathbf{q}_k) \\ \text{Subject to} \quad & B_k^{(G)}(\mathbf{q}_k) \leq \mu_k \bar{T}_{\text{fb}}. \end{aligned} \quad (48)$$

In the following, we propose an efficient algorithm to solve problem (P1- k) for user k , $k = 1, \dots, K$. Specifically, it can be shown that $\Gamma_k(\mathbf{q}_k)$ given in (41) and $B_k^{(G)}(\mathbf{q}_k)$ given in (44) are concave and convex functions over $q_{k,i}$'s, respectively. Therefore, the challenge is that we are minimizing a concave function, rather than a convex function. In this case, the majorization-minimization (MM) algorithm can be used to obtain a locally optimal solution [37].

MM is an iterative algorithm. Under the s -th iteration of the MM algorithm, we need to find a surrogate function of the objective function of problem (P1- k). Because $\Gamma_k(\mathbf{q}_k)$ is a concave function, its first-order Taylor expansion serves as its upper bound and can thus be used as its surrogate function. Specifically, given any point $\mathbf{q}_k^{(s)} = [q_{k,1}^{(s)}, \dots, q_{k,T}^{(s)}]^T$, the surrogate function of $\Gamma_k(\mathbf{q}_k)$ can be set as

$$\begin{aligned} f(\mathbf{q}_k | \mathbf{q}_k^{(s)}) &= \Gamma_k(\mathbf{q}_k^{(s)}) + (\nabla \Gamma_k(\mathbf{q}_k^{(s)}))^T (\mathbf{q}_k - \mathbf{q}_k^{(s)}) \\ &\geq \Gamma_k(\mathbf{q}_k), \end{aligned} \quad (49)$$

where $\nabla \Gamma_k(\mathbf{q}_k^{(s)})$ is the derivative of $\Gamma_k(\mathbf{q}_k)$ at $\mathbf{q}_k = \mathbf{q}_k^{(s)}$ and given by

$$\nabla \Gamma_k(\mathbf{q}_k^{(s)}) = \left[\frac{u_{k,1}^2}{(u_{k,1} + q_{k,1}^{(s)})^2}, \dots, \frac{u_{k,T}^2}{(u_{k,T} + q_{k,T}^{(s)})^2} \right]^T. \quad (50)$$

Then, under the s -th iteration of the MM algorithm, we need to solve the following problem

$$\begin{aligned} \text{(P1-}k\text{-}s\text{):} \quad & \text{Minimize}_{\{q_{k,i}\}} f(\mathbf{q}_k | \mathbf{q}_k^{(s)}) \\ \text{Subject to} \quad & B_k^{(G)}(\mathbf{q}_k) \leq \mu_k \bar{T}_{\text{fb}}. \end{aligned} \quad (51)$$

Note that problem (P1- k - s) is a convex optimization problem, and we can thus solve it globally based on the Lagrangian duality method. Specifically, the Lagrangian of problem (P1- k - s) is expressed as

$$L(\mathbf{q}_k, \eta) = f(\mathbf{q}_k | \mathbf{q}_k^{(s)}) + \eta (B_k^{(G)}(\mathbf{q}_k) - \mu_k \bar{T}_{\text{fb}}), \quad (52)$$

where $\eta \geq 0$ is the Lagrange multiplier associated with the constraint of problem (P1- k - s). The derivative of $L(\mathbf{q}_k, \eta)$ over $q_{k,i}$ is expressed as

$$\frac{\partial L(\mathbf{q}_k, \eta)}{\partial q_{k,i}} = -\frac{\eta u_{k,i}}{(u_{k,i} + q_{k,i}) q_{k,i} \ln 2} + \frac{u_{k,i}^2}{(u_{k,i} + q_{k,i}^{(s)})^2}, \quad \forall i. \quad (53)$$

By setting the derivative given in (53) as zero, it can be shown that given any $\eta \geq 0$, the Lagrangian given in (52) is minimized when

$$\bar{q}_{k,i}(\eta) = \frac{1}{2} \left[-u_{k,i} + \sqrt{u_{k,i}^2 + \frac{4\eta (q_{k,i}^{(s)} + u_{k,i})^2}{u_{k,i} \ln 2}} \right], \quad \forall k, i. \quad (54)$$

Moreover, it can be shown that under the optimal solution to problem (P1- k - s), the constraint should be satisfied with equality. Let η^* denote the optimal Lagrange multiplier to problem (P1- k - s). Then it follows that $\sum_{i=1}^T \log_2(1 + u_{k,i}/\bar{q}_{k,i}(\eta^*)) = \mu_k \bar{T}_{\text{fb}}$. The solution η^* to the above equation can be effectively obtained via the bisection method. Last, $\bar{q}_{k,i}(\eta^*)$'s will be the optimal solution to problem (P1- k - s).

After problem (P1- k - s) is solved at the s -th iteration of the MM algorithm, we set $\mathbf{q}_k^{(s)} = [\bar{q}_{k,1}(\eta^*), \dots, \bar{q}_{k,T}(\eta^*)]^T$ and solve problem (P1- k - $s+1$) for the $(s+1)$ -th iteration of the MM algorithm. As shown in [37], under the MM algorithm, the objective value of problem (P1- k) will decrease after each iteration, i.e., $\Gamma_k(\mathbf{q}_k^{(s+1)}) < \Gamma_k(\mathbf{q}_k^{(s)})$, $\forall s$. Then, the MM algorithm will converge to a locally optimal solution to problem (P1- k) [37].

Let $\mathbf{q}_k^* = [q_{k,1}^*, \dots, q_{k,T}^*]^T$ denote the solution to problem (P1- k) obtained via the above MM algorithm. Then, we show how to determine the size of codebook $\mathcal{C}_{k,i}$ in (5), i.e., $L_{k,i}$, $k = 1, \dots, K, i = 1, \dots, T$. According to (43), the theoretical number of bits to quantize $y_{k,i}$ is

$$B_{k,i}^{(G)*} = \log_2 \left(1 + \frac{u_{k,i}}{q_{k,i}^*} \right), \quad \forall k, i. \quad (55)$$

However, the above values may not be integer values. To satisfy the constraint in each problem (P1- k), in practice, we can set the number of quantization bits of user k at time sample i as

$$B_{k,i} = \lfloor B_{k,i}^{(G)*} \rfloor, \quad \forall k, i, \quad (56)$$

where $\lfloor \cdot \rfloor$ denotes the floor function. Thus, the size of $\mathcal{C}_{k,i}$ in (5) is given as

$$L_{k,i} = 2^{B_{k,i}}, \quad \forall k, i. \quad (57)$$

C. Tightness of the Gaussian Test Channel Approximation

It was analytically shown in [27] that the rate-distortion trade-off obtained from the Gaussian test channel model is a good approximation to that obtained from the Lloyd algorithm. In this sub-section, we provide a numerical example to verify this in our considered system.

In the numerical example, we assume that $D = 16$, $M = 12$ and $K = 11$. The pilot transmission overhead T is set as 32 time samples and feedback transmission overhead \bar{T}_{fb} is set as 64. Because quantization bit allocation can be performed independently over different users as shown in the previous sub-section, here we just focus on the rate-distortion trade-off of one user. In particular, we randomly generate 100 bit

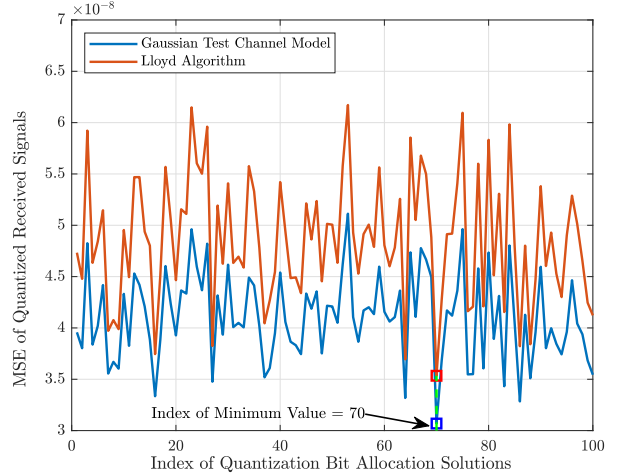


Fig. 2. MSE comparison between the Gaussian test channel model and Lloyd algorithm under quantization bit allocation solutions.

allocation solutions, each satisfying constraint (47). Given each quantization bit allocation solution, we first calculate the MSE obtained under the Gaussian test channel model according to (41), and then numerically calculate the channel estimation MSE obtained under the Lloyd algorithm based on Monte Carlo simulation.

The comparison between the above two MSEs under 100 quantization bit allocation solutions is given in Fig. 2. It is observed that the gap between the MSEs achieved by the Gaussian test channel model and the Lloyd algorithm is very small. More importantly, for any two quantization bit allocation solutions, if one solution leads to a smaller MSE compared to another solution under the Gaussian test channel model, then usually this solution also leads to smaller MSE under the Lloyd algorithm. Moreover, among the 100 quantization bit allocation solutions, the 70th quantization bit allocation solution achieves the minimum MSE under both the Gaussian test channel model and the Lloyd algorithm. To summarize, the rate-distortion performance achieved by the Gaussian test channel model is a good approximation to that achieved by the Lloyd algorithm, and it is thus feasible to obtain the quantization bit allocation policies of different users by solving problems (P1- k), $\forall k$. Note that similar observations have also been made in [27].

D. Performance Comparison Under Different Bit Allocation Solutions

In this sub-section, we show the gain of quantization bit allocation under our proposed algorithm. Two heuristic schemes are considered as the benchmark schemes. The first one is the even bit allocation scheme. Under this scheme, given the feedback overhead constraint \bar{T}_{fb} , we set $B_{k,i} = \mu_k \bar{T}_{\text{fb}}/T$, $\forall k, i$. The second one is the random bit allocation scheme. Under this scheme, for each user k , $B_{k,i}$'s are randomly generated and then normalized to satisfy the feedback overhead constraint \bar{T}_{fb} .

Fig. 3 shows the quantized received signals' normalized MSE (NMSE) over feedback transmission overhead achieved

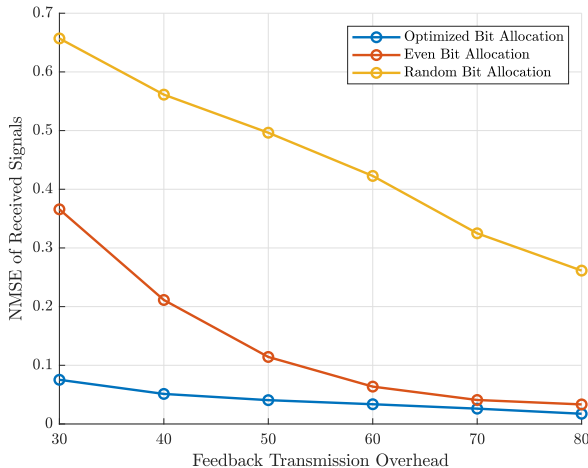


Fig. 3. NMSE performance of received signals under different bit allocation solutions.

by different quantization bit allocation schemes. The NMSE for quantizing the received signals is defined as

$$\text{NMSE}_y = \frac{\sum_{k=1}^K \mathbb{E} [\|\tilde{\mathbf{y}}_k - \mathbf{y}_k\|_2^2]}{\sum_{k=1}^K \mathbb{E} [\|\mathbf{y}_k\|_2^2]}, \quad (58)$$

where $\tilde{\mathbf{y}}_k = [\tilde{y}_{k,1}, \dots, \tilde{y}_{k,T}]^T$, $\forall k$. The numbers of antennas at the BS, IRS sub-surfaces, and users are set as $M = 16$, $D = 20$, and $K = 12$, respectively. The pilot transmission overhead T is set as 40 time samples and feedback transmission overhead \bar{T}_{fb} ranges from 30 to 80 time samples. It is observed that when the feedback overhead constraint is tight, our proposed quantization bit allocation scheme can achieve the best NMSE performance.

VI. NUMERICAL RESULTS

In this section, we provide numerical results to demonstrate the advantages of our proposed “quantize-then-estimate” scheme for IRS-assisted communication. The BS-IRS channel \mathbf{R} is modeled as a Rician fading channel with both the line-of-sight (LoS) deterministic component and the NLoS fading component:

$$\mathbf{R} = \sqrt{\frac{\kappa}{1+\kappa}} \mathbf{R}^{\text{LoS}} + \sqrt{\frac{1}{1+\kappa}} (\mathbf{C}^{\text{B}})^{\frac{1}{2}} \mathbf{R}^{\text{NLoS}} (\mathbf{C}^{\text{I}})^{\frac{1}{2}}, \quad (59)$$

where κ denotes the Rician factor set as 10dB. \mathbf{R}^{LoS} denotes the LoS component in \mathbf{R} , $\mathbf{C}^{\text{B}} \in \mathbb{C}^{M \times M} \succ \mathbf{0}$ and $\mathbf{C}^{\text{I}} \in \mathbb{C}^{D \times D} \succ \mathbf{0}$ denote the BS transmit correlation matrix and the IRS receive correlation matrix, respectively, and $\mathbf{R}^{\text{NLoS}} \sim \mathcal{CN}(\mathbf{0}, D\ell^{\text{BI}} \mathbf{I})$ denotes the i.i.d. Rayleigh fading component with ℓ^{BI} being the pass loss of \mathbf{R} . We assumed that \mathbf{C}^{B} and \mathbf{C}^{I} are generated based on the exponential correlation matrix model [1]. Next, the channel between the IRS and user k is modeled as $\mathbf{t}_k = (\mathbf{C}_k^{\text{I}})^{\frac{1}{2}} \hat{\mathbf{t}}_k$, where $\mathbf{C}_k^{\text{I}} \in \mathbb{C}^{D \times D} \succ \mathbf{0}$ denotes the IRS transmit correlation matrix for user k , which also follows the exponential correlation matrix model, $\forall k$, and $\hat{\mathbf{t}}_k \sim \mathcal{CN}(\mathbf{0}, \ell_k^{\text{IU}} \mathbf{I})$ follows the i.i.d. Rayleigh fading channel model with ℓ_k^{IU} denoting the pass loss. Moreover, the

pass loss of the BS-IRS channels \mathbf{r}_d 's and that of the IRS-user channels $t_{k,d}$'s are modeled as $\ell^{\text{BI}} = \ell_0 (\chi^{\text{BI}} / \chi_0)^{\xi_1}$ and $\ell_k^{\text{IU}} = \ell_0 (\chi_k^{\text{IU}} / \chi_0)^{\xi_2}$, $\forall k$, respectively, where ℓ_0 meter (m) denotes the reference distance, $\ell_0 = -20\text{dB}$ denotes the path loss at the reference distance, χ^{BI} and χ_k^{IU} denote the distance between the BS and the IRS, and that between the IRS and user k , respectively, and $\xi_1 = 2.2$ and $\xi_2 = 2.1$ denote the path loss factors for \mathbf{r}_d 's and $t_{k,d}$'s, respectively. The distance between the BS and the IRS is set to be 100 meters (m), and all the users are located in a circular region, whose center is 10 m away from the IRS and 105 m away from the BS and radius is 5 m. The power spectrum density of the noise at the users is assumed to be -169 dBm/Hz, and the channel bandwidth is 1 MHz. For the feedback transmission, it is assumed that each user employs 16QAM to modulate quantization bits, i.e., $\mu_k = 4$, $\forall k$. Last, we use the NMSE as the metric to evaluate the channel estimation performance. Specifically, we define $\mathbf{H}_k = [\mathbf{g}_{k,1}^T, \dots, \mathbf{g}_{k,D}^T]^T$ as the collection of user k 's cascaded channels, and $\hat{\mathbf{H}}_k = [\hat{\mathbf{g}}_{k,1}^T, \dots, \hat{\mathbf{g}}_{k,D}^T]^T$ as the collection of their estimations, $k = 1, \dots, K$. Then, the overall NMSE for estimating all the cascaded channels is defined as

$$\text{NMSE} = \frac{\sum_{k=1}^K \mathbb{E} [\|\hat{\mathbf{H}}_k - \mathbf{H}_k\|_2^2]}{\sum_{k=1}^K \mathbb{E} [\|\mathbf{H}_k\|_2^2]}. \quad (60)$$

In each numerical example, we conduct Monte Carlo simulation via 5000 channel realizations to numerically obtain the NMSE performance.

In the following, we provide two benchmark schemes for channel estimation and feedback under our considered IRS-assisted systems and compare the performance of our proposed scheme over that of the benchmark schemes.

- **Benchmark Scheme 1:** The first benchmark scheme is the conventional “estimate-then-quantize” scheme introduced in Section III-A. Under this scheme, each user k applies the LMMSE technique on its received pilot signals given in (2) to estimate its own cascaded channels, denoted by $\hat{\mathbf{g}}_{k,d}$'s, and feeds back the estimated channels to the BS with codebook designed via the Lloyd algorithm.
- **Benchmark Scheme 2:** The second benchmark scheme is an improved “estimate-then-quantize” scheme. Under this scheme, each user k still applies LMMSE technique to estimate its cascaded channels, denoted by $\hat{\mathbf{g}}_{k,d} = [\hat{g}_{k,d,1}, \dots, \hat{g}_{k,d,M}]^T$'s. However, for each IRS sub-surface d , only the reference user k_d feeds back $\hat{\mathbf{g}}_{k_d,d}$ to the BS, and each user k (including k_d) quantizes the sum of the estimated cascaded channel, i.e., $\sum_{m=1}^M \hat{g}_{k,d,m}$, and transmits the quantization bits to the BS. Then, the BS can estimate $\bar{\lambda}_{k,d}$ for $k \neq k_d$ as

$$\bar{\lambda}_{k,d} = \frac{\sum_{m=1}^M \bar{g}_{k,d,m}}{\sum_{m=1}^M \bar{g}_{k_d,d,m}}, \quad \forall d. \quad (61)$$

At last, the cascaded channels can be estimated based on (3). Note that compared to Benchmark Scheme 1, the main difference is that if a user is not a reference user, it merely feeds back the sum of its estimated channels

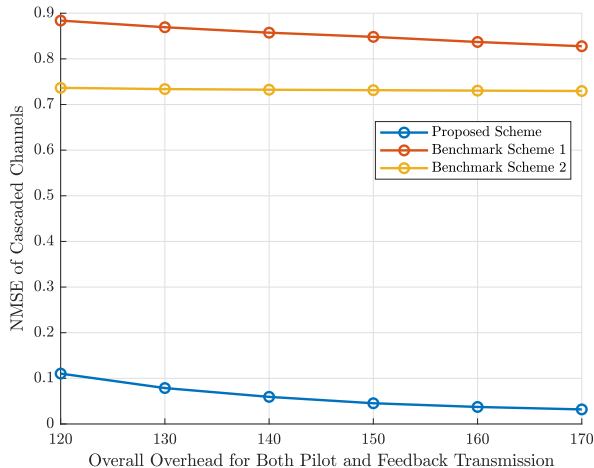


Fig. 4. NMSE performance versus overall overhead of pilot and feedback transmission when overhead of pilot transmission is fixed as 40 time samples.

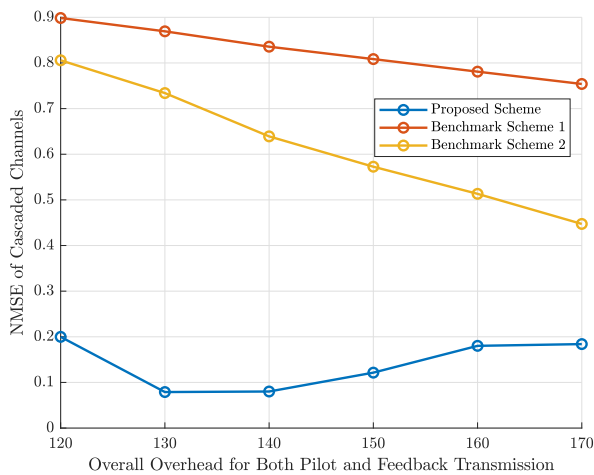


Fig. 5. NMSE performance versus overall overhead of pilot and feedback transmission when overhead of feedback transmission is fixed as 90 time samples.

thanks to (3). This can greatly reduce the feedback overhead. To characterize the overall overhead of Benchmark Scheme 2, let $\mathcal{S}_k = \{d : \forall d \text{ such that } k_d = k\}$ denote the set of sub-surfaces whose reference user is user k based on criterion (23), and $s_k = |\mathcal{S}_k|$ denote the number of times that user k is selected as the reference user, $\forall k$. Therefore, each user k needs to transmit $M s_k$ samples about $\bar{\mathbf{g}}_{k,d}$, $\forall d \in \mathcal{S}_k$, and D samples about $\bar{\lambda}_{k,d}$, $\forall d$. Then, the number of time samples for feedback transmission can be expressed as

$$T_{\text{fb,ben2}} = \max_{1 \leq k \leq K} \left\{ \frac{(M s_k + D) \lceil \log_2 L_k \rceil}{\mu_k} \right\}, \quad (62)$$

and the minimum number of overall time samples for pilot and feedback transmission is thus

$$T_{\text{tot,ben2}} = MD + T_{\text{fb,ben2}}. \quad (63)$$

In the numerical examples, the numbers of antennas at the BS, IRS sub-surfaces, and users are set as $M = 12$, $D = 16$,

and $K = 11$, respectively. The BS transmit power is 33 dBm for all the time slots. Fig. 4 shows the NMSE performance comparison between our proposed “quantize-then-estimate” scheme and the two benchmark schemes under the “estimate-then-quantize” approach. In this numerical example, the pilot transmission overhead is fixed as 40 time samples, while the feedback transmission overhead ranges from 80 to 130 time samples such that the overall overhead for pilot and feedback transmission ranges from 120 to 170 time samples. It is observed that our proposed scheme shows a significant NMSE performance gain compared to the two benchmark schemes. For example, when the overall overhead is 170 samples, the NMSE achieved by our proposed scheme is 0.0319, which is much better than that achieved by the two benchmark schemes. This is because (3) is leveraged to reduce both pilot and feedback transmission overhead. Note that Benchmark Scheme 2 shows a better performance than Benchmark Scheme 1 thanks to the reduction of feedback transmission overhead by exploiting (3). However, the performance of Benchmark Scheme 2 is much worse than that of our proposed scheme. This is because Benchmark Scheme 2 cannot utilize (3) to reduce the pilot transmission overhead. Specifically, given $M = 12$, $D = 16$, and $K = 11$, the minimum numbers of time samples required by Benchmark Scheme 2 and our proposed scheme are 192 and 32, respectively. Therefore, when the pilot transmission overhead is fixed as 40 time samples, the users cannot accurately estimate their cascaded channels under Benchmark Scheme 2.

Next, we set the feedback transmission overhead as 90 time samples, and the pilot transmission overhead ranges from 30 to 80 time samples, i.e., the overall overhead ranges from 120 to 170 time samples. Fig. 5 shows the performance comparison among different schemes. Under our proposed scheme, it is observed that as the overall overhead increases, the NMSE shows a “first-drop-then-rise” trend, while the optimal NMSE is achieved when the pilot transmission overhead is 40 time samples, i.e., the overall overhead is 130 time samples. Note that as the pilot transmission overhead increases, the BS can estimate the channels based on more received signals, but needs to quantize more pilot signals given the number of quantization bits. When the pilot transmission overhead is small, pilot transmission is the bottleneck to limit the channel estimation performance at the BS, and it is beneficial to increase the pilot transmission overhead. However, when the pilot transmission overhead is large enough, the BS has enough pilot signals to estimate the channels, and it is not a good idea to keep increasing the pilot transmission overhead because this will reduce the number of bits to quantize each pilot sample. This indicates that the pilot transmission overhead should be carefully designed. For the two benchmarks, increasing pilot transmission overhead always results in a decreasing NMSE because when the users feed back their channels, the amount of feedback is independent of the number of pilot signals, and the users can quantize better-estimated channels given the same number of quantization bits.

Moreover, Fig. 6 shows the estimated channels’ NMSE versus the number of BS antennas, with $D = 16$ and $K = 11$. The pilot transmission overhead is 64 time samples and the

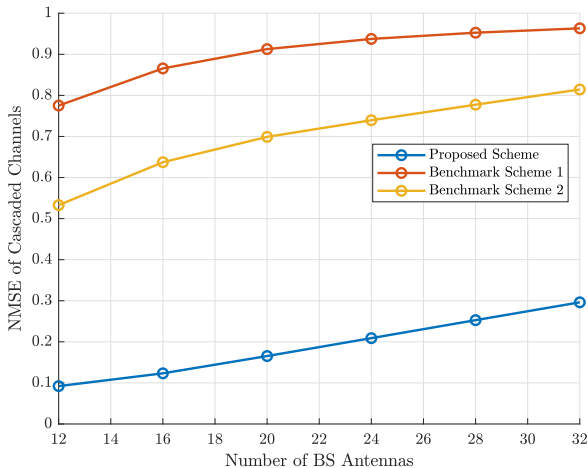


Fig. 6. NMSE performance versus number of BS antennas.

feedback transmission overhead is 100 time samples. It is observed that when the number of BS antennas increases, the NMSE increases much more slowly under our proposed scheme than the two benchmark schemes. This is attributed to exploiting (3) to reduce both pilot and feedback transmission overhead. Specifically, every additional BS antenna increases $D/K \approx 1.4545$ time samples of pilot and feedback transmission overhead, respectively, under our proposed scheme. While it causes 16 time samples of pilot and feedback transmission overhead, respectively, under Benchmark Scheme 1; and $D = 16$ time samples of pilot transmission overhead and $D/K \approx 1.4545$ time samples of feedback transmission overhead under Benchmark Scheme 2.

Last, Fig. 7 shows the estimated channels' NMSE of our proposed scheme versus different numbers of IRS sub-surfaces and signal-to-noise ratios (SNRs), where $M = 16$, $K = 11$, the pilot transmission overhead is set as 80 time samples, and the feedback transmission overhead is set as 160 time samples. It is observed that over different SNRs and numbers of IRS sub-surfaces, our proposed scheme performs better than the benchmark schemes.

VII. CONCLUSION

In this paper, we studied downlink CSI acquisition in FDD IRS-assisted communication systems. Motivated by the correlated channels among different users, we proposed a novel “quantize-then-estimate” protocol for reducing the overhead in both pilot transmission and feedback transmission. Specifically, all the users first quantize their received pilot signals and then transmit the quantization bits to the BS. After de-quantizing all the user's received signals, the BS can thus leverage the correlation embedded in users' cascaded channels to perform channel estimation. We designed efficient methods for each user to allocate the quantization bits over time and quantize the signals based on the carefully devised codebook, and for the BS to perform the LMMSE technique for estimating the channels based on the quantized signals. Moreover, we analytically characterized the minimum overhead for pilot transmission and feedback transmission under

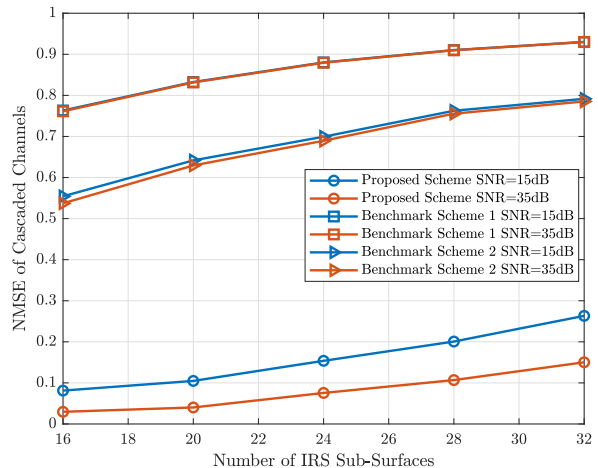


Fig. 7. NMSE performance versus SNR and number of IRS sub-surfaces.

our proposed “quantize-then-estimate” protocol and demonstrated the significant overhead reduction compared to the conventional “estimate-then-quantize” protocol. Our results open up a new solution for low-overhead communication in IRS-assisted systems.

APPENDIX

We prove the theorem in two cases: 1) $K \geq D$ and 2) $K < D$. In the case of $K \geq D$, we first prove that there exists a unique solution to (30) only if $\tau_2 \geq M - 1$. Define

$$\eta_{d,i} = \sqrt{p} \mathbf{g}_{k_d,d}^T \mathbf{x}_i, \quad d = 1, \dots, D, i = \tau_1 + 1, \dots, \tau_1 + \tau_2. \quad (64)$$

Then the received pilot signals of user k at time sample i can be expressed as

$$\tilde{y}_{k,i} = \sum_{d=1}^D \phi_{d,i} \lambda_{k,d} \eta_{d,i}, \quad \forall k, i = \tau_1 + 1, \dots, \tau_1 + \tau_2, \quad (65)$$

Define $\boldsymbol{\lambda} = [\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_D]$, then the overall received pilots at time sample i can be re-written as

$$\tilde{\mathbf{y}}_i^{(2)} = [\tilde{y}_{1,i}, \dots, \tilde{y}_{K,i}]^T = \boldsymbol{\lambda} \boldsymbol{\Gamma}_i, \quad i = \tau_1 + 1, \dots, \tau_1 + \tau_2 \quad (66)$$

where $\boldsymbol{\Gamma}_i = [\phi_{1,i} \eta_{1,i}, \dots, \phi_{D,i} \eta_{D,i}]^T$. We set $\phi_{d,i} = 1, \forall d, i$, then $\boldsymbol{\Gamma}_i = [\eta_{1,i}, \dots, \eta_{D,i}]^T$. With $\lambda_{k,d}$'s estimated in Step 1, there exist D variables $\eta_{d,i}$'s and K linear equations as given in (66). As a result, in the case of $K \geq D$, $\eta_{d,i}$'s can be perfectly estimated. Then, with the knowledge of $\eta_{d,i}$'s and $\alpha_{k_d,d}$'s, we can estimate $\mathbf{g}_{k_d,d}$'s based on (16) and (64) from the following equations

$$\begin{aligned} & [\alpha_{k_d,d}, \eta_{d,\tau_1+1}, \dots, \eta_{d,\tau_1+\tau_2}]^T \\ &= \sqrt{p} [\mathbf{x}, \mathbf{x}_{\tau_1+1}, \dots, \mathbf{x}_{\tau_1+\tau_2}]^T \mathbf{g}_{k_d,d}, \\ & d = 1, \dots, D, \end{aligned} \quad (67)$$

which characterizes a linear system with MD variables and $(\tau_2+1)D$ equations. Therefore, a unique solution to (67) exists

only when the number of equations is no smaller than the number of variables, i.e. $\tau_2 \geq M - 1$.

Next, we show that if $\tau_2 = M - 1$, there always exists a unique solution to (30) in the case of $K \geq D$. Specifically, since λ_d 's are linearly independent with each other with probability one, $\eta_{d,i}$'s can be perfectly estimated based on (66) as

$$\Gamma_i = \lambda^\dagger \tilde{\mathbf{y}}_i^{(2)}, \quad i = \tau_1 + 1, \dots, \tau_1 + M - 1. \quad (68)$$

Then we set $\mathbf{x} = [1, \dots, 1]^T$, and $\mathbf{x}_{\tau_1+1}, \dots, \mathbf{x}_{\tau_1+M-1}$ as the 2 to M columns of a $M \times M$ DFT matrix. With the knowledge of $\eta_{d,i}$'s and $\alpha_{k_d,d}$'s, $d = 1, \dots, D, i = \tau_1 + 1, \dots, \tau_1 + M - 1$, there exists a unique solution to (67), equivalently to (30) given as follows

$$\mathbf{g}_{k_d,d} = \sqrt{p} [\mathbf{x}, \mathbf{x}_{\tau_1+1}, \dots, \mathbf{x}_{\tau_1+M-1}]^* \times [\alpha_{k_d,d}, \eta_{d,\tau_1+1}, \dots, \eta_{d,\tau_1+M-1}]^T, \quad d = 1, \dots, D. \quad (69)$$

In the case of $K < D$, since the number of variables and equations in (30) are MD and $\tau_2 K + D$, respectively, there exists a unique solution to (30) only if the number of equations is no smaller than that of variables, i.e., $\tau_2 \geq \lceil \frac{(M-1)D}{K} \rceil$.

Next, we show that when $\tau_2 = \lceil \frac{(M-1)D}{K} \rceil$, there always exists a solution to (30) in the case of $K < D$. Specifically, we first set the pilot signal in Step 1 as an all-one vector, i.e., $\mathbf{x} = [1, \dots, 1]^T$, and the M -th pilot signal equal to zero in Step 2, i.e., $\mathbf{x}_{i,M} = 0, i = \tau_1 + 1, \dots, \tau_1 + \tau_2$. Then the other pilot signals, i.e., $\mathbf{x}_{i,m}, m = 1, \dots, M - 1$, as well as the IRS reflecting coefficients, i.e., $\phi_{d,i}, d = 1, \dots, D, i = \tau_1 + 1, \dots, \tau_1 + \tau_2$, are set in the same way as Theorem 2 in [1]. Then, we construct a new matrix $\bar{\Theta} \in \mathbb{C}^{(\tau_2 K + D) \times MD}$ by putting the $[(d-1)M + m]$ -th column of Θ into the $[(m-1)D + d]$ -th column of $\bar{\Theta}, \forall m, d$. Since changing the order of columns of a matrix does not change its rank, i.e., $\text{rank}(\bar{\Theta}) = \text{rank}(\Theta)$, in the following, we show that under the above construction, we have $\text{rank}(\bar{\Theta}) = MD$ when $\tau_2 = \lceil \frac{(M-1)D}{K} \rceil$. Specifically, $\bar{\Theta}$ can be re-expressed as follows:

$$\bar{\Theta} = \begin{bmatrix} \bar{\Theta}_s & \mathbf{O}_{\tau_2 K \times D} \\ \{\mathbf{I}_D\}_{M-1} & \mathbf{I}_D \end{bmatrix}, \quad (70)$$

where $\bar{\Theta}_s$ is the first $\tau_2 K$ rows and first $(M-1)D$ columns of $\bar{\Theta}$, $\mathbf{O}_{\tau_2 K \times D}$ is an all-zero matrix with dimension $\tau_2 K \times D$, \mathbf{I}_D is the identity matrix of dimension D , and $\{\mathbf{I}_D\}_{M-1} = [\mathbf{I}_D, \dots, \mathbf{I}_D] \in \mathbb{C}^{D \times (M-1)D}$. According to Theorem 2 in [1], $\text{rank}(\bar{\Theta}_s) = (M-1)D$ when $\tau_2 = \lceil \frac{(M-1)D}{K} \rceil$. Next, we derive the rank of $\bar{\Theta}$. It is observed from (70) that each of the first $\tau_2 K$ rows of $\bar{\Theta}$, whose last D elements are all zero, is linearly independent of the last D rows of $\bar{\Theta}$, i.e., $\{\mathbf{I}_D\}_M$. In other words, the row space defined by the first $\tau_2 K$ rows in $\bar{\Theta}$ does not intersect with that defined by the last D rows in $\bar{\Theta}$. In this case, $\text{rank}(\bar{\Theta}) = \text{rank}([\bar{\Theta}_s \ \mathbf{O}_{\tau_2 K \times D}]) + \text{rank}(\{\mathbf{I}_D\}_M) = MD$ [38]. Therefore, for the case of $K < D$, when $\tau_2 = \lceil \frac{(M-1)D}{K} \rceil$, there exists a unique solution to (30) given by

$$\mathbf{g} = \Theta^\dagger \tilde{\mathbf{y}}^{(2)}. \quad (71)$$

Theorem 1 is thus proved.

REFERENCES

- [1] Z. Wang, L. Liu, and S. Cui, "Channel estimation for intelligent reflecting surface assisted multiuser communications: Framework, algorithms, and analysis," *IEEE Trans. Wireless Commun.*, vol. 19, no. 10, pp. 6607–6620, Oct. 2020.
- [2] R. Wang, Z. Wang, L. Liu, S. Zhang, and S. Jin, "A quantize-then-estimate protocol for CSI acquisition in IRS-aided downlink communication," in *Proc. IEEE Global Commun. Conf. (Globecom)*, Dec. 2023.
- [3] 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. 116 753–116 773, Aug. 2019.
- [4] Q. Wu and R. Zhang, "Towards smart and reconfigurable environment: intelligent reflecting surface aided wireless network," *IEEE Commun. Mag.*, vol. 58, no. 1, pp. 106–112, Jan. 2020.
- [5] Q. Wu, S. Zhang, B. Zheng, C. You, and R. Zhang, "Intelligent reflecting surface-aided wireless communications: A tutorial," *IEEE Trans. Commun.*, vol. 69, no. 5, pp. 3313–3351, May 2021.
- [6] B. Zheng, C. You, W. Mei, and R. Zhang, "A survey on channel estimation and practical passive beamforming design for intelligent reflecting surface aided wireless communications," *IEEE Commun. Surv. Tutor.*, vol. 24, no. 2, pp. 1035–1071, Feb. 2022.
- [7] N. Jindal, "MIMO broadcast channels with finite-rate feedback," *IEEE Trans. Inf. Theory*, vol. 52, no. 11, pp. 5045–5060, Nov. 2006.
- [8] D. J. Love, R. W. Heath, V. K. Lau, D. Gesbert, B. D. Rao, and M. Andrews, "An overview of limited feedback in wireless communication systems," *IEEE J. Sel. Areas Commun.*, vol. 26, no. 8, pp. 1341–1365, Aug. 2008.
- [9] W. Shen, L. Dai, B. Shim, Z. Wang, and R. W. Heath, "Channel feedback based on AoD-adaptive subspace codebook in FDD massive MIMO systems," *IEEE Trans. Commun.*, vol. 66, no. 11, pp. 5235–5248, Nov. 2018.
- [10] C.-K. Wen, W.-T. Shih, and S. Jin, "Deep learning for massive MIMO CSI feedback," *IEEE Wireless Commun. Lett.*, vol. 7, no. 5, pp. 748–751, Mar. 2018.
- [11] J. Chen, Y.-C. Liang, H. V. Cheng, and W. Yu, "Channel estimation for reconfigurable intelligent surface aided multi-user MIMO systems," *IEEE Trans. Wireless Commun.*, vol. 22, no. 10, pp. 6853–6869, Feb. 2023.
- [12] 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, Nov. 2021.
- [13] X. Wei, D. Shen, and L. Dai, "Channel estimation for RIS assisted wireless communications: Part II - an improved solution based on double-structured sparsity," *IEEE Commun. Lett.*, vol. 25, no. 5, pp. 1403–1407, May 2021.
- [14] Z. Peng, C. Pan, G. Zhou, H. Ren, S. Jin, P. Popovski, R. Schober, and X. You, "Two-stage channel estimation for RIS-aided multiuser mmWave systems with reduced error propagation and pilot overhead," *IEEE Trans. Signal Process.*, vol. 71, pp. 3607–3622, Sep. 2023.
- [15] Q.-U.-A. Nadeem, H. Alwazani, A. Kammoun, A. Chaaban, M. Debbah, and M.-S. Alouini, "Intelligent reflecting surface-assisted multi-user MISO communication: Channel estimation and beamforming design," *IEEE Open J. Commun. Soc.*, vol. 1, pp. 661–680, May 2020.
- [16] 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.*, vol. 69, no. 6, pp. 4144–4157, Jun. 2021.
- [17] G. Zhou, C. Pan, H. Ren, P. Popovski, and A. L. Swindlehurst, "Channel estimation for RIS-aided multiuser millimeter-wave systems," *IEEE Trans. Signal Process.*, vol. 70, pp. 1478–1492, Sep. 2022.
- [18] B. Zheng, C. You, and R. Zhang, "Intelligent reflecting surface assisted multi-user OFDMA: Channel estimation and training design," *IEEE Trans. Wireless Commun.*, vol. 19, no. 12, pp. 8315–8329, Dec. 2020.
- [19] W. Chen, C.-K. Wen, X. Li, and S. Jin, "Adaptive bit partitioning for reconfigurable intelligent surface assisted FDD systems with limited feedback," *IEEE Trans. Wireless Commun.*, vol. 21, no. 4, pp. 2488–2505, Apr. 2021.
- [20] X. Ge, S. Yu, W. Shen, C. Xing, and B. Shim, "Beamforming design with partial channel estimation and feedback for FDD RIS-assisted systems," *IEEE Trans. Wireless Commun.*, Nov. 2023.
- [21] W. Chen, C.-K. Wen, X. Li, M. Matthaiou, and S. Jin, "Channel customization for limited feedback in RIS-assisted FDD systems," *IEEE Trans. Wireless Commun.*, vol. 22, no. 7, pp. 4505–4519, Jul. 2023.
- [22] J. Guo, W. Chen, C.-K. Wen, and S. Jin, "Deep learning-based two-timescale CSI feedback for beamforming design in RIS-assisted com-

munications," *IEEE Trans. Veh. Technol.*, vol. 72, no. 4, pp. 5452–5457, Apr. 2023.

- [23] D. Shen and L. Dai, "Dimension reduced channel feedback for reconfigurable intelligent surface aided wireless communications," *IEEE Trans. Commun.*, vol. 69, no. 11, pp. 7748–7760, Nov. 2021.
- [24] X. Shi, J. Wang, and J. Song, "Triple-structured sparsity-based channel feedback for RIS-assisted MU-MIMO system," *IEEE Commun. Lett.*, vol. 26, no. 5, pp. 1141–1145, Jan. 2022.
- [25] Z. Xiong, A. D. Liveris, and S. Cheng, "Distributed source coding for sensor networks," *IEEE Signal Process. Mag.*, vol. 21, no. 5, pp. 80–94, Sep. 2004.
- [26] Z. Wang, L. Liu, S. Zhang, and S. Cui, "Massive MIMO communication with intelligent reflecting surface," *IEEE Trans. Wireless Commun.*, vol. 22, no. 4, pp. 2566–2582, 2022.
- [27] P. Xia and G. B. Giannakis, "Design and analysis of transmit-beamforming based on limited-rate feedback," *IEEE Trans. Signal Process.*, vol. 54, no. 5, pp. 1853–1863, May 2006.
- [28] Y. Yang, B. Zheng, S. Zhang, and R. Zhang, "Intelligent reflecting surface meets OFDM: Protocol design and rate maximization," *IEEE Trans. Commun.*, vol. 68, no. 7, pp. 4522–4535, Jul. 2020.
- [29] X. Rao and V. K. Lau, "Distributed compressive CSIT estimation and feedback for FDD multi-user massive MIMO systems," *IEEE Trans. Signal Process.*, vol. 62, no. 12, pp. 3261–3271, Dec. 2014.
- [30] D. Mishra and H. Johansson, "Channel estimation and low-complexity beamforming design for passive intelligent surface assisted MISO wireless energy transfer," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, May 2019.
- [31] H. Alwazani, A. Kammoun, A. Chaaban, M. Debbah, M.-S. Alouini *et al.*, "Intelligent reflecting surface-assisted multi-user MISO communication: Channel estimation and beamforming design," *IEEE Open J. Commun. Soc.*, vol. 1, pp. 661–680, May 2020.
- [32] A. Gersho and R. M. Gray, *Vector quantization and signal compression*. Springer Science & Business Media, 1992.
- [33] G. Caire, N. Jindal, M. Kobayashi, and N. Ravindran, "Multiuser MIMO achievable rates with downlink training and channel state feedback," *IEEE Trans. Inf. Theory*, vol. 56, no. 6, pp. 2845–2866, Jun. 2010.
- [34] R. Wang, L. Liu, S. Zhang, and C. Yu, "A new channel estimation strategy in intelligent reflecting surface assisted networks," in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, Dec. 2021.
- [35] H. Pishro-Nik, *Introduction to probability, statistics, and random processes*. Blue Bell, PA, USA: Kappa Research, LLC, 2014.
- [36] A. El Gamal and Y.-H. Kim, *Network information theory*. Cambridge university press, 2011.
- [37] Y. Sun, P. Babu, and D. P. Palomar, "Majorization-minimization algorithms in signal processing, communications, and machine learning," *IEEE Trans. Signal Process.*, vol. 65, no. 3, pp. 794–816, Mar. 2016.
- [38] G. Matusaglia and G. P. H. Styan, "Equalities and inequalities for ranks of matrices," *Linear and Multilinear Algebra*, vol. 2, no. 3, pp. 269–292, Apr. 1974.



Rui Wang (Graduate Student Member, IEEE) received the B.S. degree in communication engineering from University of Electronic Science and Technology of China, Chengdu, China, in 2017 and the M.S. degree in information and communication engineering from Southeast University, Nanjing, China, in 2020. He is currently working towards the Ph.D. degree at the Department of Electrical and Electronic Engineering, The Hong Kong Polytechnic University, Hong Kong. His research interests include intelligent reflecting surface assisted wireless

communication, channel estimation and MIMO systems.



Zhaorui Wang (Member, IEEE) received the Ph.D. degree from The Chinese University of Hong Kong (CUHK) in 2019, and the B.S. degree from University of Electronic Science and Technology of China (UESTC) in 2015. He is now a Research Assistant Professor at the School of Science and Engineering, The Chinese University of Hong Kong, Shenzhen. He was a Postdoctoral Research Associate at The Hong Kong Polytechnic University from 2019 to 2020, and a Postdoctoral Research Associate at CUHK from 2021 to 2022. He is a recipient of the

Hong Kong PhD Fellowship from 2015 to 2018. He has been selected in the post of "Pengcheng Peacock Plan" (Type C) since 2022. His research interests include system-level design on massive machine-type communications (mMTC), and semantic communications.



Liang Liu (Senior Member, IEEE) received the B.Eng. degree from Tianjin University, China, in 2010, and the Ph.D. degree from the National University of Singapore, Singapore, in 2014. From 2015 to 2017, he was a Post-Doctoral Fellow with the Department of Electrical and Computer Engineering, University of Toronto. From 2017 to 2018, he was a Research Fellow with the Department of Electrical and Computer Engineering, National University of Singapore. He is currently an Associate Professor with the Department of Electrical and Electronic

Engineering, The Hong Kong Polytechnic University. His research interests lie in the next generation cellular technologies such as machine-type communications for the Internet of Things, integrated sensing and communication, etc. He was a recipient of the 2021 IEEE Signal Processing Society Best Paper Award, the 2017 IEEE Signal Processing Society Young Author Best Paper Award, the Best Student Award of 2022 IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP), and the Best Paper Award of the 2011 International Conference on Wireless Communications and Signal Processing. He was recognized by Clarivate Analytics as a Highly Cited Researcher in 2018. He is an Editor of IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS. He was a Leading Guest Editor of IEEE WIRELESS COMMUNICATIONS Special Issue on Massive Machine-Type Communications for IoT. He is a co-author of the book "Next Generation Multiple Access" published at Wiley-IEEE Press.



Shuowen Zhang (Member, IEEE) received the B.Eng. degree in information engineering from Chien-Shiung Wu Honors College, Southeast University, Nanjing, China, in June 2013, and the Ph.D. degree from NUS Graduate School for Integrative Sciences and Engineering (NGS), National University of Singapore, in January 2018 under the NGS scholarship. From February 2018 to July 2020, she was a Research Fellow with the Department of Electrical and Computer Engineering, National University of Singapore. Since August 2020, she has

been with The Hong Kong Polytechnic University, where she is currently an Assistant Professor at the Department of Electrical and Electronic Engineering. Her research interests include integrated sensing and communications, intelligent reflecting surface aided communications, unmanned aerial vehicles, multiple-input multiple-output (MIMO) communications, communication theory, and optimization methods. Dr. Zhang is currently serving as an Editor for IEEE Transactions on Wireless Communications. She has served as a Guest Editor for various journals such as IEEE Journal on Selected Areas in Communications, and a technical program committee (TPC) member for various IEEE flagship conferences. She also served as an IEEE Communications Society Asia/Pacific Board WICE Vice Chair, and an IEEE N2Women Mentoring Co-Chair. Dr. Zhang is the sole recipient of the 2021 Marconi Society Paul Baran Young Scholar Award, as well as a recipient of the 2022 IEEE Communications Society Young Author Best Paper Award, 2023 IEEE Communications Society Best Tutorial Paper Award, 2023 PolyU Young Innovative Researcher Award, and 2024 IEEE Communications Society Asia Pacific Outstanding Young Researcher Award.



Shi Jin (Fellow, IEEE) received the B.S. degree in communications engineering from Guilin University of Electronic Technology, Guilin, China, in 1996, the M.S. degree from Nanjing University of Posts and Telecommunications, Nanjing, China, in 2003, and the Ph.D. degree in information and communications engineering from the Southeast University, Nanjing, in 2007. From June 2007 to October 2009, he was a Research Fellow with the Adastral Park Research Campus, University College London, London, U.K. He is currently with the faculty of the National

Mobile Communications Research Laboratory, Southeast University. His research interests include wireless communications, random matrix theory, and information theory. He is serving as an Area Editor for the IEEE TRANSACTIONS ON COMMUNICATIONS and IET Electronics Letters. He was an Associate Editor for the IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS, IEEE COMMUNICATIONS LETTERS, and IET Communications. Dr. Jin and his co-authors have been awarded the 2011 IEEE Communications Society Stephen O. Rice Prize Paper Award in the field of communication theory, the 2024 IEEE Communications Society Marconi Prize Paper Award, the IEEE Vehicular Technology Society 2023 Jack Neubauer Memorial Award, a 2022 Best Paper Award and a 2010 Young Author Best Paper Award by the IEEE Signal Processing Society.