1	Deep Learning-Based Channel Extrapolation and Multi-User
2	Beamforming for RIS-aided Terahertz Massive MIMO
3	Systems over Hybrid-Field Channels
4	Yang Wang <sup>1</sup> , Zhen Gao <sup>1</sup> , and Sheng Chen <sup><math>2,3</math></sup>
5	<sup>1</sup> School of Information and Electronics, Beijing Institute of Technology, Beijing,
6	China.
7	<sup>2</sup> School of Electronics and Computer Science, University of Southampton,
8	Southampton, U.K.
9	<sup>3</sup> Faculty of Information Science and Engineering, Ocean University of China,
10	Qingdao, China

## Abstract

The reconfigurable intelligent surface (RIS) is a promising new technology for Terahertz 12 (THz) massive multiple input multiple output (MIMO) communication systems. However, due 13 to the cascaded channel structure and its lack of signal processing ability, it is difficult for RIS 14 to obtain the high-dimensional channel state information (CSI) and optimize the active/passive 15 beamforming. Therefore, in this paper, we propose a DL-based transmission scheme for RIS-16 aided THz massive MIMO systems over hybrid far-near field channels, where a channel estima-17 tion scheme with low pilot overhead and a robust beamforming scheme are conceived. Specifi-18 cally, we first propose an end-to-end deep learning (DL)-based channel estimation framework, 19 which consists of pilot design, CSI feedback, sub-channel estimation, and channel extrapolation. 20 Specifically, we firstly only turn on a fraction of all the RIS elements and estimate a sub-sampling 21 RIS channel, and then design a DL-based scheme to extrapolate the full-dimensional CSI from 22 the partial one. Moreover, to maximize the sum rate of all users under imperfect CSI, we develop 23 a DL-based scheme to simultaneously design the hybrid active beamforming at the BS and pas-24 sive beamforming at the RIS. Simulation results show that our proposed channel extrapolation 25 scheme has better CSI reconstruction performance than conventional schemes while imposing a 26 much reduced pilot overhead and our proposed beamforming scheme has superior performance 27 over conventional schemes in terms of robustness to imperfect CSI. 28

## 29 Keywords

11

Reconfigurable intelligent surface (RIS), Terahertz (THz), hybrid-field, channel extrapolation, hybrid
 beamforming, orthogonal frequency division multiplexing (OFDM), multiple-input multiple-output

<sup>32</sup> (MIMO), deep learning (DL).

## 33 1 Introduction

Over the past few years, the demand for wireless data traffic has increased significantly due 34 to the explosive growth of mobile devices and multimedia applications [1]. To accommodate these 35 demands, Terahertz (THz) communications have attracted great interest from both industry and 36 academia. However, there exists strong atmospheric attenuation and extremely high free-space 37 losses in the THz band. These disadvantages may severely degrade the service coverage of THz 38 communication systems. The deployment of massive multiple input multiple output (MIMO) or 39 even ultra-massive MIMO has been recognized as an achievable solution to provide high array gain, 40 so as to overcome the high propagation loss [2]. However, conventional massive MIMO systems 41 with fully-digital architecture require a dedicated radio frequency (RF) chain for each antenna, 42 which results in excessive power consumption and extremely high RF hardware costs. In order to 43 circumvent this technical hurdle, the hybrid analog-digital massive MIMO architecture has been 44 widely adopted to achieve large array gains with a much lower number of RF chains [3]. 45

Besides, the emerging technology of reconfigurable intelligent surface (RIS) has also been considered as a promising technique to enhance communication performance [4]. The RIS can manipulate both the phase and amplitude of the incident electromagnetic (EM) signals so as to passively reflect them towards the desired directions and provide significant beamforming gain. More importantly, RIS does not need power-hungry RF chains, which is beneficial for developing green and cost-efficient communications. Therefore, the application of RIS and massive MIMO techniques is expected to overcome the limitations of THz communications and realize its full potential.

Generally, a simplified planar-wave channel model is appropriate in the case that the user equip-53 ment (UE) works in the far-field of the base station (BS). However, since severe path loss will reduce 54 the effective coverage while the increasing array size in THz band will increase the Rayleigh distance 55 [5], both far-field and near-field need to be considered for THz massive MIMO systems. Therefore, 56 the spherical-wave is necessary to accurately analyze the propagation of THz waves under near-field 57 conditions, where the distance from each antenna of the BS to the UE needs to be considered [6]. On 58 the other hand, in THz massive MIMO systems, the number of channel parameters is proportional to 59 the number of massive antennas, which indicates that directly applying the spherical-wave channel 60 model is unrealistic. To this end, a hybrid-field (hybrid spherical- and planar-wave) channel model 61 characterized by a smaller number of parameters while maintaining high accuracy has been proposed 62 for THz massive MIMO systems [7]. For such a channel model, given the subarray structure, the 63 inter-subarray is modeled as spherical-wave and the intra-subarray is modeled as planer-wave. Al-64 though the application of RISs has been widely investigated recently [8, 9, 10, 11, 12], the utilization 65 of RIS for THz massive MIMO communications over hybrid-field channels is still at its early study 66 stage. 67

## 68 1.1 Related Work

Acquiring accurate channel state information (CSI) is critical in establishing RIS-aided commu-69 nication systems. However, it is challenging to accurately estimate high-dimensional channels with 70 a few pilots [8]. By exploiting the sparsity of the channels in the angular and/or delay domains, 71 compressive sensing (CS)-based solutions have been proposed to reduce the pilot overhead [9, 10]. 72 Nevertheless, since the dimension of the CSI to be estimated is extremely large, the involved matrix 73 inversion operations and the iterative nature of CS-based techniques result in prohibitively high 74 computational complexity and storage requirements. With the development of artificial intelligence, 75 the application of deep learning (DL) in RIS-aided communication systems has attracted extensive 76 attention. In [11], the authors designed a twin convolutional neural network (CNN) to estimate the 77 direct (BS-UE) and the cascaded (BS-RIS-UE) channels. In [12], the authors proposed a hybrid 78 passive/active RIS architecture with a few RF chains, where the orthogonal match pursuit (OMP) 79 algorithm and denoising CNN are applied to reconstruct the complete channel matrix. However, 80 the deployment of RF chains at the RIS defeats the original purpose of reducing hardware cost and 81 power consumption by deploying the RIS. 82

In fact, due to the highly-dense placement of RIS elements, there is a strong correlation between 83 the different elements of the CSI matrix, which makes it possible to extrapolate the complete channel 84 from a partial one, i.e., channel extrapolation [13]. Recently, there are some initial attempts to utilize 85 the channel extrapolation for further reducing the pilot overhead. In [14], the authors proposed an 86 antenna domain extrapolation network to perform channel extrapolation, where an antenna selection 87 network is designed to choose the optimal antennas for the extrapolation. In [15], the neural network 88 structure modified by ordinary differential equation was used to describe the latent relations between 89 different data layers and improve the performance gains of antenna extrapolation. Besides, the 90 authors of [16] adopted the element-grouping strategy to reduce the pilot overhead and the CNN-91 based channel extrapolation network to extrapolate the full-dimensional cascaded channels from the 92 partial one. However, the above extrapolation schemes only consider the extrapolation process from 93 the known sub-channels, while ignoring how to estimate the sub-channel. Moreover, the hybrid-94 field channel modeling of RISs has more complex propagation characteristics, which will hinder the 95 sub-channel acquisition and the following extrapolation of complete channels. 96

How to properly and effectively design the hybrid beamforming and RIS phase according to the 97 CSI is one of the major engineering challenges in the design of RIS-aided communication systems. 98 Recently, some work has been conducted to investigate hybrid beamforming and RIS design problems 99 [17, 18, 19, 20]. The authors in [17] proposed a geometric mean decomposition-based beamform-100 ing for RIS-assisted mmWave hybrid MIMO systems. The authors in [18] focused on the hybrid 101 beamforming and RIS phase design for RIS-aided multi-user mmWave communication systems and 102 aimed to minimize the mean squared error (MSE) using the gradient-projection method. In [19], 103 the authors proposed an iterative algorithm to jointly optimize the RIS phase and the hybrid beam-104 forming for maximizing the sum rate. Furthermore, the DL-based beamforming methods have also 105 been studied in RIS-assisted wireless communication systems. In [20], a DL-based approach was 106 proposed to jointly optimize the active beamforming at the BS and the passive beamforming at the 107

RIS for achieving rate maximization. However, a further adaptability analysis of the aforementioned
 schemes is desired because the current analysis only assumes the idealized perfect CSI condition.

## 110 1.2 Motivations

In the current research, RIS is mainly constructed to have two modes: reflective mode [17, 18, 19] 111 and transmissive mode [21, 22, 23]. At present, more researches focus on RIS-assisted communication 112 in reflective mode. The RIS in reflective mode is mainly used to solve the blind coverage problem 113 and enhance the network coverage. By contrast, the transmissive RIS is mainly used to enhance 114 the spectral and energy efficiency of the networks. Since the transmissive mode does not change 115 the propagation direction of EM waves, it is suitable to deploy transmissive RIS in the case that 116 a line-of-sight (LoS) path exists but the propagation attenuation is high, e.g., the case that the 117 outdoor BS serves indoor UEs, to improve the energy of the received signals. Therefore, we consider 118 the transmissive RIS for indoor signal enhancement service. 119

Considering the hybrid-field channel model, the authors of [24] proposed a two-stage channel 120 estimation mechanism, where a CNN-based network is designed to estimate the parameters of the 121 reference subarrays and the complete channel is reconstructed by channel extrapolation based on ge-122 ometric relationships of channel parameters. However, this parametric-based extrapolation method 123 needs to obtain a large number of accurate channel parameters as labels before training. In [25], the 124 authors proposed a sensor-based channel estimation and beamforming for RIS-aided THz systems, 125 where a LoS MIMO architecture is considered in hybrid-field. However, the channel estimation in 126 [25] is mainly based on the awareness of sensors, and it is difficult to obtain accurate CSI. Therefore, 127 similar to [14, 15, 16], we adopt a DL-based channel extrapolation method to address the perfor-128 mance limitations of conventional channel estimation methods for indoor hybrid-field propagation 129 environments. Besides, in this paper, we consider the LoS MIMO architecture under the assump-130 tion of the hybrid-field channel model, where the LoS MIMO architecture can support multi-stream 131 transmission in the pure LoS BS-RIS channel. 132

Most existing works assume that the CSIs between BS and RIS as well as between RIS and UEs 133 are perfect [17, 18, 19, 20]. However, this assumption is impractical. Therefore, it is imperative 134 to take the channel estimation error into consideration when designing RIS-aided communication 135 systems. Recently, there have been some works on beamforming with imperfect CSI [26, 27]. Specif-136 ically, in [26], the authors exploited the penalty-based alternating optimization to design secure 137 wireless communications assisted by RIS under the imperfect CSI. In [27], the authors exploited a 138 gradient projection-based alternating optimization algorithm to jointly optimize transmit beamform-139 ing, RIS placement, and reflect beamforming of the RIS under the imperfect CSI. Currently, there 140 are extensive DL-based methods for the optimization of RIS phase with the perfect CSI, but to our 141 best knowledge, there exists only few DL-based methods considering imperfect CSI [28]. Therefore, 142 this work aims to provide a DL-based hybrid beamforming and RIS phase design with imperfect 143

<sup>144</sup> CSI in RIS-aided communication systems.

## 145 **1.3** Contributions

This paper proposes a DL-based spatial-frequency domain channel extrapolation (SFDCEtra) network as well as the DL-based hybrid beamforming and RIS phase design (HBFRPD) scheme for RIS-aided downlink multi-user THz massive MIMO systems over hybrid-field channels. The main contributions of this paper are summarized as follows.

• We deploy the transmissive RIS on the window of a building to reduce the penetration loss and thus achieve indoor enhanced communication. In addition, due to the negligible non-LoS (NLoS) component energy in the THz band, the BS-RIS channel is dominated by the LoS path. To achieve multi-stream transmission in the LoS case, we consider an LoS MIMO architecture under hybrid-field channel modeling, where the BS and RIS adopt the same subarray structures, and the subarray spacing is optimized to satisfy the LoS MIMO condition.

Since the BS and the RIS are fixed and only the LoS path exists, the BS-RIS channel can be considered to be quasi-static and it is known. The RIS-UE channel by contrast is time-varying due to the mobility of the UEs. Therefore, we only need to estimate the RIS-UE channel, which significantly reduces the pilot overhead.

• To further reduce the pilot overhead for estimating the RIS-UE channel, we propose a DL-based 160 channel extrapolation scheme, where the RIS only activates part of its elements at the channel 161 estimation stage. Unlike the existing extrapolation schemes [14, 15, 16] that only focus on the 162 CSI extrapolation process, we design a complete channel extrapolation framework, including 163 the pilot design network, CSI feedback network, sub-channel estimation network, and channel 164 extrapolation network. By adopting the end-to-end (E2E) training strategy, the proposed 165 channel estimation scheme can reduce the pilot overhead while maintaining high reconstruction 166 performance. Specifically, the UE-side feeds the quantized pilot information back to the BS 167 using the CSI feedback network, and the BS estimates the sub-channel and then extrapolates 168 the complete RIS-UE channel using the channel extrapolation network. In addition, for the 169 RIS element selection, we discuss the impact of three different strategies, uniform selection, 170 random selection and learning-based selection, on the final channel estimation performance. 171

• To solve the multi-user interference problem under imperfect CSI, we propose a DL-based hybrid beamforming and RIS phase design scheme, which consists of the analog beamformer design, DL-based RIS phase design network, and knowledge-data dual-driven digital beamforming network. By adopting the sum rate as the loss function to perform E2E training, the proposed scheme can realize higher performance and better robustness than the existing state-of-the-art methods.

Notations: This paper uses lower-case letters for scalars, lower-case boldface letters for vectors, and upper-case boldface letters for matrices. Superscripts  $(\cdot)^*$ ,  $(\cdot)^T$ ,  $(\cdot)^H$ ,  $(\cdot)^{-1}$  and  $(\cdot)^{\dagger}$  denote the conjugate, transpose, conjugate transpose, inversion, and Moore-Penrose inversion operators, respectively. The operators diag $(\cdot)$ , blkdiag $(\cdot)$ ,  $\otimes$  and  $\odot$  represent the diagonalization, block diagonalization, Kronecker product and Hadamard product, respectively, while  $\|\mathbf{A}\|_F$  denotes the Frobenius norm of **A**.  $\mathbf{I}_n$  denotes the identity matrix with size  $n \times n$ , while  $\mathbf{1}_n$  ( $\mathbf{0}_n$ ) denotes the column vector of size n with all the elements being 1 (0).  $\Re\{\cdot\}$  and  $\Im\{\cdot\}$  denote the real part and imaginary part of the corresponding argument, respectively.  $\{\mathbf{A}\}_{m,n}$  denotes the m-th row and n-th column element of **A**, and  $\{\mathbf{a}\}_m$  is the m-th entry of **a**, while  $\mathbf{A}_{[:,m:n]}$  is the sub-matrix containing the m-th to n-th columns of **A**. The expectation is denoted by  $\mathbb{E}(\cdot)$ , and  $\mathcal{N}(\mu, \sigma^2)$  ( $\mathcal{CN}(\mu, \sigma^2)$ ) denotes the real (complex) Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ , while  $\mathrm{Tr}\{\cdot\}$  represents the matrix trace operator.

## <sup>190</sup> 2 System Model

<sup>191</sup> We first describe the RIS-assisted downlink THz multi-user MIMO system over frequency-<sup>192</sup> selective fading channels, and then introduce the transmission channel model of this system.



Figure 1: Model of RIS-aided THz massive MIMO system: (a) multiple indoor UEs are served by the BS with the help of a transmissive RIS deployed on the window, and (b) hardware architectures at the BS, RIS, and UEs.

## <sup>193</sup> 2.1 System Description

As shown in Figure 1, we consider a downlink transmission scenario in an indoor environment, where a transparent RIS is attached to the window surface to refract outdoor THz signals from the BS

into the room for serving U single-antenna UEs. Thus, the transparent transmissive RIS helps to 196 enhance indoor coverage. Let the BS (RIS) have  $M^{\rm B} = M_y^{\rm B} \times M_z^{\rm B}$   $(M^{\rm R} = M_y^{\rm R} \times M_z^{\rm R})$  uniformly 197 spaced subarrays, where  $M_y^{\rm B}$   $(M_y^{\rm R})$  and  $M_z^{\rm B}$   $(M_z^{\rm R})$  are the numbers of subarrays along the azimuth 198 and elevation directions, respectively. Each subarray of the BS (RIS) is a uniform planar array 199 (UPA) with  $N_{\text{sub}}^{\text{B}} = N_{y}^{\text{B}} \times N_{z}^{\text{B}}$   $(N_{\text{sub}}^{\text{R}} = N_{y}^{\text{R}} \times N_{z}^{\text{R}})$  isotropically radiating elements, where  $N_{y}^{\text{B}}$   $(N_{y}^{\text{R}})$ 200 and  $N_z^{\rm B}$  ( $N_z^{\rm R}$ ) are the numbers of antennas along the azimuth and elevation directions, respectively. 201 Therefore, the complete antenna dimension of the BS is  $N^{\rm B} = M^{\rm B} N_{\rm sub}^{\rm B}$ , and the element dimension 202 of the RIS is  $N^{\rm R} = M^{\rm R} N_{\rm sub}^{\rm R}$ . Without loss of generality, we assume that the normals of the central 203 elements of both the BS and RIS are coaxial, i.e., meeting the parallel symmetric array arrangements, 204 with a distance of D, as illustrated in Figure 1(b). 205

To reduce the power consumption and hardware cost, a sub-connected hybrid analog-digital array architecture is considered at the BS, i.e., there are only  $M^{\rm B}$  RF chains to support  $U \leq M^{\rm B}$ data streams, and each of RF chains is connected to a subarray through  $N_{\rm sub}^{\rm B}$  phase shifters. An orthogonal frequency division multiplexing (OFDM) transmission scheme with K subcarriers and sampling period  $T_s = 1/f_s$  is adopted. The cyclic prefix (CP) of length  $N_{\rm CP}T_s$  is added before each OFDM symbol to avoid inter-symbol interference. The center-carrier frequency is  $f_c$  corresponding to a wavelength  $\lambda$ .

## 213 2.2 Channel Model

## 214 2.2.1 BS-RIS Channel Model

<sup>215</sup> By considering the spherical wave propagation of THz signals, we investigate the LoS MIMO <sup>216</sup> link between the BS and RIS with only the LoS path, but it can support intra-path multiplexing <sup>217</sup> for multi-stream transmission [29]. In THz channels, the path loss of the NLoS paths is known to <sup>218</sup> be much larger than that of the LoS path. Therefore, we can neglect the NLoS paths in the channel <sup>219</sup> between the BS and the RIS. The inter-antenna spacing in each subarray is  $d = \lambda/2$ . In order to <sup>220</sup> satisfy the LoS MIMO characteristic, the BS subarray spacing  $d_{sy}^{\rm B}$  and  $d_{sz}^{\rm B}$  are set to the following <sup>221</sup> optimal LoS MIMO spacing

$$d_{sy}^{\rm B} = \sqrt{\frac{\lambda D}{M_y^{\rm B}}} - \frac{\lambda}{2} (N_y^{\rm B} - 1), d_{sz}^{\rm B} = \sqrt{\frac{\lambda D}{M_z^{\rm B}}} - \frac{\lambda}{2} (N_z^{\rm B} - 1), \tag{1}$$

i.e.,  $d_{sy}^{\rm B}$  and  $d_{sz}^{\rm B}$  should satisfy the condition  $\lambda \ll d_{sy}^{\rm B}, d_{sz}^{\rm B} \ll D$ . The detailed explanation of Equation (1) can be found in [29, 30]. The RIS subarray spacing  $d_{sy}^{\rm R}$  and  $d_{sz}^{\rm R}$  can be obtained by using a similar definition. Note that self-orthogonal LoS MIMO not only is obtained from parallel symmetric antenna arrangements but also can be obtained with symmetrical/unsymmetrical arrangements on tilted non-parallel lines/planes [30]. We have the following proposition from [31].

**Proposition 1** Let the transceiver arrays be separated by D and be communicating at a carrier wavelength of  $\lambda$  that is much smaller than D, i.e.,  $\lambda \ll D$ . If the inter-antenna spacing is in the order of  $\mathcal{O}(\lambda)$ , the first-order approximation of the spherical wave by the planner wave model can be applied. Otherwise, the spherical wave model should be used. According to Proposition 1, the subarray response vectors  $\mathbf{a}(\theta, \phi, f_k) \in \mathbb{C}^{N_{\mathrm{H}}N_{\mathrm{V}} \times 1}$  can be approximated by a planner wave model:

$$\mathbf{a}(\theta,\phi,f_k) = \mathbf{a}_h(\theta,\phi,f_k) \otimes \mathbf{a}_v(\phi,f_k)$$
$$= \begin{bmatrix} 1, \dots, e^{-j2\pi \frac{f_k}{c}d(n_h\sin\theta\cos\phi + n_v\sin\phi)}, \dots, e^{-j2\pi \frac{f_k}{c}d((N_{\rm H}-1)\sin\theta\cos\phi + (N_{\rm V}-1)\sin\phi)} \end{bmatrix}^{\rm T}, \quad (2)$$

where  $f_k = f_c - \frac{f_s}{2} + \frac{kf_s}{K}$ ,  $1 \le k \le K$ , is the k-th subcarrier frequency, c is the speed of light,  $0 \le n_h \le (N_{\rm H} - 1), 0 \le n_v \le (N_{\rm V} - 1), N_{\rm H}$  and  $N_{\rm V}$  are the numbers of antennas along the azimuth and elevation directions in the subarray, respectively, while  $\theta$  and  $\phi$  are the azimuth and elevation angles of the departure or arrival (AoD or AoA) of the path, respectively.

Since  $d_{sy}^{\text{B}}, d_{sz}^{\text{B}}, d_{sz}^{\text{R}} \ll D$ , the direction difference of the same path in different subarrays is negligible, and all the subarrays on either the BS or RIS-side share the same set of array response vectors. However, the relative phase differences among subarrays are of non-negligible values, as subarrays are widely spaced. Motivated by the above analysis, the downlink spatial-frequency channel  $\mathbf{G}[k] \in \mathbb{C}^{N^{\text{R}} \times N^{\text{B}}}$  from the BS to the RIS on the k-th subcarrier can be modeled as

$$\mathbf{G}[k] = \alpha[k]G_{\mathrm{T}}\mathbf{G}[k] \otimes \left[\mathbf{a}_{\mathrm{R}}(\theta_{\mathrm{R},\mathrm{A}},\phi_{\mathrm{R},\mathrm{A}},f_{k})\mathbf{a}_{\mathrm{B}}^{\mathrm{H}}(\theta_{\mathrm{B}},\phi_{\mathrm{B}},f_{k})\right],\tag{3}$$

where  $\alpha[k]$  is the channel attenuation coefficient on the k-th subcarrier,  $(\theta_B, \phi_B)$  and  $(\theta_{R,A}, \phi_{R,A})$  are the LoS AoD and LoS AoA of the BS-RIS channel, respectively. Without loss of generality, the LoS angles between the BS and the RIS are assumed to be fixed and known in advance based on their locations. In (3), the entries of  $\tilde{\mathbf{G}}[k] \in \mathbb{C}^{M^R \times M^B}$  are given by the spherical wave model as

$$\{\tilde{\mathbf{G}}[k]\}_{m_r,m_b} = e^{-j2\pi f_k \cdot \frac{D(m_r,m_b)}{c}},\tag{4}$$

where  $D^{(m_r,m_b)}$  denotes the distance between the  $m_r$ -th RIS-side subarray and the  $m_b$ -th BS-side subarray. Furthermore, the subarray response vectors  $\mathbf{a}_{\mathrm{R}}(\theta_{\mathrm{R},\mathrm{A}},\phi_{\mathrm{R},\mathrm{A}},f_k) \in \mathbb{C}^{N_{\mathrm{sub}}^{\mathrm{R}} \times 1}$  and  $\mathbf{a}_{\mathrm{B}}(\theta_{\mathrm{B}},\phi_{\mathrm{B}},f_k)$  $\in \mathbb{C}^{N_{\mathrm{sub}}^{\mathrm{B}} \times 1}$  are defined in Equation (2). The constant coefficient  $G_{\mathrm{T}}$  represents the antenna gain at the BS, which is the gain of the single antenna element, and is different from the array gain generated by beamforming [32]. The only unknown parameter in Equation (3) is the channel coefficient  $\alpha[k]$ , which can be obtained by placing a power detector at the RIS-side. Therefore, it is reasonable to assume that the BS-RIS channel is quasi-static and known.

## 247 2.2.2 RIS-UE Channel Model

As illustrated in Figure 1(b), we consider a multi-path THz channel model for indoor environments by using deterministic ray-tracing techniques [33]. For a RIS-UE link in a room, the multi-path propagation model consists of one LoS path and  $L_p$  NLoS paths reflected or scattered by the surrounding walls in the room, where the paths reflected by the ceiling and the floor can be neglected [34]. The three-dimensional (3D) distances for the LoS path and the NLoS paths between RIS and UE are denoted as  $d_0$  and  $d_l$ , for  $1 \leq l \leq L_p$ , respectively. The total attenuation of EM wave propagation in the THz band is composed of molecular absorption and free space path loss. The free space path loss is  $\beta_{\text{spr}}(f_k, d_l) = \frac{c}{4\pi f_k d_l}$ , and the molecular absorption is given by  $\beta_{\text{abs}}(f_k, d_l) = e^{-\frac{1}{2}\kappa(f_k)d_l}$ , where  $\kappa(f_k)$  stands for the frequency-dependent absorption coefficient [35]. Hence, the spatial-frequency channel  $\mathbf{h}[k] \in \mathbb{C}^{1 \times N^{\text{R}}}$  for the RIS-UE link can be written as

$$\mathbf{h}[k] = \beta[k]\tilde{\mathbf{h}}_{\text{LoS}}[k] \otimes \mathbf{a}_{\text{R}}^{\text{H}}(\theta_{\text{R},\text{D}}^{\text{LoS}}, \phi_{\text{R},\text{D}}^{\text{LoS}}, f_k) + \frac{1}{\sqrt{L_p}} \sum_{l=1}^{L_p} \beta_l[k]\tilde{\mathbf{h}}^l[k] \otimes \mathbf{a}_{\text{R}}^{\text{H}}(\theta_{\text{R},\text{D}}^l, \phi_{\text{R},\text{D}}^l, f_k),$$
(5)

where  $\beta[k] = \beta_{\rm spr}(f_k, d_0)\beta_{\rm abs}(f_k, d_0)$  and  $\beta_l[k] = \beta_{\rm spr}(f_k, d_l)\beta_{\rm abs}(f_k, d_l)\beta_{\rm RC}$  are the channel attenuation coefficients of the LoS component and the *l*-th NLoS component, respectively,  $(\theta_{\rm R,D}^{\rm LoS}, \phi_{\rm R,D}^{\rm LoS})$ and  $(\theta_{\rm R,D}^l, \phi_{\rm R,D}^l)$  are the LoS AoD and the NLoS AoD of the *l*-th NLoS path, respectively. We model the approximate reflection coefficient  $\beta_{\rm RC}$  in dB unit as an independent Gaussian random variable, i.e.,  $10 \log \beta_{\rm RC}[dB] \sim \min \{\mathcal{N}(\mu_{\rm R}, \sigma_{\rm R}^2), 0\}$ . The entries of  $\tilde{\mathbf{h}}_{\rm LoS}[k] \in \mathbb{C}^{1 \times M^{\rm R}}$  are given as  $\{\tilde{\mathbf{h}}_{\rm LoS}[k]\}_{m_r} = e^{-j2\pi f_k \cdot \frac{d(m_r)}{c}}$ , where  $d^{(m_r)}$  denotes the 3D distance between the UE and the  $m_r$ -th RIS-side subarray.  $\tilde{\mathbf{h}}^l[k]$  can be written by using a similar notation and assumptions.

# <sup>255</sup> 3 Problem Formulation and Proposed Channel Estimation <sup>256</sup> Solution

## 257 3.1 Problem Formulation of Channel Estimation

In this subsection, we formulate the downlink channel estimation problem for the considered RIS-assisted THz massive MIMO communication system over hybrid-field channels. As shown in Figure 2, we consider the two-stage frame structure consisting of the pilot training and data transmission stages. At the pilot training stage, the BS transmits *M* pilot OFDM symbols (i.e., *M* time



Figure 2: The schematic diagram of the frame structure, RIS element selection pattern, and RIS-UE sub-sampling channel, where the selected parts are marked in yellow blocks and the number in the yellow block indicates the index of the selected element.

slots) dedicated to channel estimation. The received signal at the UE-side<sup>1</sup> in the *m*-th time slot on the *k*-th subcarrier can be written as

$$y_m[k] = \sqrt{P_{\rm T}} \mathbf{h}[k] \mathbf{\Phi}_m \mathbf{G}[k] \mathbf{F}_{\rm RF} \mathbf{F}_{\rm BB}[k] \mathbf{s}_m[k] + n_m[k], \tag{6}$$

where  $1 \leq k \leq K$ ,  $1 \leq m \leq M$ ,  $P_{\rm T}$  is the transmit power of the BS,  $\mathbf{s}_m[k] \in \mathbb{C}^{U \times 1}$  denotes the transmitted symbol vector with  $\mathbb{E}\{\mathbf{s}_m[k]\mathbf{s}_m^{\rm H}[k]\} = \mathbf{I}_U$ , and  $n_m[k] \sim \mathcal{CN}(0, \sigma_n^2)$  is the effective complex additive white Gaussian noise (AWGN) at the UE, while  $\mathbf{h}[k] \in \mathbb{C}^{1 \times N^{\rm R}}$  and  $\mathbf{G}[k] \in \mathbb{C}^{N^{\rm R} \times N^{\rm B}}$  are the downlink RIS-UE and BS-RIS channels on the k-th subcarrier, respectively. Denote the control vector  $\mathbf{v}_{m_r,m} \in \mathbb{C}^{1 \times N_{\rm sub}^{\rm R}}$  for the  $m_r$ -th subarray elements of the RIS in the m-th time slot as

$$\mathbf{v}_{m_r,m} = \mathbf{o}_{m_r,m} \odot \tilde{\mathbf{v}}_{m_r,m} = \left[\cdots, \eta_{n_{\mathrm{sub}}^r,m_r,m},\cdots\right] \odot \left[\cdots, e^{\mathbf{j}\phi_{n_{\mathrm{sub}}^r,m_r,m}},\cdots\right],\tag{7}$$

where  $\mathbf{o}_{m_r,m} \in \mathbb{C}^{1 \times N_{sub}^{R}}$  represents the amplitude control vector,  $\tilde{\mathbf{v}}_{m_r,m} \in \mathbb{C}^{1 \times N_{sub}^{R}}$  represents the phase control vector, and  $1 \le n_{sub}^r \le N_{sub}^{R}$ , while  $\eta_{n_{sub}^r,m_r,m} \in [0, 1]$  and  $\phi_{n_{sub}^r,m_r,m} \in [0, 2\pi]$  are the amplitude control coefficient and phase control coefficient, respectively.  $\eta_{n_{sub}^r,m_r,m}^r \in [0, 2\pi]$  are the turn on/off the refraction functions of each RIS element. The entire RIS elements can be expressed as  $\mathbf{v}_m = \mathbf{o}_m \odot \tilde{\mathbf{v}}_m = [\mathbf{v}_{1,m}, \cdots, \mathbf{v}_{m_r,m}, \cdots, \mathbf{v}_{M^R,m}]^T \in \mathbb{C}^{N^R \times 1}$ , where  $\mathbf{o}_m = [\mathbf{o}_{1,m}, \cdots, \mathbf{o}_{M^R,m}]^T \in$  $\mathbb{C}^{N^R \times 1}$  and  $\tilde{\mathbf{v}}_m = [\tilde{\mathbf{v}}_{1,m}, \cdots, \tilde{\mathbf{v}}_{M^R,m}]^T \in \mathbb{C}^{N^R \times 1}$ . Then the RIS's refraction phase matrix is defined as  $\mathbf{\Phi}_m = \operatorname{diag}(\mathbf{v}_m) = \mathbf{O}_m \odot \tilde{\mathbf{V}}_m \in \mathbb{C}^{N^R \times N^R}$ , where  $\mathbf{O}_m = \operatorname{diag}(\mathbf{o}_m) \in \mathbb{C}^{N^R \times N^R}$  is the RIS selection matrix and  $\tilde{\mathbf{V}}_m = \operatorname{diag}(\tilde{\mathbf{v}}_m) \in \mathbb{C}^{N^R \times N^R}$  is the RIS phase matrix.

F<sub>RF</sub>  $\in \mathbb{C}^{N^{B} \times M^{B}}$  and  $\mathbf{F}_{BB}[k] \in \mathbb{C}^{M^{B} \times U}$  are respectively analog and digital beamforming matrices that are used at the BS to perform hybrid beamforming. Due to the sub-connected hybrid MIMO architecture, the analog beamformer implemented by the phase shifters can be expressed as

$$\mathbf{F}_{\rm RF} = \text{blkdiag}(\mathbf{f}_1, \cdots, \mathbf{f}_{m_b}, \cdots, \mathbf{f}_{M^{\rm B}}),\tag{8}$$

where  $\mathbf{f}_{m_b} = \left[\mathbf{f}_{m_b,1}, \cdots, \mathbf{f}_{m_b,n_{\text{sub}}^b}, \cdots, \mathbf{f}_{m_b,N_{\text{sub}}^B}\right]^{\text{T}} \in \mathbb{C}^{N_{\text{sub}}^{\text{B}} \times 1}$  with  $\left|\mathbf{f}_{m_b,n_{\text{sub}}^b}\right|^2 = 1/N_{\text{sub}}^{\text{B}}$ . Since the BS-RIS channel with the LoS path only is quasi-static and known, each analog beamforming vector can be designed as

$$\mathbf{f}_{m_b} = \mathbf{a}_{\mathrm{B}}(\theta_{\mathrm{B}}, \phi_{\mathrm{B}}, f_k), \ 1 \le m_b \le M^{\mathrm{B}},\tag{9}$$

where k can be set to K/2 for alleviating the beam squint problem caused by the large bandwidth [36]. The digital beamformer  $\mathbf{F}_{BB}[k]$  is designed according to the zero forcing (ZF) precoding in order to eliminate the multi-stream interference between the BS and the RIS subarrays, i.e.,

$$\mathbf{F}_{\mathrm{BB}}[k] = \zeta \tilde{\mathbf{G}}_{\mathrm{eq}}^{\dagger}[k] = \zeta \tilde{\mathbf{G}}_{\mathrm{eq}}^{\mathrm{H}}[k] \left( \tilde{\mathbf{G}}_{\mathrm{eq}}[k] \tilde{\mathbf{G}}_{\mathrm{eq}}^{\mathrm{H}}[k] \right)^{-1}, \tag{10}$$

where  $\tilde{\mathbf{G}}_{eq}[k] = [\alpha[k]G_{T}\tilde{\mathbf{G}}[k] \otimes \mathbf{a}_{B}^{H}(\theta_{B}, \phi_{B}, f_{k})]\mathbf{F}_{RF} \in \mathbb{C}^{M^{R} \times M^{B}}$  is the equivalent BS-RIS channel obtained from the perspective of the first element of different subarrays at the RIS, and  $\zeta = \sqrt{M^{B}/\text{Tr}\{\tilde{\mathbf{G}}_{eq}^{\dagger}[k](\tilde{\mathbf{G}}_{eq}^{\dagger}[k])^{H}\}}$  is a constant to meet the total transmit power constraint after beam-

<sup>&</sup>lt;sup>1</sup>Note that since each UE can perform channel estimation independently, UE subscripts are omitted.

forming. In this way, the multi-stream interference between the BS and the RIS subarrays can be eliminated, i.e.,  $\mathbf{G}_{eq}[k] = \mathbf{G}[k]\mathbf{F}_{RF}\mathbf{F}_{BB}[k] \in \mathbb{C}^{N^R \times U}$ ,  $\forall k$ , is a block diagonal constant matrix. Therefore, the equivalent pilot signal  $\mathbf{p}_m \in \mathbb{C}^{N^R \times 1}$  can be expressed as

$$\mathbf{p}_{m} = \underbrace{\left[\mathbf{O}_{m} \odot \tilde{\mathbf{V}}_{m}\right]}_{\boldsymbol{\Phi}_{m}} \underbrace{\left[\mathbf{G}[k]\mathbf{F}_{\mathrm{RF}}\mathbf{F}_{\mathrm{BB}}[k]\right]}_{\mathbf{G}_{\mathrm{eq}}[k]} \mathbf{s}_{m}[k], \tag{11}$$

where  $\mathbf{p}_m$  is identical for different subcarriers since we set the transmit symbol  $\mathbf{s}_m[k]$  to be  $\mathbf{1}_U, \forall m, k$ , 287 and the ZF digital beamformer Equation (10) for  $\mathbf{G}[k]$ . Under the assumption that the BS and RIS 288 meet the parallel symmetric array arrangements,  $\mathbf{G}_{eq}[k]$  is defined by  $\sqrt{N^{B}}\alpha[k]G_{T}$  blkdiag  $(\mathbf{1}_{N^{B}}^{1}, \cdots, \mathbf{1}_{N^{B}})$ 289  $\mathbf{1}_{N^{\mathrm{B}}_{\lambda}}^{u}, \cdots, \mathbf{1}_{N^{\mathrm{B}}_{\lambda}}^{U}$ ). Thus, the effective pilot signals can be further expressed as the RIS element vector 290 given by  $\mathbf{p}_m = \sqrt{N^{\mathrm{B}}} \alpha[k] G_{\mathrm{T}} \mathbf{v}_m \approx \sqrt{N^{\mathrm{B}}} \alpha G_{\mathrm{T}} \mathbf{v}_m = A_{\mathrm{T}} \mathbf{v}_m$ , where the approximation  $\alpha[k] \approx \alpha, \forall k$ , is 291 further applied and  $A_{\rm T} = \sqrt{N^{\rm B}} \alpha G_{\rm T}$  represents the total attenuation from the BS to the RIS. 292 By collecting  $y_m[k]$  for  $1 \leq m \leq M$  together, the aggregate received signal vector  $\mathbf{y}[k] =$ 293  $[y_1[k], \cdots, y_M[k]] \in \mathbb{C}^{1 \times M}$  can be expressed as 294

$$\mathbf{y}[k] = \sqrt{P_{\mathrm{T}}} \mathbf{h}[k] \mathbf{P} + \mathbf{n}[k], \tag{12}$$

where  $\mathbf{P} = [\mathbf{p}_1, \cdots, \mathbf{p}_M] = A_{\mathrm{T}} \mathbf{V} = A_{\mathrm{T}} [\mathbf{v}_1, \cdots, \mathbf{v}_M] \in \mathbb{C}^{N^{\mathrm{R}} \times M}$ , and  $\mathbf{n}[k] = [n_1[k], \cdots, n_M[k]] \in \mathbb{C}^{1 \times M}$ . Thus, the received signal matrix  $\mathbf{Y} = [\mathbf{y}^{\mathrm{T}}[1], \cdots, \mathbf{y}^{\mathrm{T}}[K]]^{\mathrm{T}} \in \mathbb{C}^{K \times M}$  can be expressed as

$$\mathbf{Y} = \sqrt{P_{\mathrm{T}}} \mathbf{H} \mathbf{P} + \mathbf{N},\tag{13}$$

where  $\mathbf{H} = \begin{bmatrix} \mathbf{h}^{\mathrm{T}}[1], \cdots, \mathbf{h}^{\mathrm{T}}[K] \end{bmatrix}^{\mathrm{T}} \in \mathbb{C}^{K \times N^{\mathrm{R}}}$  denotes the downlink spatial-frequency domain RIS-UE channel matrix, and  $\mathbf{N} = \begin{bmatrix} \mathbf{n}^{\mathrm{T}}[1], \cdots, \mathbf{n}^{\mathrm{T}}[K] \end{bmatrix}^{\mathrm{T}} \in \mathbb{C}^{K \times M}$ .

## <sup>299</sup> 3.2 Deep Learning Based Spatial-Frequency Domain Channel Extrapo <sup>300</sup> lation

We choose to activate only  $N_s^{\rm R} = \frac{N^{\rm R}}{\rho}$  RIS elements and estimate the sub-channels associated 301 with the activated RIS elements, where  $\rho > 1$  is the element compression ratio. Furthermore, as 302 shown in Figure 2, only  $K_s = \frac{K}{\bar{a}}$  uniformly selected subcarriers are used for pilot training, where 303  $\bar{\rho}$  is the frequency compression ratio, and the remaining subcarriers can be used for transmitting 304 control signals. We also give an example of the RIS element pattern selected uniformly and the 305 corresponding RIS-UE side sub-sampling spatial-frequency channel in Figure 2, where the yellow 306 blocks indicate the selected elements and the selected subcarriers. Thus, the practical received pilot 307 signal  $\mathbf{Y}_s \in \mathbb{C}^{K_s \times M}$  is defined as 308

$$\mathbf{Y}_s = \sqrt{P_{\mathrm{T}}} \mathbf{H}_s \mathbf{P}_s + \mathbf{N}_s, \tag{14}$$

where  $\mathbf{H}_s \in \mathbb{C}^{K_s \times N_s^{\mathrm{R}}}$  is the sub-sampling of the spatial-frequency channel,  $\mathbf{P}_s \in \mathbb{C}^{N_s^{\mathrm{R}} \times M}$  is the corresponding equivalent pilot signal, and  $\mathbf{N}_s$  is the noise. Our goal is to recover the complete channel  $\hat{\mathbf{H}} \in \mathbb{C}^{K \times N^{\mathrm{R}}}$  from the limited received pilot signals  $\mathbf{Y}_s$ , i.e., extrapolating the rest unknown



Figure 3: The overall structure of the proposed DL-based spatial-frequency domain channel extrapolation scheme.

channels from the acquired partial channels. Based on the universal approximation capability of DL, we can use DL to characterize the mapping among the channels at different space/frequency locations. Thus, we propose a DL-based spatial-frequency domain channel extrapolation network, which consists of the element selection strategy (ESS), pilot design, CSI feedback, sub-channel estimation, and spatial-frequency domain channel extrapolation modules, as illustrated in Figure 3. The overall flow of the proposed scheme is expressed as

$$\hat{\mathbf{H}} = f_{\rm SFDE}(f_{\rm SCE}(f_{\rm CsiFd}(\sqrt{P_{\rm T}}\mathbf{H}_s\mathbf{P}_s + \mathbf{N}_s))) = f_{\rm SFDE}(f_{\rm SCE}(f_{\rm CsiFd}(\sqrt{P_{\rm T}}f_{\rm ESS}(\mathbf{H})\mathbf{P}_s + \mathbf{N}_s))), (15)$$

where the mapping  $f_{\rm ESS}(\cdot)$  represents the element selection strategy for deciding the sub-sampling channel, and the equivalent pilot signal  $\mathbf{P}_{\rm s}$  can be learned as the trainable parameters, while  $f_{\rm CsiFd}(\cdot)$ ,  $f_{\rm SCE}(\cdot)$  and  $f_{\rm SFDE}(\cdot)$  represent the CSI feedback network, the sub-channel estimation network and the spatial-frequency domain extrapolation network, respectively. We now detail each component.

## 313 3.2.1 Element Selection Strategy

With only  $N_s^{\rm R}$  activated RIS elements, from (7), the RIS element selection vector  $\mathbf{o} = \mathbf{o}_m = [\mathbf{o}_{1,m}, \cdots, \mathbf{o}_{m_r,m}, \cdots, \mathbf{o}_{M^{\rm R},m}]^{\rm T} \in \{0,1\}^{N^{\rm R} \times 1}$  is an  $N_s^{\rm R}$ -hot vector with  $N_s^{\rm R}$  elements being '1' and the other elements being '0', where the subscript 'm' can be dropped since  $\mathbf{o}$  is fixed at the pilot training stage. Also since only  $K_s$  subcarriers are uniformly selected for pilot training, the frequency selection vector  $\boldsymbol{\kappa} \in \{0,1\}^{K \times 1}$  is defined by  $\{\boldsymbol{\kappa}\}_{\bar{\rho}k+1} = 1, 0 \leq k \leq K_s - 1$ , and the other elements being '0'. Thus, the selection operation of the sub-sampling function  $f_{\rm ESS}(\cdot)$  can be expressed as

$$\mathbf{H}_s = f_{\text{ESS}}(\mathbf{H}) = \mathbf{S} \odot \mathbf{H},\tag{16}$$

where  $\mathbf{S} = \boldsymbol{\kappa} \otimes \mathbf{o}^{\mathrm{T}} \in \{0, 1\}^{K \times N^{\mathrm{R}}}$  is the spatial-frequency selection matrix, and the zero rows/columns in  $\mathbf{S} \odot \mathbf{H}$  are deleted directly to yield  $\mathbf{H}_{s}$ . Note that the selection of the activated elements, i.e., the RIS element selection vector, can affect the extrapolation performance. Thus, we consider the <sup>323</sup> following three element selection strategies.

1) Uniform Selection Strategy: Since each subarray in the RIS is a UPA, its element compression ratio is expressed as  $\rho = \rho_y \times \rho_z$ . For fairness,  $\rho_y$  and  $\rho_z$  should be as close as possible. When  $\rho_y \neq \rho_z$ , the z-axis compression should be considered first to maintain the y-axis angular resolution of indoor

 $_{327}$  UEs. Thus, the *y*-*z* compression ratio allocation can be solved from the following optimization

$$\begin{array}{ll} \min_{\{\rho_y,\rho_z\}} & |\rho_z - \rho_y|, \\ \text{s.t.} & \rho_y \times \rho_z = \rho, \\ & 1 \le \rho_y \le \rho_z. \end{array} \tag{17}$$

Some allocation examples as  $\rho(2, 4, 8, 16) = \rho_y(1, 2, 2, 4) \times \rho_z(2, 2, 4, 4)$ . Given  $\rho_y$  and  $\rho_z$ , the active element index vector  $\boldsymbol{\xi}_{m_r} \in \mathbb{C}^{1 \times N_{\text{sub}}^{\text{R}}/\rho}$  of the  $m_r$ -th subarray can be expressed as

$$\{\boldsymbol{\xi}_{m_r}\}_{n_i^y N_z^{\mathrm{R}}/\rho_z + n_i^z + 1} = N_{\mathrm{sub}}^{\mathrm{R}}(m_r - 1) + N_z^{\mathrm{R}}\rho_y n_i^y + \rho_z n_i^z + 1,$$
(18)

where  $1 \leq m_r \leq M^{\mathrm{R}}$ ,  $0 \leq n_i^y \leq \frac{N_y^{\mathrm{R}}}{\rho_y} - 1$ , and  $0 \leq n_i^z \leq \frac{N_z^{\mathrm{R}}}{\rho_z} - 1$ . The entire active element index vector or set of the RIS is defined as  $\boldsymbol{\xi} = [\boldsymbol{\xi}_1, \cdots, \boldsymbol{\xi}_{m_r}, \cdots, \boldsymbol{\xi}_{M^{\mathrm{R}}}]^{\mathrm{T}} \in \mathbb{C}^{N_s^{\mathrm{R}} \times 1}$ . Thus, we set the entries of the RIS element selection vector  $\mathbf{o}$  corresponding to the index set  $\boldsymbol{\xi}$  to '1', i.e.,  $\{\mathbf{o}\}_{\boldsymbol{\xi}} = 1$  for  $\boldsymbol{\xi} \in \boldsymbol{\xi}$ , and the other elements of  $\mathbf{o}$  to '0'.

<sup>334</sup> 2) Random Selection Strategy: It randomly selects  $N_s^{\rm R}$  elements from the RIS as the random <sup>335</sup> pattern. We can use the function 'random.choice(·)' of the NumPy library in Python to generate <sup>336</sup> the active element index vector  $\boldsymbol{\xi}$ . Given that the elements of the pattern are randomly selected <sup>337</sup> from the whole RIS, if the element compression ratio  $\rho$  is not large, then the aperture of a random <sup>338</sup> pattern is usually comparable to that of the RIS.

3) Learning-Based Selection Strategy: In addition to the above two fixed selection strategies, the learning-based element selection strategy has also been widely studied. In [14], a differentiable selection network is proposed to learn the element selection vector **o**. The input of this network is a random initialization vector. By utilizing several fully-connected layers and the softmax function, a probability vector  $\mathbf{g} = \{g_1, g_2, \dots, g_{N^R}\}$  is generated, where  $g_i$  represents the probability of the *i*-th element being selected. Thus, the active element index vector  $\boldsymbol{\xi}$  can be defined as

$$\boldsymbol{\xi} = \arg \operatorname{top}_{N_{\circ}^{\mathrm{R}}} \{ \mathbf{g} \}, \tag{19}$$

where arg top<sub> $N_s^{\rm R}$ </sub>{·} is a function that finds the element index set of the first  $N_s^{\rm R}$  largest selection probabilities. The details of the antenna selection network can be found in [14].

#### 347 3.2.2 Pilot Design

As aforementioned in Equation (11), under the assumption that the BS and RIS meet the parallel symmetric array arrangements, the equivalent downlink pilots can be defined as  $\mathbf{P}_s = A_{\mathrm{T}} \mathbf{V}_s$ , where  $\mathbf{V}_s \in \mathbb{C}^{N_s^{\mathrm{R}} \times M}$  is the phase matrix of selected RIS elements at the pilot training stage. Thus, the pilot matrix  $\mathbf{P}_s$  can be obtained by adjusting the RIS phase at different time slots, which is 352 given by

$$\mathbf{P}_{s} = A_{\mathrm{T}} \exp^{(\mathbf{j}\mathbf{\Xi})} = A_{\mathrm{T}} \left( \cos(\mathbf{\Xi}) + \mathbf{j} \sin(\mathbf{\Xi}) \right), \tag{20}$$

where  $\Xi \in \mathbb{R}^{N_s^{\mathrm{R}} \times M}$  is the phase control matrix of selected RIS elements. As it is well known that complex-valued outputs are not well supported by most deep learning frameworks (e.g., Tensorflow, Pytorch), it would be difficult to directly train the complex-valued pilot matrix  $\mathbf{P}_s$ . Hence, we adopt the real-valued RIS phase control matrix  $\Xi$  whose entries take values in  $[0, 2\pi)$  as the real-valued trainable parameters of the pilot design network (PDN) and the pilot matrix  $\mathbf{P}_s$  can be obtained from Equation (20). The structure of the PDN is shown in Figure 3, where the trainable parameters of the PDN, i.e.,  $\Xi$ , are learned at the DL training stage.

## 360 3.2.3 CSI Feedback

At the uplink CSI feedback stage, the UE extracts the CSI from the received pilot signals and 361 feeds it back to the BS. However, the large number of array elements results in excessive feedback 362 overhead. Recently, DL-based solutions, such as CsiNet [37], have achieved good performance for 363 CSI feedback. Furthermore, an emerging DL architecture, known as the transformer [38], is utilized 364 as the novel CSI feedback network to further reduce the feedback overhead and obtain more efficient 365 compression performance than the CsiNet framework [39]. Therefore, in this paper, we utilize the 366 transformer as the backbone of the CSI feedback network  $f_{\rm CsiFd}(\cdot)$ . The original transformer is 367 composed of an encoder and a decoder. However, since we are dealing with the CSI without time-368 sequential information, there is no causality constraint. Thus, we only exploit the encoder module 369 of the transformer which obtains output in parallel. Since neural networks are more effective for 370 real-valued operations than complex-valued operations and the transformer can only extract the 371



Figure 4: The schematic diagram of the transformer encoder.

correlation between sequences, we reshape the received pilot signal into a real-valued two-dimensional (2D) sequence  $\bar{\mathbf{Y}}_s \in \mathbb{R}^{K_s \times 2M}$ , which can be expressed as

$$\bar{\mathbf{Y}}_s = [\Re\{\mathbf{Y}_s\}, \Im\{\mathbf{Y}_s\}], \tag{21}$$

 $_{374}$  where the number of subcarriers  $K_s$  represents the input sequence length of the transformer.

The schematic diagram of the transformer encoder is shown in Figure 4. Through the fullyconnected linear embedding layer, the input sequence  $\bar{\mathbf{Y}}_s$  is converted to  $\mathbf{X}_s \in \mathbb{R}^{K_s \times d_T}$ , which merges the relative position information of the sub-carriers using the positional embedding layer. Then, multiple encoder layers are utilized to extract features from the input sequences. Each encoder layer has the same structure which is composed of a multi-head self-attention sub-layer followed by a position-wise multi-layer perceptron (MLP) sub-layer. Layernorm is applied before every block and the residual connection is applied after every block. Among them, the multi-head attention mechanism is the key component for the performance improvement of the transformer. As shown in Figure 4, the input sequence  $\mathbf{X}_s$  is first projected onto three different sequential vectors: the queries, keys, and values with different learned linear projections, respectively, namely,  $\{\mathbf{Q}_i, \mathbf{K}_i, \mathbf{V}_i\} \in \mathbb{R}^{K_s \times d_m}, 1 \le i \le h$ , where h is the number of heads and  $d_m = d_T/h$ . Then, each value head<sub>i</sub>  $\in \mathbb{R}^{K_s \times d_m}, 1 \le i \le h$ , is outputted by performing the scaled dot-product attention in parallel, where a softmax function is applied to obtain the weight on the value, which is given by

head<sub>i</sub> = softmax 
$$\left(\frac{\mathbf{Q}_i \mathbf{K}_i^{\mathrm{T}}}{\sqrt{d_m}}\right) \mathbf{V}_i, 1 \le i \le h.$$
 (22)

These output values are concatenated and projected back to a  $d_{\rm T}$ -dimensional representation using the linear projection matrix  $\mathbf{W}^O \in \mathbb{R}^{K_s \times d_{\rm T}}$  as

$$MultiHead(\mathbf{X}_{s}) = Concat(head_{i}, \cdots, head_{h})\mathbf{W}^{O}.$$
(23)

After the transformer encoder, a fully-connected linear layer and sigmoid function are used to generate a real-valued compressed codeword. The codeword is then converted to B bits as the feedback information through a quantization layer, which is constructed by a B-bit uniform scalar quantizer. The above feedback process generates the binary vector  $\mathbf{q} \in \{0, 1\}^B$  as

$$\mathbf{q} = f_{\text{CsiFd}} \left( \bar{\mathbf{Y}}_s; \mathcal{W}_F \right), \tag{24}$$

 $_{379}$  where  $\mathcal{W}_F$  denotes the trained parameter set of the CSI feedback network.

#### 380 3.2.4 Sub-Channel Estimation

When the BS receives the feedback bits, the sub-channel estimation network is used to reconstruct the sub-sampling of the complete spatial-frequency channel. Similar to Subsection 3.2.3, we also consider the transformer encoder as the backbone of this part. As shown in Figure 3, the received CSI feedback bit vector is first inputted into a dequantization layer, which conducts the inverse operation of the quantizer and outputs a real-valued vector. Then, a fully-connected layer
is utilized to obtain an initial coarse channel estimation. Finally, the transformer encoder extracts
the spatial-frequency correlation of the channel and further improves the channel estimation performance. The output of the sub-channel estimation network is expressed as

$$\bar{\mathbf{H}}_{s} = \left[ \Re\{\hat{\mathbf{H}}_{s}\}, \Im\{\hat{\mathbf{H}}_{s}\} \right] = f_{\text{SCE}}\left(\mathbf{q}; \mathcal{W}_{S}\right), \qquad (25)$$

where  $\hat{\mathbf{H}}_{s} \in \mathbb{C}^{K_{s} \times N_{s}^{\mathrm{R}}}$  is the estimated sub-sampling channel,  $\bar{\mathbf{H}}_{s} \in \mathbb{R}^{K_{s} \times N_{s}^{\mathrm{R}} \times 2}$  is a real-valued 3D matrix, and  $\mathcal{W}_{s}$  is the trained parameter set of the sub-channel estimation network.

### <sup>391</sup> 3.2.5 Spatial-Frequency Domain Channel Extrapolation

First, the initial input  $\tilde{\mathbf{H}} \in \mathbb{R}^{K \times N^{R} \times 2}$  to the channel extrapolation network is constructed from the estimated sub-sampling channel  $\bar{\mathbf{H}}_{s} \in \mathbb{R}^{K_{s} \times N_{s}^{R} \times 2}$  with the known RIS spatial-frequency selection pattern **S**. Specifically, we copy the entries of  $\bar{\mathbf{H}}_{s}$  to the corresponding positions in  $\tilde{\mathbf{H}}$  and fill the other elements of  $\tilde{\mathbf{H}}$  with zeros according to the known RIS spatial-frequency selection pattern **S**. This initial operation can be expressed as

$$\tilde{\mathbf{H}} = f_{\mathrm{zfi}} \left( \bar{\mathbf{H}}_s; \mathbf{S} \right). \tag{26}$$

The non-zero rows/columns of  $\tilde{\mathbf{H}}$  are consistent with  $\bar{\mathbf{H}}_s$  and their locations are the same as the locations of elements '1' in **S**. The neighborhood information in the receptive field is then extracted using a convolutional layer for initial interpolation. To guarantee that the output dimensions from the convolution layer remain unchanged, we employ zero padding in the convolution layer.

Subsequently, we consider a competitive yet conceptually and technically simple architecture, called MLP-Mixer [40], as the backbone of the channel extrapolation network. The architecture of this MLP-Mixer is based entirely on MLPs, which can extract and reconstruct 2D features by repeatedly applying them to either spatial locations or feature channels. Specifically, the input  $\tilde{\mathbf{H}} \in \mathbb{R}^{K \times N^{R} \times 2}$  is rearranged as a sequence of flattened 2D patches  $\mathbf{X}_{p} \in \mathbb{R}^{N_{p} \times (2L^{2})}$ , where  $(K, N^{R})$ is the resolution of the original input, (L, L) is the resolution of each patch, and  $N_{p} = KN^{R}/L^{2}$  is the resulting number of patches. Then, all the patches are linearly projected with the same projection matrix. This results in a 2D real-valued matrix  $\tilde{\mathbf{X}} \in \mathbb{R}^{N_{p} \times d_{M}}$ . Next the input matrix  $\tilde{\mathbf{X}}$  is fed into



Figure 5: The structure of the mixer layer.

several mixer layers to extrapolate the complete channel. As illustrated in Figure 5, each mixer layer consists of two MLP blocks. The first one acts on the columns of  $\tilde{\mathbf{X}}$ , maps  $\mathbb{R}^{N_p} \mapsto \mathbb{R}^{2N_p} \mapsto \mathbb{R}^{N_p}$ , and is shared across all the columns. The second acts on the rows of  $\tilde{\mathbf{X}}$ , i.e., on the transposed input matrix  $\tilde{\mathbf{X}}^{\mathrm{T}}$ , maps  $\mathbb{R}^{d_{\mathrm{M}}} \mapsto \mathbb{R}^{2d_{\mathrm{M}}} \mapsto \mathbb{R}^{d_{\mathrm{M}}}$ , and is shared across all the rows. Each MLP block contains two fully-connected layers and a nonlinear activation function. Denote the input matrix to the *t*-th mixer layer as  $\tilde{\mathbf{X}}_t$ . The mapping of the *t*-th mixer layer can be expressed as

$$\mathbf{U} = \dot{\mathbf{X}}_{t} + \mathbf{W}_{t,2} f_{\sigma}(\mathbf{W}_{t,1} \text{LayerNorm}(\dot{\mathbf{X}}_{t})),$$
  
$$\tilde{\mathbf{X}}_{t+1} = \mathbf{U} + \left(\mathbf{W}_{t,4} f_{\sigma}(\mathbf{W}_{t,3} \text{LayerNorm}(\mathbf{U})^{\mathrm{T}})\right)^{\mathrm{T}},$$
(27)

where  $\mathbf{W}_{t,i}$ ,  $1 \le i \le 4$ , are the parameter matrices of the fully-connected layers in the *t*-th mixer layer for  $1 \le t \le L_{\mathrm{M}}$ , and  $L_{\mathrm{M}}$  is the number of mixer layers, while  $f_{\sigma}$  is the activation function.

Finally, the output of the last mixer layer is linearly projected back to the original dimension  $\mathbb{R}^{N_p \times d_M} \mapsto \mathbb{R}^{N_p \times (2L^2)}$ , and the 2D patches are rearranged back to  $\mathbb{R}^{N_p \times (2L^2)} \mapsto \mathbb{R}^{K \times N^R \times 2}$  for obtaining the final extrapolation result  $\bar{\mathbf{H}} \in \mathbb{R}^{K \times N^R \times 2}$ , which is a real-valued 3D matrix. Thus, the extrapolation process can be expressed as

$$\hat{\mathbf{H}} = \bar{\mathbf{H}}_{[:,:,1]} + j\bar{\mathbf{H}}_{[:,:,2]} = f_{\text{SFDE}}\left(\bar{\mathbf{H}}_s; \mathcal{W}_E\right),\tag{28}$$

where  $\hat{\mathbf{H}} \in \mathbb{C}^{K \times N^{\mathrm{R}}}$  is the estimated complete complex-valued channel, and  $\mathcal{W}_{E}$  is the trained parameter set of the spatial-frequency domain extrapolation network.

#### 409 3.2.6 Training Strategy

The data set for off-line training is denoted as  $\mathcal{H}$ , where  $|\mathcal{H}| = N_{\text{set}}$  is the number of off-line training samples. Furthermore, a sample in  $\mathcal{H}$  is an input-label pair written as  $(\mathbf{H}, \mathbf{H})$ , where  $\mathbf{H}$  is the extrapolation target and is also the input of the entire SFDCEtra network. The input will go through the RIS array element and subcarrier sub-sampling strategy, since we need to extrapolate the original complete channel by receiving only the pilot signal of the sub-sampling channel.

With the uniform or random ESS  $f_{\text{ESS}}(\cdot)$ , at the off-line training stage, we consider E2E training to jointly optimize the pilot design network, CSI feedback network, sub-channel estimation network, and channel extrapolation network, by minimizing the normalized mean square error (NMSE) between the output  $\hat{\mathbf{H}}$  and the target  $\mathbf{H}$ . Thus, the loss function is written as

$$\mathcal{L}_{c} = \frac{1}{B_{e}} \sum_{i=1}^{B_{e}} \frac{\|\mathbf{H} - \hat{\mathbf{H}}\|_{F}^{2}}{\|\mathbf{H}\|_{F}^{2}},$$
(29)

419 where  $B_e$  is the batch size for off-line training.

When the learning-based ESS is adopted, the parameters for the ESS and the above networks are optimized jointly. The loss function for the joint optimization problem is given by

$$\mathcal{L} = \gamma \mathcal{L}_c + (1 - \gamma) \mathcal{L}_{\text{ESS}},\tag{30}$$

where  $0 < \gamma \leq 1$  is the weight to balance the penalties of channel extrapolation and ESS with  $\gamma = 1$  denoting that the non-learning based  $f_{\text{ESS}}(\cdot)$  is selected, and  $\mathcal{L}_{\text{ESS}}$  is the loss function of the learning-based ESS. The details of  $\mathcal{L}_{\text{ESS}}$  are available in [14].

## 425 4 Proposed Beamforming Solution

## 426 4.1 Problem Formulation of RIS-aided Multi-User Beamforming

At the data transmission stage, the BS can simultaneously supports U UEs with the aid of RIS, since the LoS MIMO architecture between the BS and RIS can support multi-stream transmission via intra-path multiplexing. Similar to Equation (6), the received signal at the *u*-th UE on the *k*-th subcarrier can be expressed as

$$y[u,k] = \sqrt{P_{\mathrm{T}}} \mathbf{h}[u,k] \mathbf{\Phi} \mathbf{G}[k] \mathbf{F}_{\mathrm{RF}} \mathbf{f}_{\mathrm{BB}}[u,k] s[u,k]$$
  
+ 
$$\sum_{i=1,i\neq u}^{U} \sqrt{P_{\mathrm{T}}} \mathbf{h}[i,k] \mathbf{\Phi} \mathbf{G}[k] \mathbf{F}_{\mathrm{RF}} \mathbf{f}_{\mathrm{BB}}[i,k] s[i,k] + n[u,k],$$
(31)

where  $\mathbf{h}[u, k] \in \mathbb{C}^{1 \times N^{\mathrm{R}}}$ ,  $1 \leq u \leq U, 1 \leq k \leq K$ , represents the downlink channel vector between the RIS and the *u*-th UE on the *k*-th subcarrier,  $\mathbf{f}_{\mathrm{BB}}[u, k] \in \mathbb{C}^{M^{\mathrm{B}} \times 1}$  denotes the digital baseband beamforming vector associated with the *u*-th UE on the *k*-th subcarrier. Thus, the signal-to-interference plus-noise-ratio (SINR) of the *u*-th UE on the *k*-th subcarrier can be expressed as

$$\operatorname{SINR}[u,k] = \frac{P_{\mathrm{T}}|\mathbf{h}[u,k]\boldsymbol{\Phi}\mathbf{G}[k]\mathbf{F}_{\mathrm{RF}}\mathbf{f}_{\mathrm{BB}}[u,k]|^{2}}{P_{\mathrm{T}}\sum_{i=1,i\neq u}^{U}|\mathbf{h}[i,k]\boldsymbol{\Phi}\mathbf{G}[k]\mathbf{F}_{\mathrm{RF}}\mathbf{f}_{\mathrm{BB}}[i,k]|^{2} + \sigma_{n}^{2}}.$$
(32)

 $_{431}$  Therefore, the sum rate R in the downlink multi-user transmission can be expressed as

$$R = \frac{1}{K} \sum_{u=1}^{U} \sum_{k=1}^{K} \log_2 \left( 1 + \text{SINR}[u, k] \right).$$
(33)

Based on the estimated RIS-UE channel at the pilot training stage, the BS can design the hybrid beamformer { $\mathbf{F}_{\text{RF}}, \mathbf{F}_{\text{BB}}[k], \forall k$ } and the RIS refraction phase matrix  $\boldsymbol{\Phi}$  to maximize the sum rate R, where  $\mathbf{F}_{\text{BB}}[k] = [\mathbf{f}_{\text{BB}}[1,k], \cdots, \mathbf{f}_{\text{BB}}[U,k]]$ . This design process is formulated as the following optimization problem

$$\begin{aligned} \max_{\mathcal{F}(\cdot)} & R, \\ \text{s.t.} & \{\mathbf{F}_{\text{RF}}, \mathbf{F}_{\text{BB}}[k], \forall k, \mathbf{\Phi}\} = \mathcal{F}\left(\hat{\mathbf{H}}[u], \forall u\right), \\ & \mathbf{F}_{\text{RF}} \in (8), \\ & \|\mathbf{F}_{\text{RF}} \mathbf{F}_{\text{BB}}[k]\|_{F}^{2} = M^{\text{B}}, \forall k, \\ & \{\mathbf{\Phi}\}_{i,i} = \{\mathbf{v}\}_{i} = e^{j\phi_{i}}, \phi_{i} \in [0, \ 2\pi), \forall i, \end{aligned}$$
(34)

436 where  $\hat{\mathbf{H}}[u]$  is the estimated spatial-frequency channel between the RIS and the u-th UE, and

<sup>437</sup>  $\mathcal{F}(\cdot)$  represents a function that maps the estimated RIS-UE channels onto the hybrid beamformer <sup>438</sup> { $\mathbf{F}_{\mathrm{RF}}, \mathbf{F}_{\mathrm{BB}}[k], \forall k$ } and the RIS refraction phase matrix  $\boldsymbol{\Phi}$ .

## <sup>439</sup> 4.2 Deep Learning Based Hybrid Beamforming and RIS Phase Design

In order to solve the optimization Equation (34), some alternating iterative algorithms [17, 440 18, 19] have been proposed to obtain the analog beamformer, digital beamformer, and RIS phase. 441 respectively. Unfortunately, all the aforementioned approaches are based on the idealized case that 442 the CSI is known accurately. However, perfect CSI is usually unavailable, especially for indoor 443 channel cases where the channel characteristics are complex due to rich scatterers. Inspired by 444 the universal approximation capability of DL, it is possible to use DL to learn the complicated and 445 unknown mapping from the estimated channels to the hybrid beamformers and RIS refraction phase. 446 Thus, we propose a DL-based hybrid beamforming and RIS phase design scheme, which consists of 447 analog beamformer design, DL-based RIS refraction phase design, and knowledge-data dual-driven 448 digital beamformer design. The block diagram of the proposed scheme is shown in Figure 6. 449



Figure 6: The overall structure of the proposed DL-based hybrid beamforming and RIS refraction phase design scheme.

## 450 4.2.1 Analog Beamformer Design

Since the BS's active beamforming and RIS's passive beamforming are coupled, the optimization 451 problem is non-convex, and it is very challenging to find a global optimum. Hence, we separately 452 design the analog beamforming of the BS and the passive beamforming of the RIS. Specifically, the 453 BS analog beamforming and the RIS passive beamforming are designed to improve the received 454 SINR of UEs. However, due to the sub-connected structure of phase shifters in the LoS MIMO 455 architecture, the interference among beams from the BS subarrays to the RIS subarrays cannot 456 be eliminated. Fortunately, this part of interference can be removed by appropriately designing 457 the digital beamforming. Therefore, when designing the analog beamforming on the BS-side, it is 458 sufficient to assume that the transmit energy of the BS is focused on the RIS. 459

Since the BS-RIS channel with only LoS path is quasi-static and known, we can utilize the angle information of the BS-RIS link to design the analog beamformer. Without loss of generality, we assume that the *u*-th UE is assisted by the  $m_r$ -th subarray of the RIS. Thus, the active beamforming designed for the *u*-th UE should be aligned to the  $m_r$ -th subarray of the RIS. Therefore, the transmit beam of the  $m_b$ -th subarray  $\mathbf{f}_{m_b}$  is designed according to Equation (9).

## 465 4.2.2 DL-Based RIS Refraction Phase Design

The key challenge in the RIS-aided communication system is to optimize a common RIS phase shared by all the subcarriers. In the THz broadband case, there exists a non-negligible beam squint effect for different subcarriers [36]. Therefore, when designing the common RIS phase, it is necessary to consider this effect on all subcarriers, which makes the RIS phase design much more difficult than the narrowband case. To solve this challenging problem, we propose a transformer-based RIS phase design network (RPDN), as shown in Figure 6, to design the RIS refraction phase matrix.

We first convert all the estimated RIS-UE channels  $\hat{\mathbf{H}}[u] \in \mathbb{C}^{K \times N^{\mathrm{R}}}$  for  $1 \leq u \leq U$  into a real-valued 3D matrix  $\bar{\mathcal{H}} \in \mathbb{R}^{U \times K \times 2N^{\mathrm{R}}}$ , which is expressed as

$$\bar{\mathcal{H}} = \left[\bar{\mathbf{H}}[1], \cdots, \bar{\mathbf{H}}[u], \cdots, \bar{\mathbf{H}}[U]\right], \tag{35}$$

where  $\bar{\mathbf{H}}[u] = [\Re\{\hat{\mathbf{H}}[u]\}, \Im\{\hat{\mathbf{H}}[u]\}] \in \mathbb{R}^{K \times 2N^{R}}$  and  $\hat{\mathbf{H}}[u]$  is the estimated RIS-UE channel of the *u*-th UE obtained from the DL-based SFDCEtra network.  $\bar{\mathcal{H}}$  is inputted into the transformer encoder, which globally extracts the inter-subcarrier correlation. To consider the beam squint effect for different subcarriers, the 2D matrix  $\mathbf{X}_{r} \in \mathbb{R}^{U \times N^{R}/U}$  is obtained by the mean operation over the subcarrier dimension of the transformer encoder's output. Then  $\mathbf{X}_{r}$  is flattened as  $\mathbf{x}_{r} \in \mathbb{R}^{N^{R} \times 1}$ , and passes through the activation function to generate the RIS phase vector  $\mathbf{v} \in \mathbb{C}^{N^{R} \times 1}$  that satisfies the the constant modulus constraint. The corresponding activation function is defined as

$$\mathbf{v} = e^{\mathbf{j}2\pi \cdot \text{Sigmoid}(\mathbf{x}_r)}.$$
(36)

Finally, the RIS phase matrix  $\mathbf{\Phi} \in \mathbb{C}^{N^{R} \times N^{R}}$  is obtained through diagonalization. The overall process of the RIS refraction phase design, namely, the transformer-based RPDN, can be expressed as

$$\Phi = f_{\rm RIS} \left( \bar{\mathcal{H}}; \mathcal{W}_R \right), \tag{37}$$

where  $f_{RIS}(\cdot)$  denotes the mapping of the RPDN, whose trainable parameter set is  $\mathcal{W}_R$ .

## 484 4.2.3 Knowledge-Data Dual-Driven Digital Beamformer Design

With the known BS-RIS channel  $\mathbf{G}[k]$ , the designed RIS refraction phase matrix  $\mathbf{\Phi}$  and the analog beamforming matrix  $\mathbf{F}_{\mathrm{RF}}$  as well as the estimated RIS-UE channel  $\hat{\mathbf{h}}[u, k]$ , the BS can obtain the estimated equivalent baseband channel  $\hat{\mathbf{h}}_{\mathrm{eq}}[u, k] \in \mathbb{C}^{1 \times M^{\mathrm{B}}}$  as

$$\hat{\mathbf{h}}_{\text{eq}}[u,k] = P_T \hat{\mathbf{h}}[u,k] \mathbf{\Phi} \mathbf{G}[k] \mathbf{F}_{\text{RF}}.$$
(38)

The true equivalent baseband channel  $\mathbf{h}_{eq}[u, k]$  has the similar form to Equation (38), given the designed  $\boldsymbol{\Phi}$  and  $\mathbf{F}_{RF}$ . Thus, the optimization problem Equation (34) can be simplified as

$$\max_{\mathbf{F}_{\mathrm{BB}}[k],\forall k} \quad \frac{1}{K} \sum_{u=1}^{U} \sum_{k=1}^{K} \log_2 \left( 1 + \mathrm{SINR}[u,k] \right),$$
  
s.t. 
$$\operatorname{SINR}[u,k] = \frac{|\mathbf{h}_{\mathrm{eq}}[u,k]\mathbf{f}_{\mathrm{BB}}[u,k]|^2}{\sum_{i=1,i\neq u}^{U} |\mathbf{h}_{\mathrm{eq}}[i,k]\mathbf{f}_{\mathrm{BB}}[i,k]|^2 + \sigma_n^2},$$
  
$$\|\mathbf{F}_{\mathrm{RF}}\mathbf{F}_{\mathrm{BB}}[k]\|_F^2 = M^{\mathrm{B}}, \forall k.$$
(39)

The above problem is a classic baseband beamforming problem, which can be solved with standard
liner beamforming schemes, such as the regularized ZF (RZF) or iterative weighted minimum meansquare error (WMMSE) algorithm.

The iterative WMMSE algorithm solves the optimization (39) by solving the MMSE problem given in (40) below, which has the identical optimal solution  $\mathbf{F}_{BB}[k], \forall k$  to the problem (39).

$$\max_{\mathbf{\bar{U}},\mathbf{\bar{W}},\mathbf{F}_{BB}[k],\forall k} \quad \sum_{u=1}^{U} \sum_{k=1}^{K} \left( \bar{w}_{u,k} e_{u,k} - \log_2 \bar{w}_{u,k} \right), \\
\text{s.t.} \quad \|\mathbf{F}_{RF} \mathbf{F}_{BB}[k]\|_F^2 \leq M^{B}, \forall k,$$
(40)

where  $\bar{w}_{u,k} = \{\bar{\mathbf{W}}\}_{u,k}$  is the weight of the *u*-th user on the *k*-th subcarrier,  $e_{u,k} = \mathbb{E}\{|\hat{s}[u,k] - |\hat{s}[u]\} = |\hat{s}[u]| + |\hat{s}[u]| +$ 495  $s[u,k]|^2$  is the MSE between the transceiver symbols under the independence assumption of s[u,k]496 and n[u,k], while  $\hat{s}[u,k] = \bar{u}_{u,k}y[u,k]$  denotes the estimated data symbol at the UE-side, and 497  $\bar{u}_{u,k} = {\{\bar{\mathbf{U}}\}}_{u,k}$  is the receiver gain of the *u*-th UE on the *k*-th subcarrier. According to [41], the above 498 problem is convex in individual optimization variable. Hence each of the optimization subproblems 499 has a closed-form solution given the other optimization variables, and a block coordinate descent 500 (BCD) iterative algorithm is adopted to solve the optimization (40). This algorithm is summarized 501 in Algorithm 1, where we omit the iteration index t on the variables for clarity. 502

<sup>503</sup> However, the iterative WMMSE algorithm typically imposes a large number of iterations with

Algorithm 1 Iterative WMMSE beamforming design algorithm

- 1: Initialize  $\mathbf{F}_{BB}[k]$  that meets  $\|\mathbf{F}_{RF}\mathbf{F}_{BB}[k]\|_{F}^{2} = M^{B}$ , set the maximum iteration number  $I_{max}$ , and the current iteration index t = 0;
- 2: repeat

3: **Update** 
$$\{\bar{\mathbf{U}}\}_{u,k}$$
:  $\bar{u}_{u,k} = \left(\sum_{i=1}^{U} |\mathbf{h}_{eq}[u,k]\mathbf{f}_{BB}[i,k]|^2 + \sigma_n^2\right)^{-1} \mathbf{h}_{eq}[u,k]\mathbf{f}_{BB}[u,k], \forall u,k;$ 

4: Update 
$$\{\mathbf{W}\}_{u,k}$$
:  $\bar{w}_{u,k} = (1 - \bar{u}_{u,k}^* \mathbf{h}_{eq}[u,k]\mathbf{f}_{BB}[u,k])^{-1}, \forall u,k;$ 

5: **Update** 
$$\mathbf{f}_{BB}[u,k]$$
:  $\mathbf{f}_{BB}[u,k] = \bar{u}_{u,k}\bar{w}_{u,k} \Big(\sum_{i=1}^{U} \bar{w}_{i,k} |\bar{u}_{i,k}|^2 \mathbf{h}_{eq}^{H}[i,k] \mathbf{h}_{eq}[i,k] + \mu_k \mathbf{I} \Big)^{-1} \mathbf{h}_{eq}^{H}[u,k],$ 

where 
$$\mu_k = \sum_{j=1}^{\circ} \frac{\sigma^2}{M^{\mathrm{B}}} \bar{w}_{j,k} |\bar{u}_{j,k}|^2, \forall u, k;$$

- 6: t = t + 1;
- 7: until  $t \ge I_{\max}$
- 8: Scale  $\mathbf{F}_{BB}[k]$  to meet the transmit power constraint.

<sup>504</sup> long running time. Furthermore, the BS can only acquire the imperfect estimated CSI  $\hat{\mathbf{h}}_{eq}[u, k]$ , and <sup>505</sup> it is difficult for the traditional digital beamforming algorithms, such as Algorithm 1, to overcome the <sup>506</sup> interference induced by the imperfect CSI. Thus, we propose the knowledge-data dual-driven digital <sup>507</sup> beamforming network, as shown in Figure 6, which utilizes the transformer encoder to directly learn <sup>508</sup> the parameters of the iterative WMMSE algorithm from the imperfect CSI for better interference <sup>509</sup> elimination and shorter running time.

Specifically, the real-valued 3D matrix  $\bar{\mathcal{H}}$  is reshaped into a 2D matrix  $\bar{\mathbf{H}}_d \in \mathbb{C}^{K \times 2UN^{\mathrm{R}}}$ , which is inputted into the transformer encoder. The output of the transformer encoder  $\mathbf{X} \in \mathbb{C}^{K \times 4U}$  is converted into the weight matrix  $\bar{\mathbf{W}}$  and the receiver gain matrix  $\bar{\mathbf{U}}$ , i.e.,

$$\bar{\mathbf{W}} = \mathbf{X}_{[:,:U]}^{\mathrm{T}} + j\mathbf{X}_{[:,U:2U]}^{\mathrm{T}},\tag{41}$$

$$\bar{\mathbf{U}} = \mathbf{X}_{[:,2U:3U]}^{\mathrm{T}} + \mathbf{j}\mathbf{X}_{[:,3U:]}^{\mathrm{T}}.$$
(42)

Then, we can obtain  $\mathbf{F}_{BB}[k]$ ,  $\forall k$ , based on the learned  $\mathbf{\bar{W}}$  and  $\mathbf{\bar{U}}$  by the update function of  $\mathbf{f}_{BB}[u, k]$ , i.e., line 5 of Algorithm 1. Compared with the iterative WMMSE beamforming design, our proposed scheme does not involve an iterative process so that running time can be reduced significantly. To satisfy the transmit power constraint, the normalization operation can be expressed as

$$\mathbf{F}_{\mathrm{BB}}[k] = \frac{\sqrt{M^{\mathrm{B}}} \mathbf{F}_{\mathrm{BB}}[k]}{\|\mathbf{F}_{\mathrm{RF}} \mathbf{F}_{\mathrm{BB}}[k]\|_{F}}, \forall k.$$
(43)

<sup>514</sup> The proposed knowledge-data dual-driven digital beamformer design can be expressed as

$$\{\mathbf{F}_{\mathrm{BB}}[k], \forall k\} = f_{\mathrm{DBF}} \left( \bar{\mathbf{H}}_d; \mathcal{W}_D \right), \tag{44}$$

sis where  $f_{\text{DBF}}(\cdot)$  is the map of the digital beamforming network with a trainable parameter set  $\mathcal{W}_D$ .

### 516 4.2.4 Training Strategy

We take every U channel samples (i.e., the channels of U UEs) in the training set of the channel 517 estimation stage as a group to form a training set at the beamforming design stage, which is denoted 518 as  $\mathcal{H}_U$ . The number of off-line training samples is  $|\mathcal{H}_U| = N_{\text{set}}/U$ . A sample in  $\mathcal{H}_U$  is an UE set 519  $\{\mathbf{H}[u], 1 \le u \le U\}$ , where  $\mathbf{H}[u]$  is the spatial-frequency channel between the RIS and the u-th UE. 520  $\{\mathbf{H}[u], 1 \leq u \leq U\}$  are inputted to the trained SFDCEtra network to obtain the estimated 521 channels { $\mathbf{H}[u], 1 \leq u \leq U$ }, which form the input to the proposed network. Since imperfect CSI 522 will reduce the sum rate upper bound, to ensure a faster learning process, we apply a teacher forcing 523 technique [42] at the early stage of training by feeding the perfect CSI  $\{\mathbf{H}[u], \forall u\}$  to the proposed 524 network. At the off-line training stage, we consider E2E training to jointly optimize the proposed 525 hybrid beamforming and RIS phase design network, i.e., the parameters of the entire network are 526 trained by minimizing the negative sum rate. Thus, the loss function is written as 527

$$\mathcal{L}_b = -\frac{1}{B_b} \sum_{i=1}^{B_b} R,\tag{45}$$

where R is the sum rate defined in Equation (33) and  $B_b$  is the batch size for off-line training.

## 529 5 Results and Discussion

In this section, we evaluate the performance of the proposed spatial-frequency domain channel extrapolation scheme as well as hybrid beamforming and RIS phase design for the RIS-aided THz massive MIMO system through numerical simulations.

## 533 5.1 Simulation Settings

## 534 5.1.1 Communication Scenario Set up

In simulations, the BS is deployed on the top of a building of height 30 m, and the RIS is installed on 535 a window surface on one floor of another building. As shown in Figure 1(b), the BS (RIS) is equipped 536 with  $M^{\rm B} = M_y^{\rm B} M_z^{\rm B} = 4$   $(M^{\rm R} = M_y^{\rm R} M_z^{\rm R} = 4)$  subarrays on the *yz*-plane, where  $M_y^{\rm B} = 2$   $(M_y^{\rm R} = 2)$ 537 and  $M_z^{\rm B} = 2$  ( $M_z^{\rm R} = 2$ ). Each subarray is a UPA with  $N_{\rm sub}^{\rm B} = N_y^{\rm B} N_z^{\rm B} = 64$  ( $N_{\rm sub}^{\rm R} = N_y^{\rm R} N_z^{\rm R} = 64$ ) isotropically radiating elements, where  $N_y^{\rm B} = 8$  ( $N_y^{\rm R} = 8$ ) and  $N_z^{\rm B} = 8$  ( $N_z^{\rm R} = 8$ ). Therefore, 538 539 the number of elements of the complete array at the BS (RIS) is  $N^{\rm B}$  =  $M^{\rm B}N^{\rm B}_{\rm sub}$  = 256 ( $N^{\rm R}$  = 540  $M^{\rm R}N^{\rm R}_{\rm sub} = 256$ ). For simplicity, we assume that the BS and RIS meet the parallel symmetric array 541 arrangement with a distance of  $D = 20 \,\mathrm{m}$ . The central frequency is  $f_c = 0.3 \,\mathrm{THz}$  with bandwidth 542  $f_s = 1 \text{ GHz}$ . The number of OFDM subcarriers is set to K = 128 and the antenna gain of the BS 543 is  $G_T = 10$  dBi. Given the above parameter settings, the subarray intervals of both the BS and 544



Figure 7: Multi-ray THz channel model for the indoor scenario: the NLoS rays are reflected by the scatters.

the RIS are calculated from Equation (1) as  $d_{sy}^{\rm B}, d_{sz}^{\rm R}, d_{sz}^{\rm R} = 96.5\lambda$  for obtaining the multi-stream multiplexing gain over the LoS path.

Figure 7 depicts the schematic diagram of the fixed scattering environment, where the positions 547 of the RIS, UEs, and scatterers are marked by blue, red, and green circles, respectively. The red 548 solid line represents the LoS link between the RIS and an UE, and the black dotted line indicates the 549 NLoS link via an scatterer. We assume that U = 4 UEs are randomly distributed over the xy-plane 550 of the rectangular room  $(W_x = 5 \text{ m}, W_y = 10 \text{ m})$ , and the height of UEs is 1 m lower than the RIS. 551 Due to rich scatterers for indoor environment as well as the high scattering and diffraction losses in 552 the THz band, the number of available NLoS paths (scatterers) in the THz indoor channel is set to 553  $L_p = 5$ , implying that only a single-bounce scattering mode is considered. We set the parameters 554 of the reflection coefficient  $\beta_{\rm RC}$  as  $\mu_{\rm R} = -5$ ,  $\sigma_{\rm R} = 2$ . The noise power spectrum density at the 555 UEs is  $\sigma_{\text{NSD}}^2 = -174 \text{ dBm/Hz}$ . Thus, the power of the AWGN is  $\sigma_n^2 = \sigma_{\text{NSD}}^2 f_s / K = -105 \text{ dBm}$ . 556 The RIS-UE channel samples are generated using Equation (5), where the UEs and scatterers are 557 distributed randomly each time. 558

#### 559 5.1.2 SFDCEtra Network Parameter Configuration

In the CSI feedback network, the linear embedding layer of the transformer encoder has  $d_{\rm T} = 256$ 560 neurons. In the transformer encoder, the number of the encoder layers is  $L_{\rm T} = 3$ , where the 561 number of heads is h = 8 and the position-wise MLP sub-layer has 2 fully-connected layers with 562  $4d_{\rm T}$  and  $d_{\rm T}$  neurons, respectively, while the dimension of the output linear layer is 2M. In the 563 sub-channel estimation network, the linear layer is a  $2N_s^{\rm R}$ -dimensional fully-connected layer and the 564 hyperparameters of the transformer encoder are the same as those of the CSI feedback network. As 565 for the channel extrapolation network, the output of the zero filling is processed by the convolutional 566 layer with the kernel size of  $7 \times 7$  and the number of filters is 2. The patch size of the rearranged 567 operation is L = 16, the number of patches is  $N_p = 128$  and the number of neurons in the linear 568 layer is  $d_{\rm M} = 512$ . We set the number of mixer layers as  $L_{\rm M} = 6$ , where each mixer layer consists 569 of two MLP blocks, and the numbers of neurons in the MLP blocks are set to  $2N_p$ ,  $N_p$ ,  $2d_M$ , and 570  $d_{\rm M}$ , respectively. The above structural parameters of the SFDCEtra network are empirically found 571 to be appropriate. 572

We divide the data set into the training set, validation set, and test set, which contain 102400, 10240, and 10240 samples, respectively. Unless otherwise specified, the uniform element selection strategy is adopted in the simulations. When considering the learning-based element selection strategy, the weight factor  $\gamma$  is set to 0.9. At the network training stage, the Adam optimizer is adopted to update the network weight parameters and the learning rate varies depending on the *warmup* mechanism [38]. We set the batch size of the training set to 512, and 200 epochs for training.

## 579 5.1.3 HBFRPD Network Parameter Configuration

Again we determine appropriate structural parameters of the HBFRPD network empirically. Specifically, in the RIS phase design network, the linear embedding layer of the transformer encoder has  $d_{\rm B} = 128$  neurons. In the transformer encoder, the number of the encoder layers is  $L_{\rm B} = 3$ , where the number of heads is h = 8 and the position-wise MLP sub-layer has 2 fully-connected layers with  $4d_{\rm B}$  and  $d_{\rm B}$  neurons, respectively, while the output linear layer of the transformer encoder has  $N^{\rm R}/U = 64$  neurons. In the digital beamforming network, the hyperparameters of the transformer encoder are the same as those of the RIS phase design network, and the output linear layer of the transformer encoder has 4U neurons.

We take each U channel samples as a group to form a data set, which is divided into the training set, validation set, and testing set, which contain 25600, 2560, and 2560 samples, respectively. We set the batch size of the training set to 32, and 180 epochs for training.



Figure 8: NMSE performance comparison of different channel estimation schemes versus transmit power  $P_{\rm T}$ .

## <sup>591</sup> 5.2 DL-Based Spatial-Frequency Domain Channel Extrapolation

Since the RIS element can only passively receive EM waves, selecting partial elements of array would 592 reduce the signal energy radiated into the room. For the fair comparison between different schemes, 593 we adopt the same transmit power instead of the same SNR as the comparison criterion to avoid 594 ignoring the performance differences induced by the number of activated RIS elements. Specifically, 595 as shown in Figure 8, we plot the NMSE performance of the different schemes as a function of 596 transmit power  $P_{\rm T}$ . The number of NLoS paths is  $L_p = 5$ . We consider three model-based chan-597 nel estimation benchmark algorithms, namely, the simultaneous orthogonal match pursuit (SOMP) 598 algorithm [43], the multiple-measurement-vector approximate message passing (MMV-AMP) algo-599 rithm [44], and the model-driven DL-based channel estimation scheme using the MMV learned AMP 600 (MMV-LAMP) network [45], which utilize M = 64 OFDM symbols on all subcarriers and then di-601

rectly estimate the complete channel. For the SOMP and MMV-LAMP schemes, the redundant 602 dictionary with an oversampling ratio of 4 is considered to further improve the performance, i.e., the 603 number of codewords is  $G_d = 1024$ . Note that the MMV-AMP scheme requires the measurement 604 matrix's elements to be independent and identically distributed, and hence we cannot consider the 605 redundant dictionary (i.e.,  $G_d = N^{\rm R} = 256$ ). Since data-driven DL algorithms have the potential to 606 achieve better performance, we also compare our proposed DL-based SFDCEtra network with the 607 transformer-based channel estimation network (Transformer-CEN) [39] and the CNN-based chan-608 nel extrapolation network (CNN-CEtraN) [16]. For these methods, we set the number of OFDM 609 symbols to M = 16 and the subcarrier compression ratio to  $\bar{\rho} = 16$ . The transformer-based scheme 610 collects the signals from all RIS elements, i.e.,  $\rho = 1$ , and directly estimates the complete channel. 611 Both the CNN-based and our proposed channel extrapolation schemes consider the element com-612 pression ratio of  $\rho = 4$  to perform partial channel extrapolation. Note that for fairness, the above 613 model- and data-driven algorithms do not consider the quantization of CSI feedback information. 614 Therefore, we additionally consider the proposed scheme with B = 256 feedback bits generated via 615 a 2-bit quantizer, denoted as 'Proposed-Q'. 616

It can be observed from Figure 8 that the proposed channel extrapolation scheme outperforms the other channel estimation schemes considerably in terms of NMSE performance while imposing a smaller pilot overhead. This is because exploiting the spatial-frequency correlation allows our DL-based channel extrapolation scheme to recover the unobserved channel part from the estimated low-dimensional sub-channel, thus reducing the training overhead while improving the NMSE performance. In particular, our extrapolation scheme significantly improves the NMSE performance compared with the state-of-the-art CNN-based channel extrapolation scheme. Unlike local percep-



Figure 9: NMSE performance comparison of the proposed scheme versus the number of multipath  $L_p$ , given  $\rho = 4$ ,  $\bar{\rho} = 16$  and M = 16. Offline training is based on the channel samples with  $L_p = 5$  multipath components.

tion in CNN, the MLP-mixer is utilized as the backbone of our channel extrapolation module and it can extract the global features of the channel for enhanced extrapolation accuracy. Considering the actual situation of finite quantized feedback, we can see that our proposed scheme with 2-bit quantizer, 'Proposed-Q', can still achieve very good performance. These results demonstrate that the proposed channel extrapolation scheme can learn latent features from the data more effectively to achieve better channel estimation accuracy with less pilot and feedback overhead.

We further investigate the robustness of the proposed channel extrapolation scheme as a function of the number of multipath  $L_p$  in Figure 9. Note that the proposed DL-based channel extrapolation scheme is trained at the offline training stage, based on the channel samples with  $L_p = 5$  multipath components. It can be clearly seen from Figure 9 that at the online estimation stage, the proposed scheme can adapt to estimate multipath channels with  $L_p \neq 5$ , without having to retrain the entire network architecture. Therefore, the proposed DL-based channel extrapolation enjoys better robustness and generalization capability to various channel conditions.



Figure 10: NMSE performance comparison of the proposed scheme with different pilot numbers versus transmit power  $P_{\rm T}$ , given  $\rho = 4$ ,  $\bar{\rho} = 16$ ,  $L_p = 5$ .

In Figure 10, we investigate the channel extrapolation NMSE performance of our proposed scheme 637 with different numbers of pilot OFDM symbols, M = 4, 8, 16, 32 and 64. As expected, the channel 638 extrapolation performance improves with the increase of the number of pilot OFDM symbols. This 639 is because more pilot OFDM symbols can improve the accuracy of sub-channel estimation, thus 640 reducing the error propagation and improving the reconstruction of the extrapolation module. Fur-641 thermore, we can see that the proposed scheme can provide more significant performance gain by 642 increasing the number of pilot OFDM symbols in the case of low transmit power. This is because 643 the increase in the number of observations can improve the received SNR. 644

Figure 11 depicts the NMSE performance of the proposed DL-based channel extrapolation scheme versus the element compression ratio  $\rho$ , with three ESEs. Specifically, the curve labeled by 'Uniform'



Figure 11: NMSE performance comparison of the proposed scheme with different element selection strategies versus element compression ratio  $\rho$ , given  $\bar{\rho} = 16$ ,  $L_p = 5$ , M = 16,  $P_T = 44 \text{ dBm}$ .

corresponds to the uniform selection strategy, the curve labeled by 'Random' represents the random 647 selection strategy, while the other two marked by 'DL-based with 200 epochs' and 'DL-based with 648 300 epochs' use the DL-based element selection strategy. As expected, the NMSE improves as the 649 element compression ratio  $\rho$  decreases. This is largely due to two reasons: 1) As the number of 650 selected RIS elements increases, or the element compression ratio  $\rho$  decreases, the received signal 651 power will increase, thus improving the estimation accuracy of channel extrapolation input (i.e., sub-652 channel estimate), and 2) The received pilot signal can provide more channel information when more 653 RIS elements are selected. However, this does not imply that we can obtain the best performance by 654 choosing the lowest element compression ratio (or performing complete observations directly without 655 extrapolation). Indeed, the channel extrapolation performance heavily depends on the number of 656 transmission resources, the accuracy of the sub-channel estimate, and the number of selected RIS 657 elements (i.e., the dimension of the sub-channel). Only when the transmission resources are sufficient, 658 can the gain provided by more selected RIS elements be seen clearly. Moreover, it can be seen 659 that the performance difference between different element selection strategies is not obvious at low 660 compression ratios. Only at high compression ratio ( $\rho > 8$ ), can the performance difference be seen 661 clearly as 'Uniform' < 'Random' < 'DL-based'. Since the aperture of the random pattern is usually 662 larger than that of the uniform pattern, the random selection strategy is slightly better than that of 663 the uniform selection strategy. The performance of the DL-based approach is relatively better than 664 that of the first two approaches only when its training reaches sufficient epochs, as the learning of 665 the selection network needs sufficient number of epochs to converge. 666

To fully illustrate the effectiveness of the proposed DL-based channel extrapolation solution, we verify its channel extrapolation module separately. To do so, we fix the compression ratio of RIS elements to 4, i.e.,  $N_s^{\rm R} = 64$ . First, the least squares (LS), the SOMP, and the proposed



Figure 12: (a) NMSE performance comparison of different sub-channel estimation schemes versus transmit power  $P_{\rm T}$ ; and (b) NMSE performance of channel extrapolation versus transmit power  $P_{\rm T}$  for different sub-channel estimation schemes.

transformer-based algorithm are utilized for sub-channel estimation, and the results are shown in 670 Figure 12(a). Observe that the NMSE of the SOMP-based sub-channel estimation with M = 64 pilot 671 symbols is significantly better than that of the LS-based sub-channel estimation with M = 64 pilot 672 symbols, particularly at low transmit power  $P_{\rm T}$ . Furthermore, the NMSE of our transformer-based 673 sub-channel estimation algorithm with only M = 16 pilot symbols is considerably better than that 674 of the SOMP-based sub-channel estimation with M = 64 pilot symbols. Then, we input the sub-675 channels estimated by different algorithms into the trained channel extrapolation network  $f_{\rm SFDE}(\cdot)$ , 676 which outputs the estimation of the complete channel. The corresponding results are shown in 677 Figure 12(b). Observe that the NMSE performance of the complete channel extrapolated from our 678 channel extrapolation network is even better than the NMSE of the estimated low-dimensional sub-679 channel, without any additional pilot overhead. This shows that our proposed channel extrapolation 680 network can not only be used for DL-based communication architecture, but also be combined with 681 traditional algorithms to significantly reduce resource overhead. Therefore, we conclude that the 682 proposed DL-based spatial-frequency domain channel extrapolation scheme can significantly reduce 683 the pilot overhead while achieving the same or better channel estimation NMSE performance. 684

## 5.3 DL-Based Hybrid Beamforming and RIS Phase Design

Figure 13 shows the sum rates achieved by different schemes under the perfect CSI case. We considered two comparison schemes, both of which adopt the analog beamforming design discussed in Subsection 4.2.1 as well as the beam alignment-based RIS phase design. In the beam alignmentbased RIS phase design, each subarray of the RIS selects one UE and performs beam alignment according to the RIS-UE CSI on the central subcarrier to concentrate and refract the signal energy to it, while ignoring interference to other UEs. For digital beamforming design, these two comparison



Figure 13: Sum rates achieved by different schemes versus transmit power  $P_{\rm T}$  given the perfect CSI. The actual transmit power of the 'w/o LoS MIMO' case is  $UP_{\rm T} = 4P_{\rm T}$ .

schemes adopt the RZF and iterative WMMSE algorithms, respectively, thus they are abbreviated 692 as 'RZF' and 'WMMSE', respectively. It can be observed that our proposed HBFRPD scheme has 693 the better performance than other schemes and the superiority is more evident as the transmit power 694 increases. In addition, the proposed HBFRPD scheme does not require to obtain  $\mathbf{F}_{BB}[k], \forall k$ , in an 695 iterative manner. Thus, it runs much faster than the iterative WMMSE algorithm. We also analyze 696 the performance gain provided by LoS MIMO architecture. Considering the case of the BS and the 697 RIS w/o LoS MIMO array structure (i.e., both use UPA arrays), the BS-RIS channel is a single 698 LoS path with rank 1, which only provides single stream data transmission. To ensure fairness, the 699 transmit power of the 'w/o LoS MIMO' case is equal to the total transmit power of the 'w/ LoS 700 MIMO' cases, i.e., the transmit power of the 'w/o LoS MIMO' case is actually  $UP_{\rm T} = 4P_{\rm T}$ . By 701 calculating the sum rate, we obtain the green curve in Figure 13. It can be seen that the sum rate 702 in the 'w/ LoS MIMO' case is much higher than that of 'w/o LoS MIMO' case. This is because the 703 LoS MIMO architecture can increase the sum rate linearly benefited from extra spatial multiplexing 704 gain, while 'w/o LoS MIMO' case can only provide log-level growth as the SINR increases. 705

Although most schemes can achieve high sum rate performance under perfect CSI, the sum rate 706 of multi-users is actually limited by the inter-user interference induced by CSI error. Figure 14(a) 707 illustrates the sum rate performance of the different schemes with imperfect CSIs estimated at two 708 different transmit powers  $P_{\rm T}$  (CE). Compared with the case of perfect CSI, the sum rate degrades 709 significantly with the decrease of CSI estimation accuracy, i.e., with the decrease of the transmit 710 power at the channel estimation stage. It can be clearly seen that due to the inter-user interference 711 induced by CSI errors, the sum rates of the RZF and iterative WMMSE schemes barely increases 712 with transmit power. Moreover, our proposed HBFRPD scheme exhibits a significant performance 713 gain over the RZF and iterative WMMSE algorithms in the presence of CSI estimation errors. This 714



Figure 14: (a) Sum rates achieved by different schemes versus transmit power  $P_{\rm T}$  under the imperfect CSI case, and (b) The CDFs of the sum rates achieved by different schemes under the imperfect CSI case, given  $P_{\rm T} = 44$  dBm. We have  $\rho = 4$ ,  $\bar{\rho} = 16$ ,  $L_p = 5$  and M = 16.

result indicates that our proposed scheme is capable of mitigating the interference caused by CSI
 errors and hence has better robustness to inaccurate CSI than the other schemes.

Figure 14(b) shows the cumulative distribution functions (CDFs) that characterize the sum rate 717 performance achieved by the different schemes. Here, we consider the transmit power  $P_{\rm T} = 44 \, {\rm dBm}$ 718 at the data transmission stage. It can be seen from Figure 14(b) that when the transmit power is 719  $P_{\rm T}(\rm CE) = 34 \, \rm dBm$  at the channel estimation stage, the proposed HBFRPD network has a probability 720 of about 64.6% to achieve a sum rate exceeding 30 bps/Hz, while the other two schemes can only 721 achieve 16.3%. When the transmit power is  $P_{\rm T}({\rm CE}) = 44 \, {\rm dBm}$  at the channel estimation stage, 722 our HBFRPD network has a probability of about 68.8% to achieve a sum rate exceeding 40 bps/Hz, 723 which is significantly better than the other two schemes. This result again confirms the superior 724 performance of the proposed DL-based HBFRPD network over the existing conventional schemes. 725

## 726 5.4 Computational Complexity Analysis

We now investigate the computational complexity. For the DL-based schemes, since there is no strict time limit at the offline training stage, we only consider the computational complexity at the inference stage. The computational complexity analysis of different schemes is presented in Table 1. All the numerical results are obtained on a PC with Intel(R) Core(TM) i9-10980XE CPU @ 3.00GHz and an Nvidia GeForce RTX 3090 GPU. The DL-based methods and the existing solutions are implemented on the PyCharm framework. The details are further elaborated as follows.

<sup>733</sup> 1) Channel estimation schemes: In the SOMP algorithm [43], correlation operation imposes <sup>734</sup> significant computational complexity, where  $G_d$  is the dimension of the redundant dictionary and I<sup>735</sup> is the number of iterations. The MMV-AMP algorithm [44] mainly requires matrix multiplication <sup>736</sup> operations, but a large number of iterations I increases its computational complexity. The MMV-

Channel estimation scheme	Complexity	FLOPs	Run time/s
SOMP	$\mathcal{O}\left(G_d KMI + G_d^2 KI\right)$	4.707 G	0.1130
MMV-AMP	$O\left(MKN^{\mathrm{R}}I\right)$	$5.337~{ m G}$	0.7482
MMV-LAMP	$\mathcal{O}\left(MG_{d}KI ight)$	0.341 G	0.0689
Transformer-CEN	$\mathcal{O}\left(L_{\mathrm{T}}(Kd_{\mathrm{T}}^{2}+K^{2}d_{\mathrm{T}}) ight)$	$0.362~{ m G}$	0.0103
CNN-CEtraN	$\mathcal{O}\left(Z_5^2 N^{\mathrm{R}} K C_{32}^2\right)$	10.16 G	0.6039
Proposed	$\mathcal{O}\left(L_{\mathrm{M}}(N_{p}^{2}d_{\mathrm{M}}+N_{p}d_{\mathrm{M}}^{2}) ight)$	1.066 G	0.0248
Beamforming scheme	Complexity	FLOPs	Run time/s
RZF	$\mathcal{O}\left((2U(M^{\mathrm{B}})^{2} + (M^{\mathrm{B}})^{3})K\right)$	$24.58~{\rm K}$	0.1389
WMMSE	$\mathcal{O}\left(IK(U^2(M^{\mathrm{B}})^2 + U(M^{\mathrm{B}})^3)\right)$	8.192 M	0.9184
Proposed	$\mathcal{O}\left(L_{\rm B}U(Kd_{\rm B}^2+K^2d_{\rm B})\right)$	0.512 G	0.0616

Table 1: Computational Complexity of Different Schemes.

LAMP algorithm [45] has a low computational complexity because DL reduces the required number 737 of iterations. The Transformer-CEN [39] also has a low computational complexity, and the main 738 sources of its computational complexity come from self-attention and MLP sublayers. In the CNN-739 CEtraN [16], convolutional layers introduce significant computational complexity. By contrast, the 740 MLP-mixer layers provide the majority of the computational complexity in our proposed SFDCEtra 741 network, which is much lower than that of the CNN-CEtraN. We further meticulously count the 742 numbers of floating-point operations per second (FLOPs) and run times per sample on CPU for 743 different schemes in Table 1. Observe that at the inference stage, the FLOPs and run time per 744 sample of the proposed scheme are lower than most benchmarks. Specifically, our SFDCEtra network 745 imposes the second lowest run time per sample, and only the MMV-LAMP and Transformer-CEN 746 have lower FLOPs than our proposed scheme. 747

2) Beamforming schemes: A matrix inversion is required in the RZF algorithm, which is its main 748 source of computational complexity. In the iterative WMMSE algorithm [41], a large number of 749 iterations increases the computational complexity and the run time per sample. In the proposed DL-750 based HBFRPD Network, self-attention and MLP sublayers impose higher computational complexity 751 and FLOPs than the other two algorithms. However, the run time per sample of our proposed scheme 752 is significantly lower than that of the two model-based schemes. This is due to the fact that the 753 DL-based HBFRPD network just needs matrix multiplication operations and does not requires an 754 iterative procedure. This is a superior advantage of our DL-based HBFRPD network. 755

## 756 6 Conclusions

In this paper, we have proposed a DL-based transmission scheme for RIS-aided THz massive 757 MIMO systems over hybrid-field channels. Our novel twofold contribution has been to develop a 758 channel estimation scheme with low pilot overhead and to design a robust beamforming scheme. 759 More specifically, we have first proposed an E2E DL-based channel estimation framework, which 760 consists of pilot design, CSI feedback, sub-channel estimation, and channel extrapolation. Then, to 761 maximize the sum rate of all UEs under imperfect CSI, we have developed a DL-based scheme to 762 simultaneously design the hybrid beamforming and RIS phase. Simulation results have shown that 763 our proposed channel extrapolation scheme significantly outperforms the existing state-of-the-art 764

schemes, in terms of reconstruction performance, while imposing a much reduced pilot overhead. Moreover, the results have also demonstrated that our proposed beamforming scheme is superior over the existing designs in terms of achievable sum rate performance and robustness to imperfect CSI. Potential future research directions based on the outcomes of this paper include the practical phase shift model of reflecting elements, the analysis of hardware impairments, the analysis of the complex near-field channel, and sensing-aided communications.

## 771 References

- [1] H. Elayan, O. Amin, B. Shihada, R. M. Shubair, and M.-S. Alouini, "Terahertz band: The last piece of RF spectrum puzzle for communication systems," *IEEE Open J. Commun. Soc.*, vol. 1, pp. 1-32, 2020.
- [2] C. Lin and G. Y. Li, "Indoor terahertz communications: How many antenna arrays are needed?," *IEEE Trans. Wirel. Commun.*, vol. 14, no. 6, pp. 3097-3107, Jun. 2015.
- [3] F. Sohrabi and W. Yu, "Hybrid digital and analog beamforming design for large-scale antenna arrays," *IEEE J. Sel. Top. Signal Process.*, vol. 10, no. 3, pp. 501-513, Apr. 2016.
- [4] M. Di Renzo, et al., "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, Nov. 2020.
- [5] K. T. Selvan and R. Janaswamy, "Fraunhofer and fresnel distances: Unified derivation for
   aperture antennas," *IEEE Antennas Propag. Mag.*, vol. 59, no. 4, pp. 12–15, Aug. 2017.
- [6] M. Cui and L. Dai, "Channel estimation for extremely large-scale MIMO: Far-field or near-field?," *IEEE Trans. Commun.*, vol. 70, no. 4, pp. 2663-2677, Apr. 2022.
- [7] L. Yan, Y. Chen, C. Han, and J. Yuan, "Joint inter-path and intra-path multiplexing for
  terahertz widely-spaced multi-subarray hybrid beamforming systems," *IEEE Trans. Commun.*,
  vol. 70, no. 2, pp. 1391-1406, Feb. 2022.
- [8] D. Mishra and H. Johansson, "Channel estimation and low-complexity beamforming design for passive intelligent surface assisted MISO wireless energy transfer," *Proc. ICASSP 2019* (Brighton, UK), May 12-17, 2019, pp. 4659-4663.
- [9] P. Wang, J. Fang, H. Duan, and H. Li, "Compressed channel estimation for intelligent reflecting surface-assisted millimeter wave systems," *IEEE Signal Process. Lett.*, vol. 27, pp. 905–909, May 2020.
- [10] L. Wei, C. Huang, G. C. Alexandropoulos, C. Yuen, Z. Zhang, and M. Debbah, "Channel esti mation for RIS-empowered multi-user MISO wireless communications", *IEEE Trans. Commun.*,
   vol. 69, no. 6, pp. 4144–4157, Jun. 2021.

- [11] A. M. Elbir, A. Papazafeiropoulos, P. Kourtessis, and S. Chatzinotas, "Deep channel learning
   for large intelligent surfaces aided mm-Wave massive MIMO systems," *IEEE Wirel. Commun. Lett.*, vol. 9, no. 9, pp. 1447-1451, Sep. 2020.
- [12] S. Liu, Z. Gao, J. Zhang, M. D. Renzo, and M. -S. Alouini, "Deep denoising neural network assisted compressive channel estimation for mmWave intelligent reflecting surfaces," *IEEE Trans. Veh. Technol.*, vol. 69, no. 8, pp. 9223-9228, Aug. 2020.
- [13] S. Zhang, et al., "Deep learning based channel extrapolation for large-scale antenna systems:
  Opportunities, challenges and solutions," *IEEE Wirel. Commun.*, vol. 28, no. 6, pp. 160-167,
  Dec. 2021.
- [14] B. Lin, F. Gao, S. Zhang, T. Zhou, and A. Alkhateeb, "Deep learning-based antenna selection
  and CSI extrapolation in massive MIMO systems," *IEEE Trans. Wirel. Commun.*, vol. 20,
  no. 11, pp. 7669-7681, Nov. 2021.
- [15] M. Xu, S. Zhang, C. Zhong, J. Ma, and O. A. Dobre, "Ordinary differential equation-based CNN
  for channel extrapolation over RIS-assisted communication," *IEEE Commun. Lett.*, vol. 25, no. 6, pp. 1921-1925, Jun. 2021.
- [16] S. Zhang, S. Zhang, F. Gao, J. Ma, and O. A. Dobre, "Deep learning-based RIS channel
  extrapolation with element-grouping," *IEEE Wirel. Commun. Lett.*, vol. 10, no. 12, pp. 26442648, Dec. 2021.
- [17] K. Ying, et al., "GMD-based hybrid beamforming for large reconfigurable intelligent surface
  assisted millimeter-wave massive MIMO," *IEEE Access*, vol. 8, pp. 19530-19539, Jan. 2020.
- [18] C. Pradhan, A. Li, L. Song, B. Vucetic, and Y. Li, "Hybrid precoding design for reconfigurable
  intelligent surface aided mmWave communication systems," *IEEE Wirel. Commun. Lett.*, vol. 9,
  no. 7, pp. 1041-1045, Jul. 2020.
- [19] B. Di, et al., "Hybrid beamforming for reconfigurable intelligent surface based multi-user communications: Achievable rates with limited discrete phase shifts," *IEEE J. Sel. Areas Commun.*,
  vol. 38, no. 8, pp. 1809-1822, Aug. 2020.
- [20] Y. Ahn and B. Shim, "Deep learning-based beamforming for intelligent reflecting surfaceassisted mmWave systems," *Proc. ICTC 2021* (Jeju Island, Korea), Oct. 20-22, 2021, pp. 17311734.
- [21] S. Zhang, et al., "Beyond intelligent reflecting surfaces: Reflective-transmissive metasurface
  aided communications for full-dimensional coverage extension," *IEEE Trans. Veh. Technol.*,
  vol. 69, no. 11, pp. 13905–13909, Sep. 2020.
- [22] Y. Youn, et al., "Demo: Transparent intelligent surfaces for sub-6 GHz and mmWave B5G/6G
  systems," in Proc. ICC Workshops 2022 (Seoul, Korea), May 16-20, 2022, pp. 1-2.

- [23] D. Kitayama, *et al.*, "Transparent dynamic metasurface for a visually unaffected reconfigurable
   intelligent surface: Controlling transmission/reflection and making a window into an RF lens,"
   *Opt. Express*, vol. 29, no. 18, pp. 29292-29307, Aug. 2021.
- [24] Y. Chen, L. Yan, and C. Han, "Hybrid spherical- and planar-wave modeling and DCNN-powered
  estimation of terahertz ultra-massive MIMO channels," *IEEE Trans. Commun.*, vol. 69, no. 10,
  pp. 7063-7076, Oct. 2021.
- [25] X. Wang, Z. Lin, F. Lin, and L. Hanzo, "Joint hybrid 3D beamforming relying on sensor-based training for reconfigurable intelligent surface aided terahertz-based multiuser massive MIMO systems," *IEEE Sens. J.*, vol. 22, no. 14, pp. 14540-14552, Jul. 2022.
- [26] S. Hong, et al., "Robust transmission design for intelligent reflecting surface-aided secure communication systems with imperfect cascaded CSI," *IEEE Trans. Wirel. Commun.*, vol. 20, no. 4, pp. 2487-2501, Apr. 2021.
- <sup>844</sup> [27] Z. Chen, et al., "Robust hybrid beamforming design for multi-RIS assisted MIMO system with imperfect CSI," *IEEE Trans. Wirel. Commun.*, Early Access, Nov. 2022.
  <sup>846</sup> DOI:10.1109/TWC.2022.3222218
- <sup>847</sup> [28] W. Xu, L. Gan, and C. Huang, "A robust deep learning-based beamforming design for RIS<sup>848</sup> assisted multiuser MISO communications with practical constraints," *IEEE Trans. Cogn. Com-*<sup>849</sup> mun. Netw., vol. 8, no. 2, pp. 694-706, Jun. 2022.
- [29] P. Larsson, "Lattice array receiver and sender for spatially orthonormal MIMO communication,"
  in *Proc. VTC 2005*, (Stockholm, Sweden), May. 30-Jun. 1, 2005, pp. 192-196.
- [30] F. Bohagen, P. Orten, and G. E. Oien, "Optimal design of uniform planar antenna arrays for
  strong line-of-sight MIMO channels," *Proc. SPAWC 2006* (Cannes, France), Jul. 02-05, 2006,
  pp. 1-5.
- [31] X. Song, W. Rave, N. Babu, S. Majhi, and G. Fettweis, "Two-level spatial multiplexing using
  hybrid beamforming for millimeter-wave backhaul," *IEEE Trans. Wirel. Commun.*, vol. 17, no.
  7, pp. 4830-4844, Jul. 2018.
- [32] L. Yan, C. Han, and J. Yuan, "Energy-efficient dynamic-subarray with fixed true-time-delay
  design for terahertz wideband hybrid beamforming" *IEEE J. Sel. Areas Commun.*, vol. 40,
  no. 10, pp. 2840-2854, Oct. 2022.
- [33] C. Han, A. O. Bicen, and I. F. Akyildiz, "Multi-ray channel modeling and wideband characterization for wireless communications in the terahertz band," *IEEE Trans. Wirel. Commun.*, vol. 14, no. 5, pp. 2402–2412, May 2015.
- Y. Wu, J. Kokkoniemi, C. Han, and M. Juntti, "Interference and coverage analysis for terahertz networks with indoor blockage effects and line-of-sight access point association," *IEEE Trans. Wirel. Commun.*, vol. 20, no. 3, pp. 1472–1486, Mar. 2021.

- [35] J. M. Jornet and I. F. Akyildiz, "Channel modeling and capacity analysis for electromagnetic
  wireless nanonetworks in the terahertz band," *IEEE Trans. Wirel. Commun.*, vol. 10, no. 10,
  pp. 3211–3221, Oct. 2011.
- B. Wang, F. Gao, S. Jin, H. Lin, and G. Y. Li, "Spatial- and frequency-wideband effects
  in millimeter-wave massive MIMO systems," *IEEE Trans. Signal Process.*, vol. 66, no. 13,
  pp. 3393-3406, Jul. 2018.
- <sup>873</sup> [37] C. Wen, W. Shih, and S. Jin, "Deep learning for massive MIMO CSI feedback," *IEEE Wirel.* <sup>874</sup> *Commun. Lett.*, vol. 7, no. 5, pp. 748-751, Oct. 2018.
- [38] A. Vaswani, et al., "Attention is all you need," in Proc. NIPS 2017 (Long Beach, CA, USA),
   Dec. 4-9, 2017, pp. 6000–6010.
- [39] Y. Wang, et al., "Transformer-empowered 6G intelligent networks: From massive MIMO processing to semantic communication," *IEEE Wirel. Commun.*, Early Access, Nov. 2022.
  DOI:10.1109/MWC.008.2200157
- [40] I. O. Tolstikhin, *et al.*, "Mlp-mixer: An all-mlp architecture for vision," in *Proc. NIPS 2021*,
   Dec. 6-14, 2021, pp. 24261-24272.
- <sup>882</sup> [41] Q. Shi, M. Razaviyayn, Z. Luo, and C. He, "An iteratively weighted MMSE approach to distributed sum-utility maximization for a MIMO interfering broadcast channel," *IEEE Trans. Signal Process.*, vol. 59, no. 9, pp. 4331-4340, Sep. 2011.
- <sup>885</sup> [42] R. J. Williams and D. Zipser, "A learning algorithm for continually running fully recurrent <sup>886</sup> neural networks," *Neural Computation*, vol. 1, no. 2, pp. 270-280, Jun. 1989.
- [43] C. Tsai, Y. Liu, and A. Wu, "Efficient compressive channel estimation for millimeter-wave large-scale antenna systems," *IEEE Trans. Signal Process.*, vol. 66, no. 9, pp. 2414-2428, May 2018.
- [44] M. Ke, Z. Gao, Y. Wu, X. Gao, and R. Schober, "Compressive sensing based adaptive active
  user detection and channel estimation: Massive access meets massive MIMO," *IEEE Trans. Signal Process.*, vol. 68, pp. 764–779, Jan. 2020.
- [45] X. Ma, Z. Gao, F. Gao, and M. Di Renzo, "Model-driven deep learning based channel estimation and feedback for millimeter-wave massive hybrid MIMO systems," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 8, pp. 2388-2406, Aug. 2021.