

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 Wang¹, Zhen Gao^{1,2,3*}, Sheng Chen^{3,4}, Chun Hu¹, and Dezhi Zheng¹

5 ¹MIT Key Laboratory of Complex-Field Intelligent Sensing, Beijing Institute of
6 Technology, Beijing, China

7 ²Yangtze Delta Region Academy of Beijing Institute of Technology, Beijing Institute
8 of Technology, Jiaxing, China

9 ³Advanced Research Institute of Multidisciplinary Science, Beijing Institute of
10 Technology, Jinan, China

11 ⁴School of Electronics and Computer Science, University of Southampton,
12 Southampton, U.K.

13 ⁵Faculty of Information Science and Engineering, Ocean University of China,
14 Qingdao, China

15 *Address correspondence to: gaozhen16@bit.edu.cn

16 **Abstract**

17 The reconfigurable intelligent surface (RIS) is a promising novel technology for Terahertz
18 (THz) massive multiple-input multiple-output (MIMO) communication systems. However,
19 the acquirement of the high-dimensional channel state information (CSI) and the efficient ac-
20 tive/passive beamforming for RIS are challenging due to its cascaded channel structure and its
21 lack of signal processing units. To address this, we propose a deep learning (DL)-based phys-
22 ical signal processing scheme for RIS-aided THz massive MIMO systems over hybrid far-near
23 field channels, where a channel estimation scheme with low pilot overhead and a robust beam-
24 forming scheme are conceived. Specifically, we first propose an end-to-end DL-based channel
25 estimation framework, which consists of pilot design, CSI feedback, sub-channel estimation, and
26 channel extrapolation. Specifically, we first only activate partial RIS elements and estimate a
27 sub-sampling RIS channel, and then utilize a DL-based extrapolation network to reconstruct
28 the full-dimensional CSI. Moreover, to maximize the sum rate under imperfect CSI, a DL-based
29 scheme is developed to simultaneously design the hybrid active beamforming at the BS and
30 passive beamforming at the RIS. Simulation results show that our proposed channel extrapola-
31 tion scheme has better CSI reconstruction performance than conventional schemes while greatly

32 reducing pilot overhead and our proposed beamforming scheme has superior performance over
33 conventional schemes in terms of robustness to imperfect CSI.

34 **Keywords**

35 Reconfigurable intelligent surface (RIS), Terahertz (THz), hybrid-field, channel extrapolation, hybrid
36 beamforming, deep learning (DL).

37 **Introduction**

38 **1 Introduction**

39 The surge in demand for wireless data traffic in recent years, owing to the exponential growth
40 of massive internet-of-things (IoT) devices and broadband multimedia applications, has necessitated
41 the exploration of Terahertz (THz) communications as a viable solution [1]. However, extremely
42 high free-space losses and strong atmospheric attenuation in the THz band pose a challenge to
43 the long-range coverage of THz communication systems. To overcome this problem, the massive or
44 ultra-massive multiple-input multiple-output (MIMO) technique has been considered to achieve high
45 array gain and mitigate the high propagation loss [2]. Conventional massive MIMO systems require a
46 dedicated radio frequency (RF) chain for each antenna (i.e., fully-digital architecture), which suffers
47 from extremely high power consumption and hardware costs. To circumvent this technical hurdle,
48 the hybrid analog-digital massive MIMO architecture has been widely adopted to reduce the number
49 of RF chains while ensuring high array gains [3].

50 Besides, the advent of reconfigurable intelligent surfaces (RIS) has garnered attention as a poten-
51 tially transformative technology for improving communication performance [4–10]. By manipulating
52 the phase and amplitude of RIS phase shifters, RIS passively reflects incident electromagnetic (EM)
53 signals towards desired directions and provides significant beamforming gain. More importantly, RIS
54 does not require power-intensive RF chains, which contributes to a more environmentally friendly
55 and cost-effective communication solution. Therefore, the integration of RIS and massive MIMO
56 techniques holds promise for overcoming the limitations of THz communications and realizing its
57 full potential.

58 Generally, a simplified planar-wave channel model is appropriate in the case that the user equip-
59 ment (UE) works in the far field of the base station (BS). However, since severe path loss will reduce
60 the effective coverage while the increasing array size in the THz band will increase the Rayleigh
61 distance [11], both far and near-field need to be considered for THz massive MIMO systems. There-
62 fore, the distance from each antenna of the BS to the UE needs to be considered under near-field
63 conditions by the spherical-wave channel model [12]. On the other hand, the number of spherical-
64 wave channel parameters is proportional to the number of massive antennas, which indicates that
65 directly adopting the spherical-wave channel model in THz massive MIMO systems is unrealistic.
66 To this end, a hybrid-field (hybrid spherical- and planar-wave) channel model characterized by a

67 smaller number of parameters while maintaining high accuracy has been proposed for THz massive
68 MIMO systems [13]. For such a channel model, the EM signal is modeled as a spherical wave for the
69 inter-subarray and a planer wave for the intra-subarray, based on different subarray architectures.
70 Although the application of RISs has been widely researched recently [14–18], the utilization of RIS
71 for THz massive MIMO communications over hybrid-field channels is still at its early study stage.

72 1.1 Related Work

73 Acquiring accurate channel state information (CSI) is critical in establishing RIS-aided commu-
74 nication systems [19–22]. However, accurately estimating high-dimensional CSI with limited pilot
75 signals remains a formidable challenge [14]. To address this challenge, compressive sensing (CS)-
76 based solutions have been proposed to reduce the pilot overhead by leveraging the channel sparsity
77 [15, 16]. However, these solutions present challenges with regard to computational complexity and
78 storage requirements, as the corresponding matrix inversion and iterative operations. Recently, the
79 integration of deep learning (DL) in communication systems has garnered extensive attention. For
80 instance, in [23], the authors proposed an effective pilot reduction technique by gradually pruning
81 less significant neurons from the dense layers during training. In [17], the authors designed a DL-
82 based channel estimation network to acquire the RIS-aided channel and the non-RIS-aided channel.
83 In [18], a semi-passive RIS architecture was proposed, where the orthogonal match pursuit (OMP)
84 algorithm and a denoising convolutional neural network (CNN) are applied to reconstruct the CSI.
85 However, the deployment of RF chains negates the key benefits of the RIS, i.e., reducing hardware
86 costs and power consumption.

87 In fact, due to the highly-dense arrangement of RIS elements [24], there is a strong correlation
88 between the different elements of the CSI matrix, which makes it possible to extrapolate the complete
89 channel from a partial one, i.e., channel extrapolation [25]. Recently, there are some initial attempts
90 to utilize the channel extrapolation for further reducing the pilot overhead. In [26], the authors
91 proposed a DL-based extrapolation network to extrapolate the complete CSI by exploiting the
92 correlation of the antenna domain, where partial antennas are activated by a selection network. In
93 [27], the authors utilized a neural network structure modified by ordinary differential equations to
94 improve the performance of extrapolation. Besides, the authors of [28] adopted a grouping strategy
95 to reduce the dimension of the estimated channel and designed a CNN-based network to extrapolate
96 the full-dimensional cascaded channel as well as eliminate the grouping interference. However, the
97 above extrapolation schemes only consider the extrapolation process from the known sub-channels,
98 while ignoring how to estimate the sub-channel. Moreover, the hybrid-field channel modeling of
99 RISs has more complex EM wave propagation characteristics, which will hinder the sub-channel
100 acquisition and the following extrapolation of complete channels.

101 How to properly and effectively design the hybrid beamforming and RIS phase according to the
102 CSI is one of the major engineering challenges in the design of RIS-aided communication systems.
103 Recently, some work has been conducted to investigate hybrid beamforming and RIS design problems
104 [29–31]. In [29], simultaneous orthogonal matching pursuit (SOMP)-based hybrid beamforming was
105 proposed for RIS-aided mmWave MIMO systems. In [30], an iteration-based jointly active/passive

106 beamforming algorithm was designed to maximize the sum rate of systems. Furthermore, the DL-
107 based beamforming methods have also been studied in RIS-aided wireless communication systems.
108 In [31], a deep neural network (DNN)-based beamforming approach was developed to jointly optimize
109 the transmit/reflect beamforming vectors for achieving data rate maximization. However, further
110 analysis of the aforementioned schemes with regard to adaptability is necessary, as the current
111 analysis only considers the idealized CSI assumption.

112 1.2 Motivations

113 The current research on RIS has primarily centered on the development of two modes of op-
114 eration, namely, reflective mode [29, 30, 32] and transmissive mode [33–35]. A number of studies
115 have been conducted on RIS-aided communication in reflective mode, which is primarily utilized to
116 address the blind coverage problem. By contrast, the main purpose of transmissive RIS is to improve
117 the spectral efficiency of the networks, as the transmissive mode does not alter the direction of EM
118 waves. Therefore, it is suitable to deploy transmissive RIS in the case that a line-of-sight (LoS) path
119 exists but the propagation attenuation is high, e.g., the case that the outdoor BS serves indoor UEs,
120 to improve the energy of the received signals. In view of this, the transmissive RIS has the potential
121 to provide indoor signal enhancement service.

122 Considering the hybrid-field channel model, the authors of [36] presented a two-stage channel
123 estimation mechanism, where a CNN-based network is designed to estimate channel parameters
124 and the complete channel is reconstructed by channel extrapolation based on geometric relation-
125 ships of channel parameters. However, this parametric-based extrapolation method requires a large
126 number of training labels containing accurate channel parameters. In [37], the authors proposed a
127 sensor-assisted channel estimation and beamforming technique, where a LoS MIMO architecture is
128 considered in the hybrid field. However, the channel estimation in [37] relies heavily on the aware-
129 ness of sensors, which can prove challenging in obtaining accurate CSI. Therefore, similar to [26–28],
130 we adopt a DL-based channel extrapolation method to address the performance limitations of con-
131 ventional channel estimation methods for indoor hybrid-field propagation environments. Besides, in
132 this paper, we consider the LoS MIMO architecture under the assumption of the hybrid-field channel
133 model, where the LoS MIMO architecture can support multi-stream transmission in the pure LoS
134 BS-RIS channel.

135 Most existing works in the field of RIS-aided communication systems have made the assumption
136 that the BS-RIS and RIS-UEs CSIs are perfect [29–32]. However, this assumption is impractical.
137 Therefore, the channel estimation error should be considered when designing these systems. Re-
138 cently, imperfect CSI conditions have been considered in some works [38, 39]. For instance, the
139 authors of [38] utilized a penalty-based alternating algorithm to jointly design active beamform-
140 ing and RIS phase under the presence of imperfect CSI. Similarly, the authors of [39] exploited a
141 gradient projection-based alternating optimization algorithm to jointly design active beamforming,
142 RIS placement, and RIS phase under imperfect CSI. While there are numerous DL-based methods
143 available for RIS-aided communication systems with the perfect CSI, there are only a few DL-based
144 methods that consider imperfect CSI [40]. Therefore, this work aims to provide a DL-based hybrid

145 beamforming and RIS phase design solution that incorporates imperfect CSI in RIS-aided commu-
146 nication systems.

147 1.3 Contributions

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

- 152 • We deploy a transmissive RIS on the window to reduce the penetration loss and thus achieve
153 indoor enhanced communication. In addition, due to the negligible non-LoS (NLoS) component
154 energy in the THz band, the BS-RIS channel is dominated by the LoS path. To achieve multi-
155 stream transmission in the LoS case, we consider a LoS MIMO architecture under hybrid-field
156 channel modeling, where the BS and RIS adopt the same subarray structures, and the subarray
157 spacing is optimized to satisfy the LoS MIMO condition.
- 158 • Since the BS and the RIS are fixed as well as only one LoS path exists, the BS-RIS channel can
159 be considered to be quasi-static and known. In contrast, due to the mobility of the UE, the
160 RIS-UE channel is time-varying. Therefore, we only focus on estimating the RIS-UE channel,
161 which significantly reduces the pilot overhead.
- 162 • To further reduce the pilot overhead for estimating the RIS-UE channel, we propose a DL-
163 based channel extrapolation scheme, where the RIS only activates part of its elements at the
164 channel estimation stage. Unlike the existing extrapolation schemes [26–28] that only focus
165 on the CSI extrapolation process, we design a complete channel extrapolation framework,
166 including the pilot design network, CSI feedback network, sub-channel estimation network,
167 and channel extrapolation network. By adopting the end-to-end (E2E) training strategy, the
168 proposed channel estimation scheme can maintain high reconstruction performance with a
169 few pilot overhead. Specifically, by using the CSI feedback network, the UE-side feeds the
170 quantized pilot information back to the BS, and the BS estimates the sub-sampling RIS-UE
171 channel and then extrapolates the complete RIS-UE channel using the channel extrapolation
172 network. In addition, for the RIS element selection, we discuss the impact of three differ-
173 ent strategies, uniform selection, random selection, and learning-based selection, on the final
174 channel estimation performance.
- 175 • To solve the multi-user interference problem under imperfect CSI, we propose a DL-based
176 hybrid beamforming and RIS phase design scheme, which consists of the analog beamformer
177 design, DL-based RIS phase design network, and knowledge-data dual-driven digital beam-
178 forming network. By maximizing the sum rate with E2E training, the proposed scheme can
179 realize higher performance and better robustness than the existing state-of-the-art methods.

180 *Notations:* In this paper, scalars are denoted as lower-case letters, vectors are denoted as lower-
181 case boldface letters, and matrices are denoted as upper-case boldface letters. The conjugate, trans-

182 pose, conjugate transpose, inversion, and Moore-Penrose inversion operators are denoted as the
 183 superscripts $(\cdot)^*$, $(\cdot)^T$, $(\cdot)^H$, $(\cdot)^{-1}$, and $(\cdot)^\dagger$, respectively. The diagonalization, block diagonalization,
 184 Kronecker product, and Hadamard product are represented by the operators $\text{diag}(\cdot)$, $\text{blkdiag}(\cdot)$, \otimes ,
 185 and \odot , respectively. The Frobenius norm of \mathbf{A} is denoted as $|\mathbf{A}|_F$. The identity matrix with size
 186 $n \times n$ is represented by \mathbf{I}_n , while the column vector of size n with all elements equal to 1 (0) is
 187 represented by $\mathbf{1}_n$ ($\mathbf{0}_n$). The real and imaginary parts of the corresponding argument are denoted
 188 as $\Re\{\cdot\}$ and $\Im\{\cdot\}$, respectively. The m -th row and n -th column element of \mathbf{A} is represented by
 189 $\{\mathbf{A}\}_{m,n}$, and the m -th entry of \mathbf{a} is represented by $\{\mathbf{a}\}_m$. The sub-matrix containing the m -th
 190 to n -th columns of \mathbf{A} is represented by $\mathbf{A}_{[:,m:n]}$. The expectation operator is represented by $\mathbb{E}(\cdot)$,
 191 and the real (complex) Gaussian distribution with mean μ and variance σ^2 is denoted as $\mathcal{N}(\mu, \sigma^2)$
 192 $(\mathcal{CN}(\mu, \sigma^2))$, where the matrix trace operator is represented by $\text{Tr}\{\cdot\}$.

193 Materials and Methods

194 2 System Model

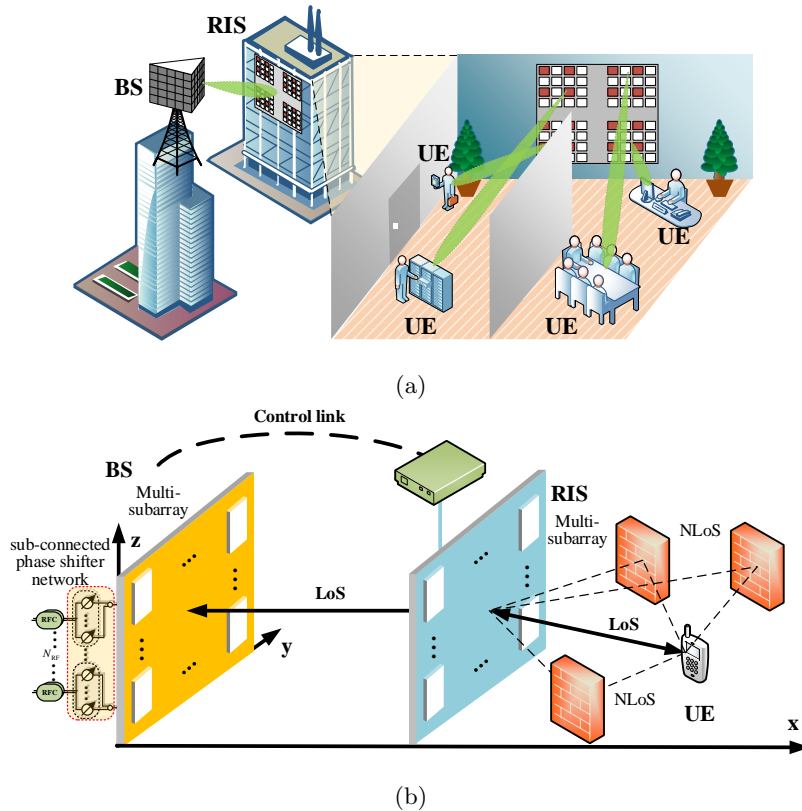


Figure 1: Schematic diagram of a RIS-aided THz massive MIMO system: (a) multiple indoor UEs are served by the BS with the help of a transmissive RIS deployed on the window, and (b) hardware architectures at the BS, RIS, and UEs.

2.1 System Description

As shown in Figure 1, we consider a downlink RIS-aided MIMO orthogonal frequency division multiplexing (OFDM) transmission system 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 enhance indoor coverage. Let the BS (RIS) have $M^B = M_y^B \times M_z^B$ ($M^R = M_y^R \times M_z^R$) uniformly spaced subarrays, where M_y^B (M_y^R) and M_z^B (M_z^R) are the numbers of BS (RIS)-side subarrays along the horizontal and vertical directions, respectively. Each subarray of the BS (RIS) is a uniform planar array (UPA) with $N_{\text{sub}}^B = N_y^B \times N_z^B$ ($N_{\text{sub}}^R = N_y^R \times N_z^R$) isotropically radiating elements, where N_y^B (N_y^R) and N_z^B (N_z^R) are the numbers of BS (RIS)-side subarray antennas along the horizontal and vertical directions, respectively. Therefore, the complete antenna dimension of the BS is $N^B = M^B N_{\text{sub}}^B$, and the element dimension of the RIS is $N^R = M^R N_{\text{sub}}^R$. To simplify the analysis, we assume that the normals of central elements of both the BS and RIS are coaxial, i.e., meeting the parallel symmetric array arrangements, with a distance of D , as illustrated in Figure 1(b).

In this paper, we consider a BS-side sub-connected hybrid analog-digital array architecture. This architecture consists of M^B RF chains, which are capable of supporting $U \leq M^B$ data streams. Each of these RF chains is connected to a subarray through N_{sub}^B phase shifters. Furthermore, We set the number of subcarriers to K and sampling frequency (i.e., bandwidth) to f_s . The carrier frequency is f_c , which corresponds to central wavelength λ .

2.2 Channel Model

2.2.1 BS-RIS Channel Model

Due to the negligible non-LoS (NLoS) component energy in the THz band, we only consider the LoS path in the analysis of the BS-RIS channel. By utilizing the spherical wave propagation characteristic, we construct the LoS MIMO link between the BS and RIS with only one single LoS path, but it can support intra-path multiplexing for multi-stream transmission [41]. The inter-antenna spacing in each subarray is $d = \lambda/2$. In order to satisfy the LoS MIMO characteristic, the BS subarray spacing d_{sy}^B and d_{sz}^B are set to the following optimal LoS MIMO spacing

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

i.e., d_{sy}^B and d_{sz}^B should satisfy the condition $\lambda \ll d_{sy}^B, d_{sz}^B \ll D$. The detailed explanation of Equation (1) can be found in [41, 42]. The RIS subarray spacing d_{sy}^R and d_{sz}^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 [42]. We have the following proposition from [43].

Proposition 1 *Let the transceiver arrays be placed with a separation distance of D and be working at a carrier wavelength λ ($\lambda \ll D$). If the inter-antenna spacing and carrier wavelength λ are in*

229 *the same order of magnitude, the planner wave model can be applied. Otherwise, the spherical wave*
 230 *model should be exploited.*

According to Proposition 1, the subarray response vectors $\mathbf{a}(\theta, \phi, f_k) \in \mathbb{C}^{N_H N_V \times 1}$ can be approximated by a planner wave model:

$$\begin{aligned} \mathbf{a}(\theta, \phi, f_k) &= \mathbf{a}_h(\theta, \phi, f_k) \otimes \mathbf{a}_v(\phi, f_k) \\ &= [1, \dots, e^{-j2\pi \frac{f_k}{c} d(n_h \sin \theta \cos \phi + n_v \sin \phi)}, \dots, e^{-j2\pi \frac{f_k}{c} d((N_H-1) \sin \theta \cos \phi + (N_V-1) \sin \phi)}]^T, \end{aligned} \quad (2)$$

231 where $f_k = f_c - \frac{f_s}{2} + \frac{kf_s}{K}$, $1 \leq k \leq K$, is the k -th subcarrier frequency, c is the speed of light,
 232 $0 \leq n_h \leq (N_H - 1)$, $0 \leq n_v \leq (N_V - 1)$, N_H and N_V are the numbers of horizontal and vertical
 233 antennas, respectively, while θ and ϕ are the horizontal and vertical angles of the departure or arrival
 234 (AoD or AoA) of the path, respectively.

235 Since $d_{sy}^B, d_{sz}^B, d_{sy}^R, d_{sz}^R \ll D$, the same path's direction difference in different subarrays is negligible.
 236 Therefore, all subarrays on either the BS or RIS-side can be assumed to share the identical array
 237 response vectors. However, as subarrays are widely spaced, the relative phase differences among
 238 subarrays are non-negligible [43]. Motivated by the above analysis, the downlink spatial-frequency
 239 BS-RIS channel $\mathbf{G}[k] \in \mathbb{C}^{N^R \times N^B}$ on the k -th subcarrier can be modeled as

$$\mathbf{G}[k] = \alpha[k] G_T \tilde{\mathbf{G}}[k] \otimes [\mathbf{a}_R(\theta_{R,A}, \phi_{R,A}, f_k) \mathbf{a}_B^H(\theta_B, \phi_B, f_k)], \quad (3)$$

where $\alpha[k]$ is the channel attenuation coefficient on the k -th subcarrier, (θ_B, ϕ_B) and $(\theta_{R,A}, \phi_{R,A})$ are
 AoD and AoA of the LoS path, respectively. Without loss of generality, we assume the LoS angles are
 fixed and known in advance since the BS and RIS are fixed. In (3), the entries of $\tilde{\mathbf{G}}[k] \in \mathbb{C}^{M^R \times M^B}$
 are defined according to the spherical wave model as

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

240 where $D(m_r, m_b)$ represents the distance between the m_r -th RIS-side subarray and the m_b -th BS-side
 241 subarray. Furthermore, the subarray response vectors $\mathbf{a}_R(\theta_{R,A}, \phi_{R,A}, f_k) \in \mathbb{C}^{N_{\text{sub}}^R \times 1}$ and $\mathbf{a}_B(\theta_B, \phi_B, f_k)$
 242 $\in \mathbb{C}^{N_{\text{sub}}^B \times 1}$ are defined in Equation (2). The constant coefficient G_T represents the antenna gain at
 243 the BS, which is different from the array gain generated by beamforming [44]. The only unknown
 244 parameter in Equation (3) is the channel coefficient $\alpha[k]$, which can be obtained by placing a power
 245 detector at the RIS side. Therefore, it is reasonable to assume that the quasi-static BS-RIS channel
 246 is 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 envi-
 ronments [45]. The indoor RIS-UE channel model consists of one LoS path and L_p NLoS paths,
 where their three-dimensional (3D) distances are represented as d_0 and d_l , for $1 \leq l \leq L_p$, respec-
 tively [46]. The total EM wave propagation loss mainly consists of two parts: the free space path
 loss $\beta_{\text{spr}}(f_k, d_l) = \frac{c}{4\pi f_k d_l}$ and the molecular absorption loss $\beta_{\text{abs}}(f_k, d_l) = e^{-\frac{1}{2}\kappa(f_k)d_l}$, where $\kappa(f_k)$

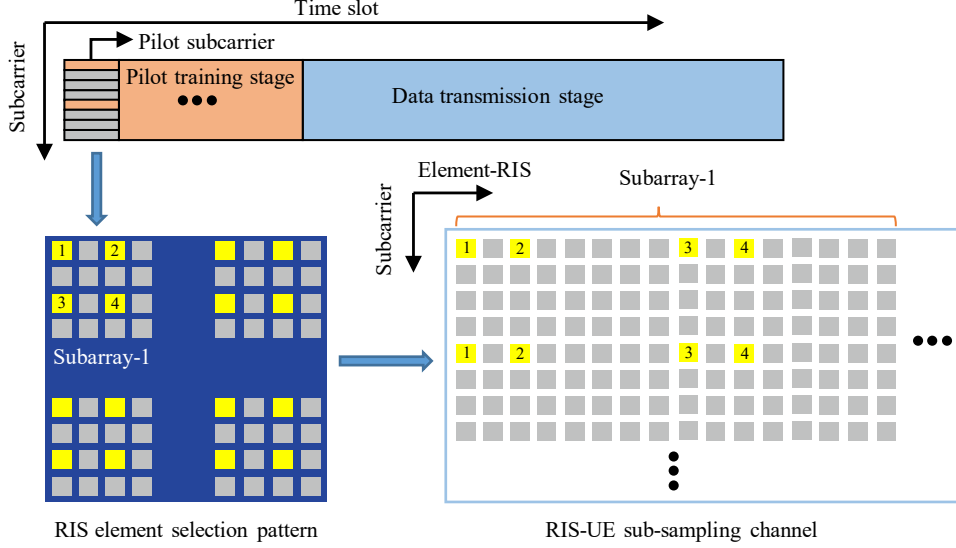


Figure 2: Block diagram of the frame structure, RIS element selection pattern, and RIS-UE sub-sampling channel, where the selected parts are marked in yellow blocks and the number in the yellow block indicates the index of the selected element.

denotes the frequency-dependent absorption coefficient [47]. Hence, the spatial-frequency channel $\mathbf{h}[k] \in \mathbb{C}^{1 \times N^R}$ for the RIS-UE link is

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

248 where $\beta[k] = \beta_{\text{spr}}(f_k, d_0) \beta_{\text{abs}}(f_k, d_0)$ and $\beta_l[k] = \beta_{\text{spr}}(f_k, d_l) \beta_{\text{abs}}(f_k, d_l) \beta_{\text{RC}}$ are the channel attenua-
 249 tion coefficients of the LoS path and the l -th NLoS path, respectively, $(\theta_{\text{R,D}}^{\text{LoS}}, \phi_{\text{R,D}}^{\text{LoS}})$ and $(\theta_{\text{R,D}}^l, \phi_{\text{R,D}}^l)$
 250 are the LoS AoD and the NLoS AoD of the l -th NLoS path, respectively. Additionally, the reflection
 251 coefficient β_{RC} is a Gaussian random variable, i.e., $10 \log \beta_{\text{RC}}[\text{dB}] \sim \min\{\mathcal{N}(\mu_{\text{R}}, \sigma_{\text{R}}^2), 0\}$. The entries
 252 of $\tilde{\mathbf{h}}_{\text{LoS}}[k] \in \mathbb{C}^{1 \times M^R}$ are given as $\{\tilde{\mathbf{h}}_{\text{LoS}}[k]\}_{m_r} = e^{-j2\pi f_k \cdot \frac{d^{(m_r)}}{c}}$, where $d^{(m_r)}$ denotes the 3D distance
 253 between the UE and the m_r -th RIS-side subarray. $\tilde{\mathbf{h}}^l[k]$ has a similar notation and assumptions.

254 3 Problem Formulation and Proposed Channel Estimation 255 Solution

256 3.1 Problem Formulation of Channel Estimation

257 In this subsection, the downlink channel estimation problem is formulated based on the con-
 258 sidered RIS-aided THz massive MIMO communication system over hybrid-field channels. As shown
 259 in Figure 2, we consider the two-stage frame structure consisting of the pilot training and data
 260 transmission stages. At the pilot training stage, the BS transmits M pilot OFDM symbols (i.e., M
 261 time slots) dedicated to channel estimation. The m -th received signal at the UE-side¹ on the k -th

¹Note that since each UE can perform channel estimation independently, UE subscripts are omitted.

262 subcarrier is represented by

$$y_m[k] = \sqrt{P_T} \mathbf{h}[k] \mathbf{\Phi}_m \mathbf{G}[k] \mathbf{F}_{RF} \mathbf{F}_{BB}[k] \mathbf{s}_m[k] + n_m[k], \quad (6)$$

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

$$\mathbf{v}_{m_r, m} = \mathbf{o}_{m_r, m} \odot \tilde{\mathbf{v}}_{m_r, m} = [\cdots, \eta_{n_{\text{sub}, m_r, m}}^r, \cdots] \odot [\cdots, e^{j\phi_{n_{\text{sub}, m_r, m}}^r}, \cdots], \quad (7)$$

263 where $\mathbf{o}_{m_r, m} \in \mathbb{C}^{1 \times N_{\text{sub}}^R}$ represents the amplitude control vector, $\tilde{\mathbf{v}}_{m_r, m} \in \mathbb{C}^{1 \times N_{\text{sub}}^R}$ represents the
 264 phase control vector, and $1 \leq n_{\text{sub}}^r \leq N_{\text{sub}}^R$, while $\eta_{n_{\text{sub}, m_r, m}}^r \in [0, 1]$ and $\phi_{n_{\text{sub}, m_r, m}}^r \in [0, 2\pi]$
 265 are the amplitude/phase control coefficient, respectively. $\eta_{n_{\text{sub}, m_r, m}}^r$ can control the switch of the
 266 refraction function for each RIS element. The entire RIS elements can be expressed as $\mathbf{v}_m =$
 267 $\mathbf{o}_m \odot \tilde{\mathbf{v}}_m = [\mathbf{v}_{1, m}, \cdots, \mathbf{v}_{m_r, m}, \cdots, \mathbf{v}_{M^R, m}]^T \in \mathbb{C}^{N^R \times 1}$, where $\mathbf{o}_m = [\mathbf{o}_{1, m}, \cdots, \mathbf{o}_{M^R, m}]^T \in \mathbb{C}^{N^R \times 1}$
 268 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
 269 $\mathbf{\Phi}_m = \text{diag}(\mathbf{v}_m) = \mathbf{O}_m \odot \tilde{\mathbf{V}}_m \in \mathbb{C}^{N^R \times N^R}$, where $\mathbf{O}_m = \text{diag}(\mathbf{o}_m) \in \mathbb{C}^{N^R \times N^R}$ is the RIS selection
 270 matrix and $\tilde{\mathbf{V}}_m = \text{diag}(\tilde{\mathbf{v}}_m) \in \mathbb{C}^{N^R \times N^R}$ is the RIS phase matrix.

271 $\mathbf{F}_{RF} \in \mathbb{C}^{N^B \times M^B}$ and $\mathbf{F}_{BB}[k] \in \mathbb{C}^{M^B \times U}$ are respectively analog and digital beamforming matrices
 272 that are used at the BS to provide array gain and eliminate the multi-stream interference. Since the
 273 sub-connected architecture, the analog beamformer implemented by phase shifters is written as

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

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

$$\mathbf{f}_{m_b} = \mathbf{a}_B(\theta_B, \phi_B, f_k), 1 \leq m_b \leq M^B, \quad (9)$$

277 where k can be set to $K/2$ for alleviating the beam squint problem induced by the large bandwidth
 278 [48]. The digital beamformer $\mathbf{F}_{BB}[k]$ is designed according to the zero-forcing (ZF) precoding in
 279 order to eliminate the multi-stream interference between the BS and RIS subarrays, i.e.,

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

280 where $\tilde{\mathbf{G}}_{\text{eq}}[k] = [\alpha[k] G_T \tilde{\mathbf{G}}[k] \otimes \mathbf{a}_B^H(\theta_B, \phi_B, f_k)] \mathbf{F}_{RF} \in \mathbb{C}^{M^R \times M^B}$ is the equivalent BS-RIS chan-
 281 nel obtained from the perspective of the first element of different subarrays at the RIS, and $\zeta =$
 282 $\sqrt{M^B / \text{Tr}\{\tilde{\mathbf{G}}_{\text{eq}}^\dagger[k] (\tilde{\mathbf{G}}_{\text{eq}}^\dagger[k])^H\}}$ is a constant to meet the total transmit power constraint after beam-
 283 forming. In this way, the multi-stream interference between the BS and the RIS subarrays can be
 284 eliminated, i.e., $\mathbf{G}_{\text{eq}}[k] = \mathbf{G}[k] \mathbf{F}_{RF} \mathbf{F}_{BB}[k] \in \mathbb{C}^{N^R \times U}$, $\forall k$, is a block diagonal constant matrix.

285 Therefore, the equivalent pilot signal $\mathbf{p}_m \in \mathbb{C}^{N^R \times 1}$ can be written as

$$\mathbf{p}_m = \underbrace{\left[\mathbf{O}_m \odot \tilde{\mathbf{V}}_m \right]}_{\Phi_m} \underbrace{\left[\mathbf{G}[k] \mathbf{F}_{\text{RF}} \mathbf{F}_{\text{BB}}[k] \right]}_{\mathbf{G}_{\text{eq}}[k]} \mathbf{s}_m[k], \quad (11)$$

286 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 in Equation (10) for $\mathbf{G}[k]$. Under the assumption that the BS and RIS
 288 meet the parallel symmetric array arrangements, $\mathbf{G}_{\text{eq}}[k]$ is defined by $\sqrt{N^B} \alpha[k] G_{\text{T}} \text{blkdiag}(\mathbf{1}_{N_{\text{sub}}^{\text{R}}}, \dots,$
 289 $\mathbf{1}_{N_{\text{sub}}^{\text{R}}}, \dots, \mathbf{1}_{N_{\text{sub}}^{\text{R}}})$. Thus, the effective pilot signals can be further expressed as the RIS element vector
 290 given by $\mathbf{p}_m = \sqrt{N^B} \alpha[k] G_{\text{T}} \mathbf{v}_m \approx \sqrt{N^B} \alpha G_{\text{T}} \mathbf{v}_m = A_{\text{T}} \mathbf{v}_m$, where the approximation $\alpha[k] \approx \alpha$, $\forall k$, is
 291 further applied and $A_{\text{T}} = \sqrt{N^B} \alpha G_{\text{T}}$ represents the total attenuation from the BS to the RIS.

292 By collecting continuous measurements of M time slots, the aggregate received signal vector
 293 $\mathbf{y}[k] = [y_1[k], \dots, y_M[k]] \in \mathbb{C}^{1 \times M}$ is written as

$$\mathbf{y}[k] = \sqrt{P_{\text{T}}} \mathbf{h}[k] \mathbf{P} + \mathbf{n}[k], \quad (12)$$

294 where $\mathbf{P} = [\mathbf{p}_1, \dots, \mathbf{p}_M] = A_{\text{T}} \mathbf{V} = A_{\text{T}} [\mathbf{v}_1, \dots, \mathbf{v}_M] \in \mathbb{C}^{N^R \times M}$, and $\mathbf{n}[k] = [n_1[k], \dots, n_M[k]] \in$
 295 $\mathbb{C}^{1 \times M}$. Thus, the received signal matrix $\mathbf{Y} = [\mathbf{y}^{\text{T}}[1], \dots, \mathbf{y}^{\text{T}}[K]]^{\text{T}} \in \mathbb{C}^{K \times M}$ can be written as

$$\mathbf{Y} = \sqrt{P_{\text{T}}} \mathbf{H} \mathbf{P} + \mathbf{N}, \quad (13)$$

296 where $\mathbf{H} = [\mathbf{h}^{\text{T}}[1], \dots, \mathbf{h}^{\text{T}}[K]]^{\text{T}} \in \mathbb{C}^{K \times N^R}$ represents the downlink spatial-frequency domain RIS-
 297 UE channel matrix, and $\mathbf{N} = [\mathbf{n}^{\text{T}}[1], \dots, \mathbf{n}^{\text{T}}[K]]^{\text{T}} \in \mathbb{C}^{K \times M}$.

298 3.2 Deep Learning Based Spatial-Frequency Domain Channel Extrapolation

299

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

$$\mathbf{Y}_s = \sqrt{P_{\text{T}}} \mathbf{H}_s \mathbf{P}_s + \mathbf{N}_s, \quad (14)$$

where $\mathbf{H}_s \in \mathbb{C}^{K_s \times N_s^{\text{R}}}$ is the sub-sampling of the spatial-frequency channel, $\mathbf{P}_s \in \mathbb{C}^{N_s^{\text{R}} \times M}$ is the
 corresponding equivalent pilot signal, and \mathbf{N}_s is the noise. Our goal is to recover complete channel
 $\hat{\mathbf{H}} \in \mathbb{C}^{K \times N^{\text{R}}}$ based on limited received pilot signals \mathbf{Y}_s , i.e., extrapolating the rest unknown channels
 from the acquired partial channels. Based on the non-linear function fitting capability of DL, We can

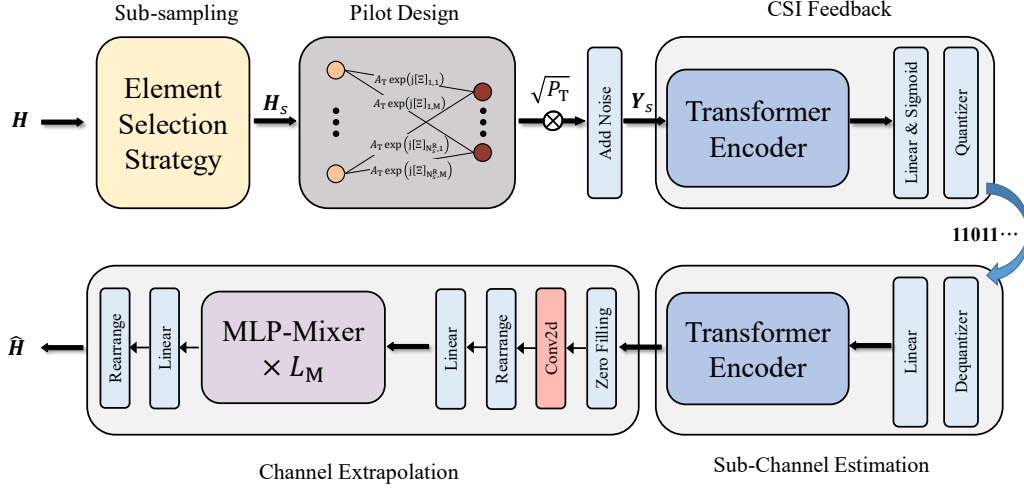


Figure 3: The overall block diagram of the proposed DL-based spatial-frequency domain channel extrapolation scheme.

learn a mapping to represent proximity correlations between different spatial/frequency locations of channels. 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 complete process of the proposed scheme can be represented as

$$\hat{\mathbf{H}} = f_{\text{SFDE}}(f_{\text{SCE}}(f_{\text{CSIFd}}(\sqrt{P_T} f_{\text{ESS}}(\mathbf{H}) \mathbf{P}_s + \mathbf{N}_s))), \quad (15)$$

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

3.2.1 Element Selection Strategy

With only N_s^R activated RIS elements, from Equation (7), the RIS element selection vector $\mathbf{o} = \mathbf{o}_m = [\mathbf{o}_{1,m}, \dots, \mathbf{o}_{m_r,m}, \dots, \mathbf{o}_{M^R,m}]^T \in \{0, 1\}^{N_s^R \times 1}$ is an N_s^R -hot vector with N_s^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 $f_{\text{ESS}}(\cdot)$ can be expressed as

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

where $\mathbf{S} = \boldsymbol{\kappa} \otimes \mathbf{o}^T \in \{0, 1\}^{K \times N_s^R}$ 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 different RIS element selection vectors can affect the extrapolation performance. Thus, we consider the following three element selection strategies.

325 1) Uniform Selection Strategy: Since each subarray in the RIS is a UPA, its element compression
 326 ratio is expressed as $\rho = \rho_y \times \rho_z$, where ρ_y and ρ_z are the compression ratios along the azimuth
 327 and elevation directions, respectively. To fairly sound the channel and ensure a balanced estimation
 328 performance along two directions, ρ_y and ρ_z are expected to be as close as possible. However, $\rho_y = \rho_z$
 329 cannot always be guaranteed under all system parameter configurations. In cases where $\rho_y \neq \rho_z$, it
 330 is desirable to allocate more activated elements along the azimuth (y -axis) direction rather than the
 331 z -axis direction (i.e., $\rho_y \leq \rho_z$). This strategic choice aligns with the consideration of indoor UEs,
 332 which are more likely to be distributed across a wide azimuth range, as opposed to the elevation
 333 range, as indoor UEs are typically stationary in the vertical dimension. In light of this, the y - z
 334 compression ratio allocation can be solved from the following optimization problem

$$\begin{aligned} \min_{\{\rho_y, \rho_z\}} \quad & |\rho_z - \rho_y|, \\ \text{s.t.} \quad & \rho_y \times \rho_z = \rho, \\ & 1 \leq \rho_y \leq \rho_z. \end{aligned} \quad (17)$$

335 Some allocation examples as $\rho(2, 4, 8, 16) = \rho_y(1, 2, 2, 4) \times \rho_z(2, 2, 4, 4)$. Given ρ_y and ρ_z , the active
 336 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

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

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

341 2) Random Selection Strategy: It randomly selects N_s^{R} elements from the RIS as the random
 342 pattern and generates the active element index vector $\boldsymbol{\xi}$. If the element compression ratio ρ is not
 343 large, then the aperture of a random pattern is usually comparable to that of the RIS.

344 3) Learning-Based Selection Strategy: In addition to the above two fixed selection strategies,
 345 the learning-based element selection strategy has also been widely studied. In [26], a differentiable
 346 selection network is proposed to learn the element selection vector \mathbf{o} . The input of this network is
 347 a random initialization vector. By utilizing several fully-connected layers and the softmax function,
 348 a probability vector $\mathbf{g} = [g_1, g_2, \dots, g_{N^{\text{R}}}]^{\text{T}} \in \mathbb{C}^{N^{\text{R}} \times 1}$ is generated, where g_i denotes the probability
 349 that the i -th element is selected. Thus, the active element index vector $\boldsymbol{\xi}$ can be defined as

$$\boldsymbol{\xi} = \arg \text{top}_{N_s^{\text{R}}} \{\mathbf{g}\}, \quad (19)$$

350 where $\arg \text{top}_{N_s^{\text{R}}} \{\cdot\}$ is a function that finds the element index set of the first N_s^{R} largest selection
 351 probabilities. The details of the selection network can be found in [26].

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 $\mathbf{P}_s = A_T \mathbf{V}_s$, where $\mathbf{V}_s \in \mathbb{C}^{N_s^R \times M}$ denotes the RIS phase matrix of selected 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 given by

$$\mathbf{P}_s = A_T \exp(j\mathbf{\Xi}) = A_T (\cos(\mathbf{\Xi}) + j \sin(\mathbf{\Xi})), \quad (20)$$

where $\mathbf{\Xi} \in \mathbb{R}^{N_s^R \times M}$ is the phase control matrix of selected RIS elements. Since most DL frameworks, such as Tensorflow and Pytorch, have limited support for complex-valued operations, it becomes challenging to train the complex-valued pilot matrix \mathbf{P}_s directly. To circumvent this issue, we adopt the real-valued RIS phase control matrix $\mathbf{\Xi}$ whose entries take values in $[0, 2\pi)$ as 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 trainable parameters of the PDN, i.e., $\mathbf{\Xi}$, are learned at the DL training stage.

3.2.3 CSI Feedback

Recently, DL-based solutions, such as CsiNet [49], have achieved good performance for CSI feedback. Furthermore, an emerging CSI feedback architecture based on the transformer [50] has been demonstrated to further reduce the feedback overhead and obtain more efficient compression performance than the CsiNet framework [51]. Therefore, we utilize the transformer as the backbone of CSI feedback network $f_{\text{CsiFd}}(\cdot)$. The original transformer is divided into an encoder and a decoder. However, since we are dealing with the CSI without time-sequential information, there is no causality

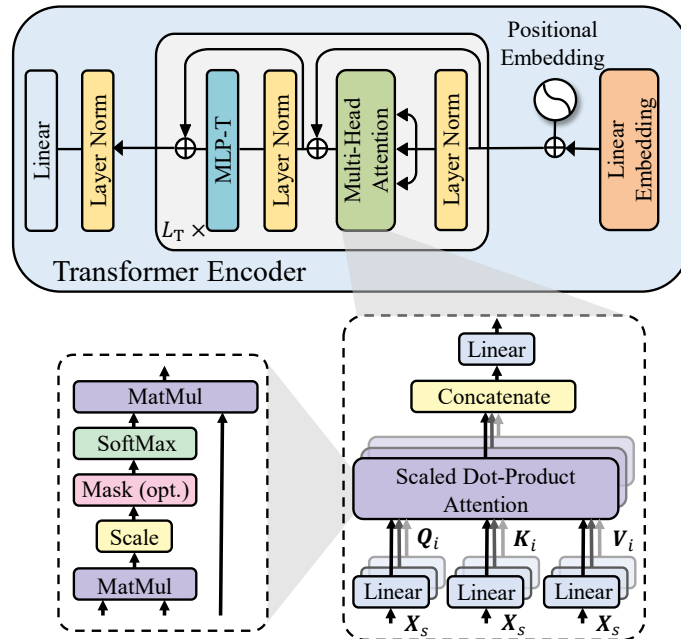


Figure 4: The structure of the transformer encoder.

372 constraint. Thus, we only exploit the encoder module of the transformer which obtains output in
 373 parallel. Since real-valued operations are more effective and the transformer can only extract the
 374 correlation between sequences, we convert the received pilot signal into a real-valued two-dimensional
 375 (2D) sequence $\bar{\mathbf{Y}}_s \in \mathbb{R}^{K_s \times 2M}$, which can be expressed as

$$\bar{\mathbf{Y}}_s = [\Re\{\mathbf{Y}_s\}, \Im\{\mathbf{Y}_s\}], \quad (21)$$

376 where the number of subcarriers K_s represents the length of the input sequence.

The schematic diagram of the transformer encoder is shown in Figure 4. Through the fully-connected linear embedding layer, input sequence $\bar{\mathbf{Y}}_s$ can be converted into $\mathbf{X}_s \in \mathbb{R}^{K_s \times d_T}$, which merges the relative position information of the sub-carriers using the positional embedding layer. Then, multiple encoder layers are utilized to extract correlations between 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 plays a key role in 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, $\{\mathbf{Q}_i, \mathbf{K}_i, \mathbf{V}_i\} \in \mathbb{R}^{K_s \times d_m}$, $1 \leq i \leq h$, where $d_m = d_T/h$ and h is the number of heads. Then, each value head $i \in \mathbb{R}^{K_s \times d_m}$, $1 \leq i \leq h$, is outputted by performing the scaled dot-product attention simultaneously, where the weights on values can be obtained from a softmax function, which is given by

$$\text{head}_i = \text{softmax} \left(\frac{\mathbf{Q}_i \mathbf{K}_i^T}{\sqrt{d_m}} \right) \mathbf{V}_i, \quad 1 \leq i \leq h. \quad (22)$$

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

$$\text{MultiHead}(\mathbf{X}_s) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) \mathbf{W}^O. \quad (23)$$

377 After the transformer encoder, a linear layer followed by a sigmoid function is used to generate a
 378 compressed codeword, which is then transformed into B bits as the feedback information through
 379 a uniform scalar quantization layer. The above feedback process generates the binary vector $\mathbf{q} \in$
 380 $\{0, 1\}^B$ as

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

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

382 3.2.4 Sub-Channel Estimation

383 When the BS receives the feedback bits, the sub-channel estimation network is used to recon-
 384 struct the sub-sampling of the complete spatial-frequency channel. Similar to Subsection 3.2.3, we
 385 also consider the transformer encoder as the backbone of this part. As shown in Figure 3, received
 386 CSI feedback bits are initially processed by a dequantization layer, which conducts the inverse op-

387 eration of the quantizer and outputs a real-valued vector. Then, the initial coarse channel estimate
 388 is obtained by a linear layer. Finally, the transformer encoder extracts the spatial-frequency corre-
 389 lation of the channel and further improves the channel estimation performance. The sub-channel
 390 estimation process can be represented by

$$\bar{\mathbf{H}}_s = \left[\Re\{\hat{\mathbf{H}}_s\}, \Im\{\hat{\mathbf{H}}_s\} \right] = f_{\text{SCE}}(\mathbf{q}; \mathcal{W}_S), \quad (25)$$

391 where $\hat{\mathbf{H}}_s \in \mathbb{C}^{K_s \times N_s^{\text{R}}}$ is the estimated sub-sampling channel, $\bar{\mathbf{H}}_s \in \mathbb{R}^{K_s \times N_s^{\text{R}} \times 2}$ is a real-valued 3D
 392 matrix, and \mathcal{W}_S is the trained parameter set of the sub-channel estimation network.

393 3.2.5 Spatial-Frequency Domain Channel Extrapolation

394 First, the initial input $\tilde{\mathbf{H}} \in \mathbb{R}^{K \times N^{\text{R}} \times 2}$ to the channel extrapolation network is constructed from
 395 the estimated sub-sampling channel $\bar{\mathbf{H}}_s \in \mathbb{R}^{K_s \times N_s^{\text{R}} \times 2}$ with the known RIS spatial-frequency selection
 396 pattern \mathbf{S} . Specifically, we copy the entries of $\bar{\mathbf{H}}_s$ to the corresponding positions in $\tilde{\mathbf{H}}$ and fill the
 397 other elements of $\tilde{\mathbf{H}}$ with zeros according to the known RIS spatial-frequency selection pattern \mathbf{S} .
 398 This initial operation is represented by

$$\tilde{\mathbf{H}} = f_{\text{zfi}}(\bar{\mathbf{H}}_s; \mathbf{S}). \quad (26)$$

399 The non-zero rows/columns in $\tilde{\mathbf{H}}$ are consistent with $\bar{\mathbf{H}}_s$ and their locations are the same as the
 400 locations of elements ‘1’ in \mathbf{S} . The neighborhood information in the receptive field is then extracted
 401 using a convolutional layer for initial interpolation. To guarantee that the output dimensions from
 402 the convolution layer remain unchanged, we employ zero padding, i.e., adding additional zeros around
 403 the input feature map.

Subsequently, we consider a competitive yet conceptually and technically simple architecture,
 called MLP-Mixer [52], 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
 $\tilde{\mathbf{H}} \in \mathbb{R}^{K \times N^{\text{R}} \times 2}$ is rearranged as a series of flattened 2D patches $\mathbf{X}_p \in \mathbb{R}^{N_p \times (2L^2)}$, where (K, N^{R})
 represents the size of the original input, (L, L) represents each patch’s length and width, as well
 as $N_p = KN^{\text{R}}/L^2$ represents the number of patches. Then, all the patches are linearly projected

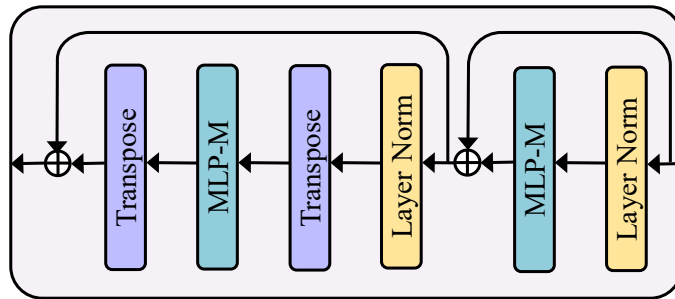


Figure 5: The structure of the mixer layer.

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 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 $\tilde{\mathbf{X}}^T$, maps $\mathbb{R}^{d_M} \mapsto \mathbb{R}^{2d_M} \mapsto \mathbb{R}^{d_M}$, and is shared across all the rows. Each MLP block contains two fully-connected layers and a nonlinear activation function. The mapping of the t -th mixer layer can be expressed as

$$\begin{aligned} \mathbf{U} &= \tilde{\mathbf{X}}_t + \mathbf{W}_{t,2} f_\sigma(\mathbf{W}_{t,1} \text{LayerNorm}(\tilde{\mathbf{X}}_t)), \\ \tilde{\mathbf{X}}_{t+1} &= \mathbf{U} + (\mathbf{W}_{t,4} f_\sigma(\mathbf{W}_{t,3} \text{LayerNorm}(\mathbf{U})^T))^T, \end{aligned} \quad (27)$$

where $\tilde{\mathbf{X}}_t$ denotes the input matrix to the t -th mixer layer, $\mathbf{W}_{t,i}$, $1 \leq i \leq 4$, are the parameter matrices of the fully-connected layers in the t -th mixer layer for $1 \leq t \leq L_M$, and L_M is the number of mixer layers, while f_σ denotes an 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 is represented by

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

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

3.2.6 Training Strategy

The off-line training dataset is represented as \mathcal{H} , which comprises of $|\mathcal{H}| = N_{\text{set}}$ samples. Each sample in \mathcal{H} is an input-label pair represented by (\mathbf{H}, \mathbf{H}) , where \mathbf{H} serves as both the extrapolation target and the input for the 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 an E2E training for the pilot design network, CSI feedback network, sub-channel estimation network, and channel extrapolation network. Thus, the loss function can be represented by minimizing the normalized mean square error (NMSE) between the output $\hat{\mathbf{H}}$ and the target \mathbf{H} , i.e.,

$$\mathcal{L}_c = \frac{1}{B_e} \sum_{i=1}^{B_e} \frac{\|\mathbf{H} - \hat{\mathbf{H}}\|_F^2}{\|\mathbf{H}\|_F^2}, \quad (29)$$

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

425 are optimized jointly, i.e., the loss function can be represented by

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

426 where $0 < \gamma \leq 1$ represents the weight used to balance channel extrapolation and ESS with $\gamma = 1$
 427 denoting that the non-learning based $f_{\text{ESS}}(\cdot)$ is selected, and \mathcal{L}_{ESS} is the loss function of the learning-
 428 based ESS. The details of \mathcal{L}_{ESS} are available in [26].

429 4 Proposed Beamforming Solution

430 4.1 Problem Formulation of RIS-aided Multi-User Beamforming

The BS can simultaneously support U UEs with the aid of RIS at the data transmission stage, since the LoS MIMO architecture 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 represented by

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

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

$$\text{SINR}[u, k] = \frac{P_T |\mathbf{h}[u, k] \mathbf{\Phi} \mathbf{G}[k] \mathbf{F}_{\text{RF}} \mathbf{f}_{\text{BB}}[u, k]|^2}{P_T \sum_{i=1, i \neq u}^U |\mathbf{h}[u, k] \mathbf{\Phi} \mathbf{G}[k] \mathbf{F}_{\text{RF}} \mathbf{f}_{\text{BB}}[i, k]|^2 + \sigma_n^2}. \quad (32)$$

435 Therefore, the sum rate R of total UEs is represented by

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

436 By utilizing the estimated RIS-UE channel at the pilot training stage, the BS can design the
 437 hybrid beamformer $\{\mathbf{F}_{\text{RF}}, \mathbf{F}_{\text{BB}}[k], \forall k\}$ and the RIS refraction phase matrix $\mathbf{\Phi}$ to maximize the sum

438 rate R , where $\mathbf{F}_{\text{BB}}[k] = [\mathbf{f}_{\text{BB}}[1, k], \dots, \mathbf{f}_{\text{BB}}[U, k]]$. This design process is illustrated as

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

439 where $\hat{\mathbf{H}}[u]$ is the estimated spatial-frequency RIS-UE channel of the u -th UE, and $\mathcal{F}(\cdot)$ represents a
440 function that maps the estimated RIS-UE channels onto the hybrid beamformer $\{\mathbf{F}_{\text{RF}}, \mathbf{F}_{\text{BB}}[k], \forall k\}$
441 and the RIS refraction phase matrix Φ .

442 4.2 Deep Learning Based Hybrid Beamforming and RIS Phase Design

443 To solve this optimization Equation (34), some alternating iterative algorithms [29, 30, 32] have
444 been proposed to obtain the analog beamformer, digital beamformer, and RIS phase, respectively.
445 Unfortunately, all the aforementioned approaches are based on the idealized case that the CSI
446 is known accurately. However, perfect CSI is usually unavailable, especially for indoor channel
447 cases where the channel characteristics are complex due to rich scatterers. By using the non-linear
448 function fitting capability of DL, we can learn the complicated and unknown mapping from the
449 estimated channels to the hybrid beamformers and RIS refraction phase. Thus, we propose a DL-
450 based hybrid beamforming and RIS phase design scheme, which consists of analog beamformer
451 design, DL-based RIS refraction phase design, and knowledge-data dual-driven digital beamformer
452 design. The diagram depicting the design of the proposed scheme is presented in Figure 6.

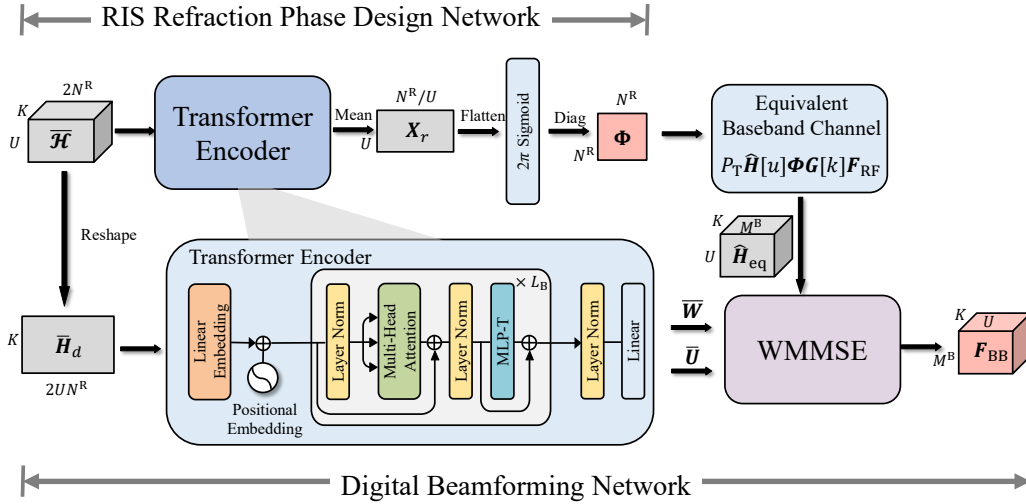


Figure 6: The overall structure of the proposed DL-based hybrid beamforming and RIS refraction phase design scheme.

4.2.1 Analog Beamformer Design

The integration of the active/passive beamforming at the BS and RIS is a non-convex optimization problem, which poses significant difficulties in finding a global optimum solution. Hence, we separately design the analog beamforming and the passive beamforming. Specifically, both BS analog beamforming and RIS passive beamforming are designed to focus energy for improving the received SINR of UEs. However, since the sub-connected structure in the LoS MIMO architecture, the interference among beams from the BS subarrays to the RIS subarrays cannot be eliminated. Fortunately, this part of interference can be removed by appropriately designing the digital beamforming. Therefore, when designing the analog beamforming on the BS-side, it is sufficient to assume that the transmit energy is focused on the RIS.

Since the BS-RIS channel with only the LoS path is quasi-static and known, we can utilize the angle information of the BS-RIS link to design the analog beamformer. Specifically, the transmit beam of the m_b -th subarray designed for the u -th UE should be aligned to the m_r -th subarray of the RIS, where the u -th UE is assisted by the m_r -th subarray of the RIS. Therefore, analog beamformer $\mathbf{F}_{\text{RF}} = \text{blkdiag}(\mathbf{f}_1, \dots, \mathbf{f}_{m_b}, \dots, \mathbf{f}_{M^{\text{B}}})$ can be simply designed for alignment between BS and RIS subarrays according to Equation (9).

4.2.2 DL-Based RIS Refraction Phase Design

Optimizing a common RIS phase shared by all the subcarriers is a crucial challenge in a RIS-aided OFDM system. In the THz broadband case, there exists a non-negligible beam squint effect for different subcarriers [48]. 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, a transformer-based RIS phase design network (RPDN) is proposed in Figure 6, to design the RIS refraction phase matrix.

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

$$\bar{\mathcal{H}} = [\bar{\mathbf{H}}[1], \dots, \bar{\mathbf{H}}[u], \dots, \bar{\mathbf{H}}[U]], \quad (35)$$

where $\bar{\mathbf{H}}[u] = [\Re\{\hat{\mathbf{H}}[u]\}, \Im\{\hat{\mathbf{H}}[u]\}] \in \mathbb{R}^{K \times 2N^{\text{R}}}$ 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 $\mathbf{X}_r \in \mathbb{R}^{U \times N^{\text{R}}/U}$ is obtained by the mean operation over the subcarrier dimension of the transformer encoder's output. Then \mathbf{X}_r is flattened as $\mathbf{x}_r \in \mathbb{R}^{N^{\text{R}} \times 1}$, and passes through the activation function to generate the RIS phase vector $\mathbf{v} \in \mathbb{C}^{N^{\text{R}} \times 1}$ that satisfies the constant modulus constraint, i.e.,

$$\mathbf{v} = e^{j2\pi \cdot \text{Sigmoid}(\mathbf{x}_r)}. \quad (36)$$

Finally, the RIS phase matrix $\Phi \in \mathbb{C}^{N^{\text{R}} \times N^{\text{R}}}$ is obtained through diagonalization. The overall process

486 of the RIS refraction phase design, namely, the transformer-based RPDN, can be expressed as

$$\Phi = f_{\text{RIS}}(\bar{\mathcal{H}}; \mathcal{W}_R), \quad (37)$$

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

488 4.2.3 Knowledge-Data Dual-Driven Digital Beamformer Design

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

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

492 The true equivalent baseband channel $\mathbf{h}_{\text{eq}}[u, k]$ has a similar form to Equation (38), given the
 493 designed Φ and \mathbf{F}_{RF} . Therefore, the Equation (34) can be simplified as

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

494 Note that since the analog beamformer and the RIS refraction phase have been specifically designed
 495 in Subsections 4.2.1 and 4.2.2 respectively, the above problem (39) is a classic baseband beamforming
 496 problem, which can be solved with standard linear beamforming schemes, such as the regularized
 497 ZF (RZF) or iterative weighted minimum mean-square error (WMMSE) algorithm. Taking the
 498 latter as an example, the iterative WMMSE algorithm is designed to solve the optimization (39) by
 499 addressing the equivalent MMSE problem specified in (40) below, which has the identical optimal
 500 solution $\mathbf{F}_{\text{BB}}[k], \forall k$ to the problem (39).

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

501 where $\bar{w}_{u,k} = \{\bar{\mathbf{W}}\}_{u,k}$ is the weight of the u -th user on the k -th subcarrier, $e_{u,k} = \mathbb{E}\{|\hat{s}[u, k] -$
 502 $s[u, k]|^2\}$ is the MSE between the transceiver symbols under the independence assumption of $s[u, k]$
 503 and $n[u, k]$, while $\hat{s}[u, k] = \bar{u}_{u,k} y[u, k]$ is the estimated data symbol at the UE-side, and $\bar{u}_{u,k} = \{\bar{\mathbf{U}}\}_{u,k}$
 504 is the receiver gain of the u -th UE on the k -th subcarrier. According to [53], the above problem
 505 is convex in each individual optimization variable. This property enables each subproblem to have
 506 a closed-form solution, given the other optimization variables. Then, the optimization (40) can
 507 be solved by a block coordinate descent (BCD) iterative algorithm. Here, we depict the iterative
 508 WMMSE beamforming design algorithm in Algorithm 1.

509 However, the iterative WMMSE algorithm typically imposes a large number of iterations with
 510 long running time. Furthermore, the BS can only acquire the imperfect estimated CSI $\hat{\mathbf{h}}_{\text{eq}}[u, k]$, and
 511 it is difficult for the traditional digital beamforming algorithms, such as Algorithm 1, to overcome the
 512 interference induced by the imperfect CSI. Thus, we propose the knowledge-data dual-driven digital
 513 beamforming network, as shown in Figure 6, which utilizes the transformer encoder to directly learn
 514 the parameters of the iterative WMMSE algorithm from the imperfect CSI for better interference
 515 elimination and shorter running time.

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

$$\bar{\mathbf{W}} = \mathbf{X}_{[:, :U]}^{\text{T}} + \text{j} \mathbf{X}_{[:, U:2U]}^{\text{T}}, \quad (41)$$

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

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

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

520 The proposed knowledge-data dual-driven digital beamformer design can be represented by

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

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

Algorithm 1 Iterative WMMSE beamforming design algorithm

- 1: **Initialize** $\mathbf{F}_{\text{BB}}[k]$ that meets $\|\mathbf{F}_{\text{RF}} \mathbf{F}_{\text{BB}}[k]\|_F^2 = M^{\text{B}}$, set the maximum iteration number I_{max} , and the current iteration index $t = 0$;
 - 2: **repeat**
 - 3: **Update** $\{\bar{\mathbf{U}}\}_{u,k}$: $\bar{u}_{u,k} = \left(\sum_{i=1}^U |\mathbf{h}_{\text{eq}}[u, k] \mathbf{f}_{\text{BB}}[i, k]|^2 + \sigma_n^2 \right)^{-1} \mathbf{h}_{\text{eq}}[u, k] \mathbf{f}_{\text{BB}}[u, k], \forall u, k$;
 - 4: **Update** $\{\bar{\mathbf{W}}\}_{u,k}$: $\bar{w}_{u,k} = (1 - \bar{u}_{u,k}^* \mathbf{h}_{\text{eq}}[u, k] \mathbf{f}_{\text{BB}}[u, k])^{-1}, \forall u, k$;
 - 5: **Update** $\mathbf{f}_{\text{BB}}[u, k]$: $\mathbf{f}_{\text{BB}}[u, k] = \bar{u}_{u,k} \bar{w}_{u,k} \left(\sum_{i=1}^U \bar{w}_{i,k} |\bar{u}_{i,k}|^2 \mathbf{h}_{\text{eq}}^{\text{H}}[i, k] \mathbf{h}_{\text{eq}}[i, k] + \mu_k \mathbf{I} \right)^{-1} \mathbf{h}_{\text{eq}}^{\text{H}}[u, k]$,
 where $\mu_k = \sum_{j=1}^U \frac{\sigma_n^2}{M^{\text{B}}} \bar{w}_{j,k} |\bar{u}_{j,k}|^2, \forall u, k$;
 - 6: $t = t + 1$;
 - 7: **until** $t \geq I_{\text{max}}$
 - 8: Scale $\mathbf{F}_{\text{BB}}[k]$ to meet the transmit power constraint.
-

522 4.2.4 Training Strategy

523 We take every U channel samples (i.e., the channels of U UEs) in the training set of the channel
 524 estimation stage as a group to form a training set at the beamforming design stage, which is denoted
 525 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
 526 $\{\mathbf{H}[u], 1 \leq u \leq U\}$, where $\mathbf{H}[u]$ is the spatial-frequency RIS-UE channel of the u -th UE.

527 $\{\mathbf{H}[u], 1 \leq u \leq U\}$ are inputted to the trained SFDCetra network to obtain the estimated
 528 channels $\{\hat{\mathbf{H}}[u], 1 \leq u \leq U\}$, which form the input to the proposed network. Since imperfect CSI
 529 will reduce the sum rate upper bound, to ensure a faster learning process, we apply a teacher forcing
 530 technique [54] at the early stage of training by feeding the perfect CSI $\{\mathbf{H}[u], \forall u\}$ to the proposed
 531 network. At the off-line training stage, we consider E2E training to jointly optimize the hybrid
 532 beamforming and RIS phase, i.e., the parameters of the entire network are trained by minimizing
 533 the negative sum rate. Thus, the loss function is written as

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

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

535 Results and Discussion

536 5 Numerical Results

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

540 5.1 Simulation Settings

541 5.1.1 Communication Scenario Set up

542 In simulations, the BS is deployed on the top of a building of height 30 m, and the RIS is installed
 543 on a window surface on one floor of another building. As shown in Figure 1(b), the BS (RIS) is
 544 equipped with $M^B = M_y^B M_z^B = 4$ ($M^R = M_y^R M_z^R = 4$) subarrays on the yz -plane, where $M_y^B = 2$
 545 ($M_y^R = 2$) and $M_z^B = 2$ ($M_z^R = 2$). Each subarray is a UPA with $N_{\text{sub}}^B = N_y^B N_z^B = 64$ ($N_{\text{sub}}^R =$
 546 $N_y^R N_z^R = 64$) isotropically radiating elements, where $N_y^B = 8$ ($N_y^R = 8$) and $N_z^B = 8$ ($N_z^R = 8$).
 547 Therefore, the number of elements of the complete array at the BS (RIS) is $N^B = M^B N_{\text{sub}}^B = 256$
 548 ($N^R = M^R N_{\text{sub}}^R = 256$). For simplicity, the BS and RIS are assumed to meet the parallel symmetric
 549 array arrangement with a distance of $D = 20$ m. The central frequency is $f_c = 0.3$ THz with
 550 bandwidth $f_s = 1$ GHz. The number of OFDM subcarriers is $K = 128$ and the BS's antenna gain
 551 is $G_T = 10$ dBi. Given the above parameter settings, the subarray intervals of both the BS and
 552 the RIS are calculated from Equation (1) as $d_{sy}^B, d_{sz}^B, d_{sy}^R, d_{sz}^R = 96.5\lambda$ for obtaining the multi-stream
 553 multiplexing gain over the LoS path.

554 Figure 7 depicts the schematic diagram of RIS-UE channel model for the indoor environment,
 555 where the positions of the RIS, UEs, and scatterers are indicated by blue, red, and green circles,
 556 respectively. The RIS-UE LoS path is depicted by a red solid line, and the NLoS link via a scatterer
 557 is represented by a black dotted line. We assume that $U = 4$ UEs are randomly distributed over
 558 the xy -plane of the rectangular room ($W_x = 5$ m, $W_y = 10$ m), and the height of UEs is 1 m lower
 559 than the RIS. The number of available NLoS paths (scatterers) is set to $L_p = 5$, implying that
 560 only a single-bounce scattering mode is considered. The reflection coefficient parameters β_{RC} are
 561 set to $\mu_R = -5$, $\sigma_R = 2$. The noise power spectrum density at the UEs is $\sigma_{\text{NSD}}^2 = -174$ dBm/Hz.
 562 Thus, the power of the AWGN is $\sigma_n^2 = \sigma_{\text{NSD}}^2 f_s / K = -105$ dBm. The RIS-UE channel samples are
 563 generated using Equation (5), where the UEs and scatterers are distributed randomly each time.

564 5.1.2 SFDCetra Network Parameter Configuration

565 In the CSI feedback network, the linear embedding layer of the transformer encoder has $d_T = 256$
 566 neurons. In the transformer encoder, the number of the encoder layers is $L_T = 3$, where the
 567 number of heads is $h = 8$ and the position-wise MLP sub-layer has 2 fully-connected layers with
 568 $4d_T$ and d_T neurons, respectively, while the dimension of the output linear layer is $2M$. In the
 569 sub-channel estimation network, the linear layer is a $2N_s^R$ -dimensional fully-connected layer and the
 570 hyperparameters of the transformer encoder are the same as those of the CSI feedback network.
 571 As for the channel extrapolation network, a convolutional layer has 7×7 kernel and 2 filters. The
 572 parameters of the rearranged operation are $L = 16$ and $N_p = 128$, as well as the number of neurons
 573 in the linear layer is $d_M = 512$. The number of mixer layers is $L_M = 6$, where each mixer layer
 574 consists of two MLP blocks, and the numbers of neurons in the MLP blocks are set to $2N_p$, N_p , $2d_M$,

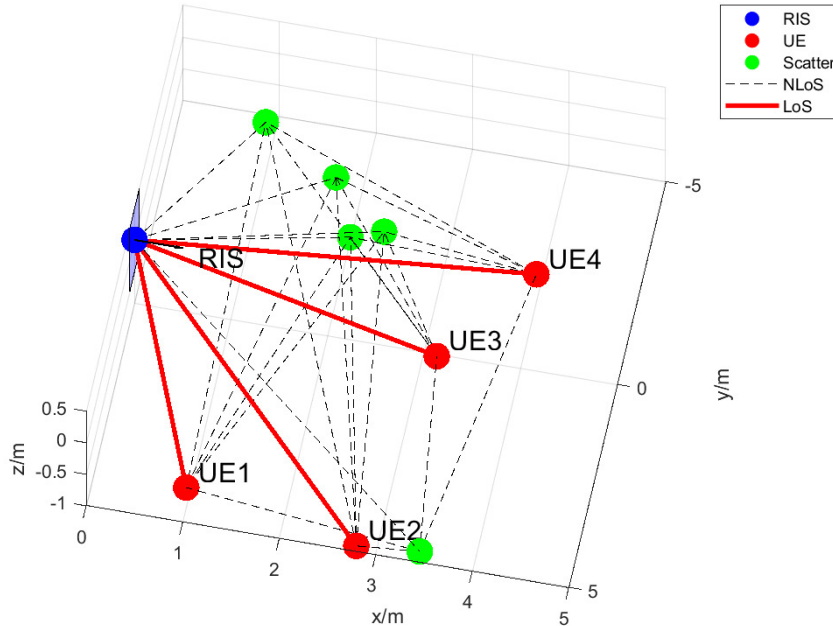


Figure 7: The schematic diagram of RIS-UE channel model for the indoor environment.

575 and d_M , respectively. The above structural parameters of the SFDCEtra network are empirically
 576 found to be appropriate.

577 The dataset is divided into three distinct subsets, namely the training set, validation set, and
 578 testing set, which consist of 102400, 10240, and 10240 samples respectively. Unless otherwise speci-
 579 fied, simulations adopt the uniform element selection strategy. When considering the learning-based
 580 element selection strategy, the weight factor γ is 0.9. At the network training stage, the Adam op-
 581 timizer is adopted to update the network weight parameters and the learning rate varies depending
 582 on the *warmup* mechanism [50]. The batch size is set to 512 with 200 epochs.

583 5.1.3 HBFRPD Network Parameter Configuration

584 Again we determine the appropriate structural parameters of the HBFRPD network empirically.
 585 Specifically, in the RIS phase design network, the linear embedding layer of the transformer encoder
 586 has $d_B = 128$ neurons. In the transformer encoder, the number of the encoder layers is $L_B = 3$,
 587 where the number of heads is $h = 8$ and the position-wise MLP sub-layer has 2 fully-connected layers
 588 with $4d_B$ and d_B neurons, respectively, while the output linear layer of the transformer encoder has
 589 $N^R/U = 64$ neurons. In the digital beamforming network, the hyperparameters of the transformer
 590 encoder are the same as those of the RIS phase design network, and the output linear layer of the
 591 transformer encoder has $4U$ neurons.

592 We take each U channel samples as a group to form a dataset, which is composed of three parts,
 593 with a sample size of 25600, 2560, and 2560 respectively. The batch size is set to 32 with 180 epochs.

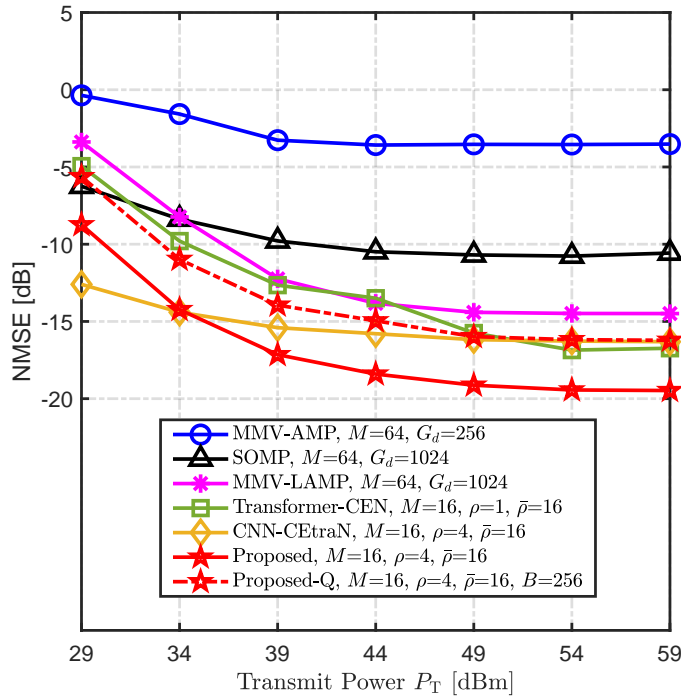


Figure 8: NMSE performance comparison of different channel estimation schemes versus transmit power P_T .

5.2 DL-Based Spatial-Frequency Domain Channel Extrapolation

Since the RIS element can only passively receive EM waves, selecting partial elements of the RIS array would reduce the signal energy radiated into the room. For the fair comparison between different schemes, we adopt the same transmit power instead of the same SNR as the comparison criterion to avoid ignoring the performance differences induced by the number of activated RIS elements. Specifically, as illustrated in Figure 8, we show the NMSE performance of the different schemes with different transmit power P_T . The number of NLoS paths is $L_p = 5$. We consider three model-based channel estimation benchmark algorithms, namely, the SOMP algorithm [55], the multiple-measurement-vector approximate message passing (MMV-AMP) algorithm [56], and the model-driven MMV learned AMP (MMV-LAMP) network [57], which utilize $M = 64$ OFDM symbols on all subcarriers and then directly estimate the complete channel. For the SOMP and MMV-LAMP schemes, the redundant dictionary with an oversampling ratio of 4 is utilized to further improve the performance, i.e., the number of codewords is $G_d = 1024$. However, for the MMV-AMP approach, the requirement of independent and identically distributed elements in the measurement matrix precludes the use of a redundant dictionary (i.e., $G_d = N^R = 256$). Since data-driven DL algorithms have the potential to achieve better performance, we also compare our proposed DL-based SFDCEtra network with the transformer-based channel estimation network (Transformer-CEN) [51] and the CNN-based channel extrapolation network (CNN-CEtraN) [28]. For these methods, we set $M = 16$ OFDM symbols and $\bar{\rho} = 16$ subcarrier compression ratio. The transformer-based scheme turn on all RIS elements, i.e., $\rho = 1$, and directly estimates the complete channel. Both the CNN-based and our proposed channel extrapolation schemes consider the element compression ratio of $\rho = 4$ to perform partial channel extrapolation. Note that for fairness, the above model- and data-driven algorithms do not consider the quantization of CSI feedback information. Therefore, we additionally consider the proposed scheme with $B = 256$ feedback bits generated via a 2-bit quantizer, denoted as ‘Proposed-Q’.

It can be observed from Figure 8 that our proposed channel extrapolation scheme outperforms the other schemes considerably in terms of NMSE performance while imposing a smaller pilot overhead. This is because exploiting spatial-frequency correlations 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 perception 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 a 2-bit quantizer, ‘Proposed-Q’, can still achieve very good performance. These results demonstrate that the proposed channel extrapolation scheme can achieve high reconstruction performance while ensuring low pilot and feedback overheads.

We further investigate the robustness of the proposed DL-based channel extrapolation scheme with respect to the number of multipath L_p in Figure 9. The proposed DL-based channel extrapolation scheme is trained offline using channel samples that contain $L_p = 5$ multipath components. As

634 depicted in Figure 9, at the online estimation stage, the proposed scheme demonstrates its ability
 635 to estimate channels with different L_p without the need for retraining the entire network. Thus,
 636 our proposed scheme exhibits superior robustness and generalization capabilities in various channel
 637 conditions.

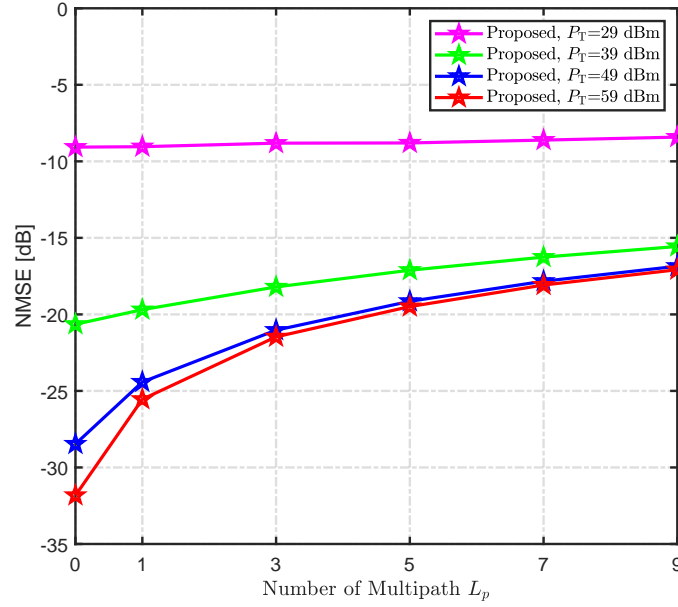


Figure 9: NMSE performance comparison of the proposed scheme versus the number of multipath L_p , given $\rho = 4$, $\bar{\rho} = 16$ and $M = 16$. Offline training is based on the channel samples with $L_p = 5$ multipath components.

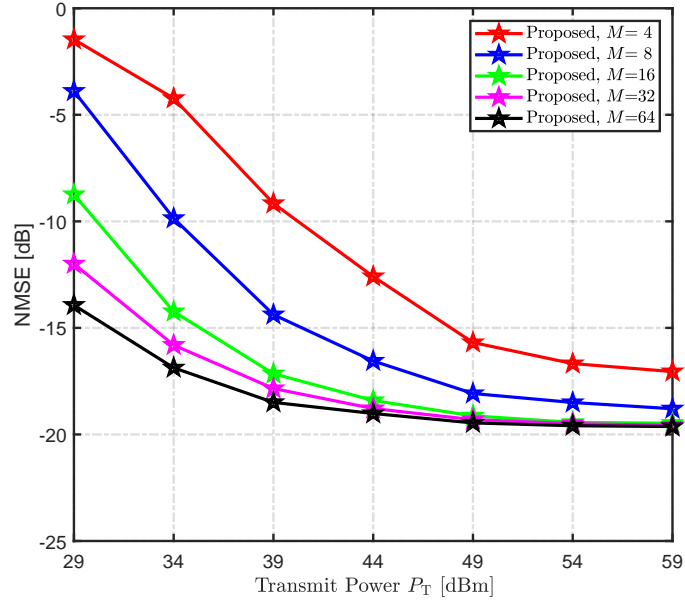


Figure 10: NMSE performance comparison of the proposed scheme with different pilot numbers versus transmit power P_T , given $\rho = 4$, $\bar{\rho} = 16$, $L_p = 5$.

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

646 Figure 11 depicts the NMSE performance of the proposed DL-based channel extrapolation scheme
647 versus the element compression ratio ρ , with three ESEs. Specifically, the curve labeled by ‘Uniform’
648 corresponds to the uniform selection strategy, the curve labeled by ‘Random’ represents the random
649 selection strategy, while the other three labeled by ‘DL-based with 200 epochs’, ‘DL-based with 300
650 epochs’, and ‘DL-based with 400 epochs’ use the DL-based element selection strategy. As expected,
651 the NMSE improves as the element compression ratio ρ decreases. This is largely due to two reasons:
652 1) As the number of selected RIS elements increases, or the element compression ratio ρ decreases,
653 the received signal power will increase, thus improving the estimation accuracy of channel extrap-
654 olation input (i.e., sub-channel estimate), and 2) The received pilot signal can provide more channel
655 information when more RIS elements are selected. However, this does not imply that we can obtain
656 the best performance by choosing the lowest element compression ratio (or performing complete ob-
657 servations directly without extrapolation). Indeed, the channel extrapolation performance heavily
658 depends on the amount of wireless communication transmission resources, the accuracy of the sub-
659 channel estimation, and the amount of selected RIS elements (i.e., the dimension of the sub-channel).
660 Only when the transmission resources are sufficient, can the gain provided by more selected RIS ele-
661 ments be seen clearly. Moreover, we can observe that the performance gap between different element
662 selection strategies is not obvious at low compression ratios. However, at a high compression ratio
663 (e.g., $\rho > 8$), the performance difference can be seen clearly as ‘Uniform’ < ‘Random’ < ‘DL-based’,
664 which demonstrates the effectiveness of the proposed approach. Since the aperture of the random
665 pattern is statistically larger than that of the fixed uniform pattern, the random selection strategy
666 is better than that of the uniform selection strategy especially at a high compression ratio. The
667 performance of the DL-based approach is relatively better than that of the first two approaches
668 after reaching a sufficient number of training epochs—specifically 300 epochs in this scenario, as the
669 learning of the selection network requires more epochs to converge.

670 To fully illustrate the effectiveness of our proposed scheme, its channel extrapolation module is
671 verified separately. To do so, we fix the compression ratio of RIS elements to 4, i.e., $N_s^R = 64$. First,
672 the least squares (LS), the SOMP, and the proposed transformer-based algorithm are utilized for
673 sub-channel estimation, and the results are shown in Figure 12(a). Observe that the NMSE of the
674 SOMP-based sub-channel estimation with $M = 64$ pilot symbols is significantly better than that of
675 the LS-based sub-channel estimation with $M = 64$ pilot symbols, particularly at low transmit power
676 P_T . Furthermore, the NMSE of our transformer-based sub-channel estimation algorithm with only
677 $M = 16$ pilot symbols is considerably better than that of the SOMP-based sub-channel estimation
678 with $M = 64$ pilot symbols. Then, we input the sub-channels estimated by different algorithms into

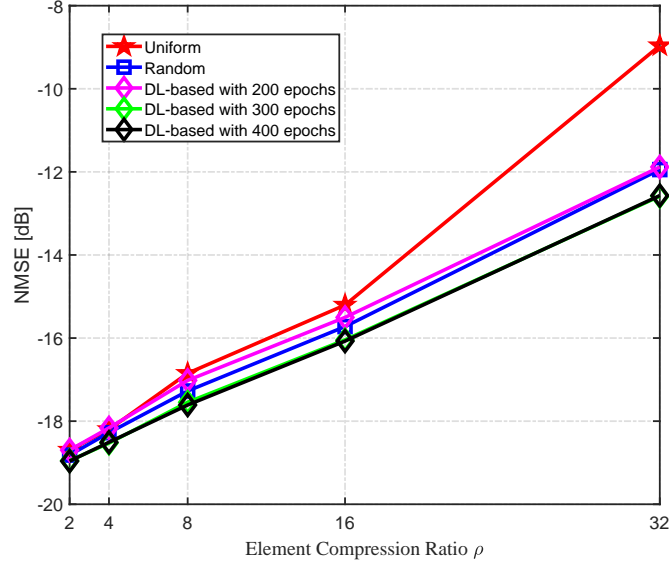


Figure 11: NMSE performance comparison of the proposed scheme with different element selection strategies versus element compression ratio ρ , given $\bar{\rho} = 16$, $L_p = 5$, $M = 16$, $P_T = 44$ dBm.

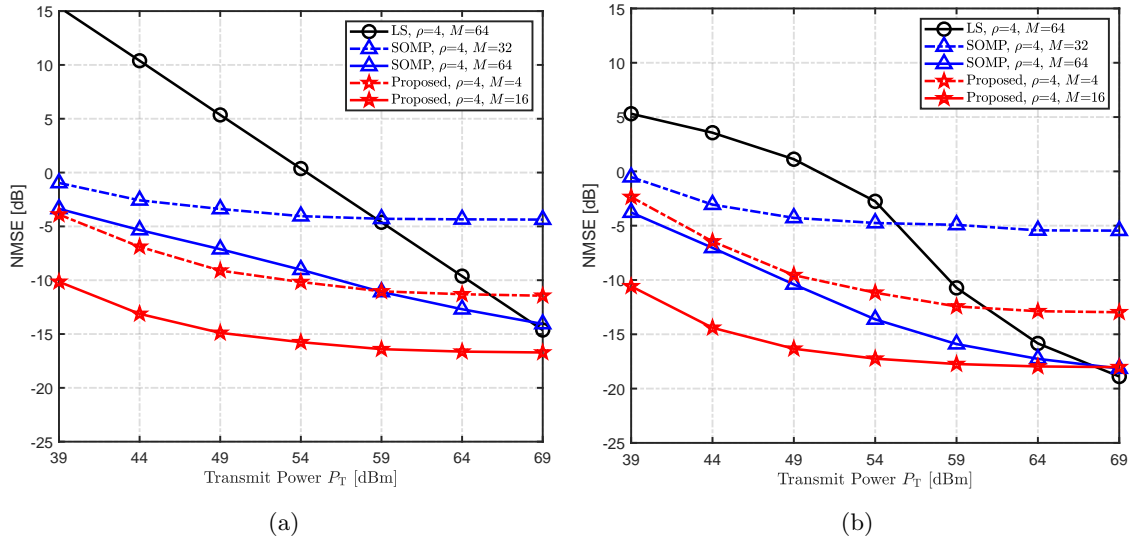


Figure 12: (a) NMSE performance comparison of different sub-channel estimation schemes versus transmit power P_T ; and (b) NMSE performance of channel extrapolation versus transmit power P_T for different sub-channel estimation schemes.

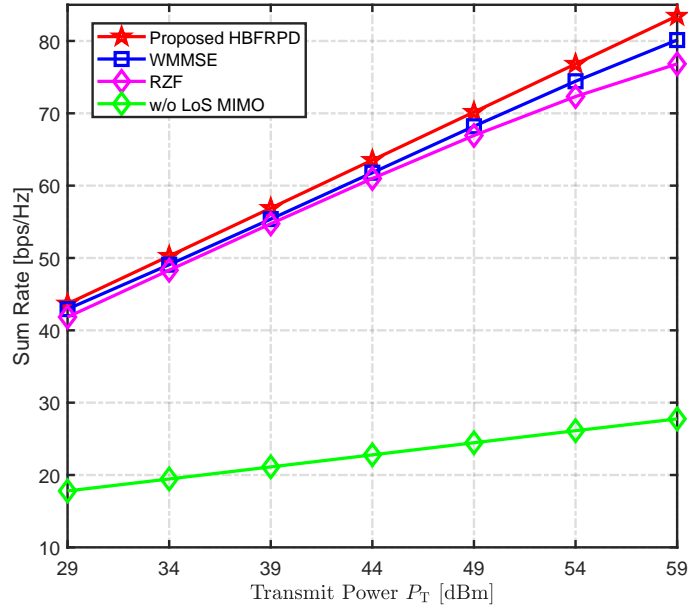


Figure 13: Sum rates achieved by different schemes versus transmit power P_T given the perfect CSI. The actual transmit power of the ‘w/o LoS MIMO’ case is $UP_T = 4P_T$.

679 the trained channel extrapolation network $f_{\text{SFDE}}(\cdot)$, which outputs the estimation of the complete
 680 channel. The corresponding results are shown in Figure 12(b). Observe that the NMSE performance
 681 of the complete channel extrapolated from our channel extrapolation network is even better than
 682 the NMSE of the estimated low-dimensional sub-channel, without any additional pilot overhead.
 683 This shows that our proposed channel extrapolation network can not only be used for DL-based
 684 communication architecture, but also be combined with traditional algorithms to significantly reduce
 685 resource overhead. Therefore, we conclude that the proposed DL-based spatial-frequency domain
 686 channel extrapolation scheme can learn a latent mapping among channel elements to significantly
 687 reduce the pilot overhead while achieving the same or better channel estimation performance.

688 5.3 DL-Based Hybrid Beamforming and RIS Phase Design

689 Figure 13 shows the sum rates of total UEs achieved by different schemes under the perfect
 690 CSI case. We considered two comparison schemes, both of which adopt the analog beamforming
 691 design discussed in Subsection 4.2.1 as well as the beam alignment-based RIS phase design. In the
 692 beam alignment-based RIS phase design, the beam of each subarray is aligned to the corresponding
 693 associated UE. For digital beamforming design, these two comparison schemes adopt the RZF and
 694 iterative WMMSE algorithms, respectively, thus they are abbreviated as ‘RZF’ and ‘WMMSE’,
 695 respectively. We can observe that our proposed HBFRPD scheme has better performance than
 696 other schemes and the superiority is more evident as the transmit power increases. Besides, another
 697 advantage of our proposed HBFRPD scheme is that it does not require $\mathbf{F}_{\text{BB}}[k], \forall k$, in an iterative
 698 manner. Thus, it runs much faster than the iterative WMMSE algorithm. We also analyze the
 699 performance gain provided by LoS MIMO architecture. Considering the case without LoS MIMO

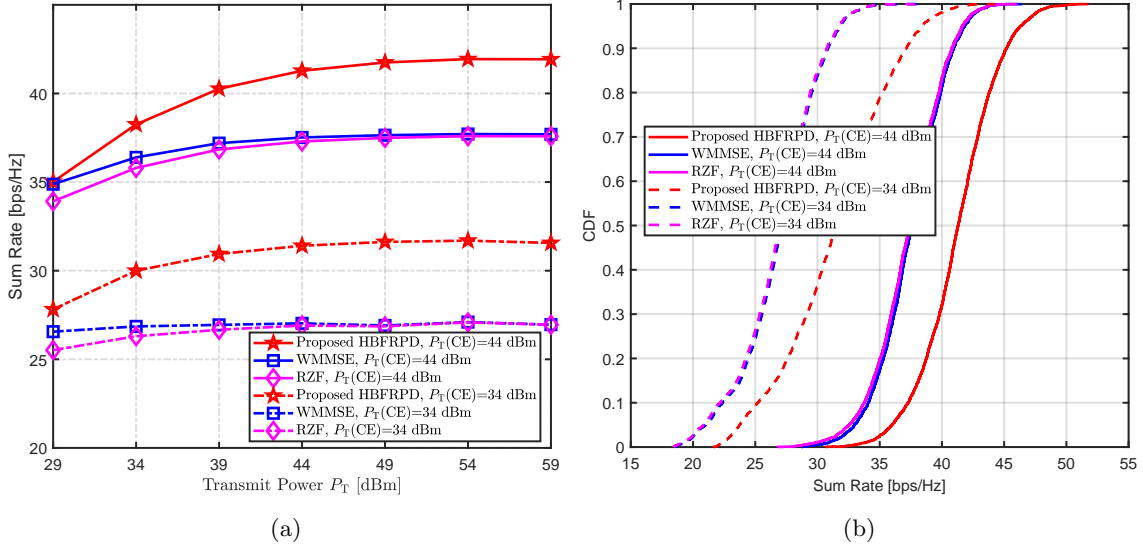


Figure 14: (a) Sum rates achieved by different schemes versus transmit power P_T under the imperfect CSI case, and (b) The CDFs of the sum rates achieved by different schemes under the imperfect CSI case, given $P_T = 44$ dBm. We have $\rho = 4$, $\bar{\rho} = 16$, $L_p = 5$ and $M = 16$.

700 array structure (i.e., both the BS and RIS use conventional UPA arrays), the BS-RIS channel is a
 701 single LoS path with rank 1, which only provides single stream data transmission. To ensure fairness,
 702 the transmit power in the absence of LoS MIMO is equal to that with LoS MIMO, i.e., the transmit
 703 power in the absence of LoS MIMO is actually $UP_T = 4P_T$. By calculating the sum rate, we obtain
 704 the green curve in Figure 13. It can be seen that the sum rate with LoS MIMO is much higher than
 705 that without LoS MIMO. This is because the LoS MIMO architecture can increase the sum rate
 706 linearly benefited from extra spatial multiplexing gain, while the conventional array architecture can
 707 only provide log-level growth as the SINR increases.

708 Although most schemes can achieve good sum rate performance under perfect CSI, the sum rate
 709 of multi-users is actually degraded due to inter-user interference induced by CSI error. Figure 14(a)
 710 illustrates the sum rate performance of the different schemes with imperfect CSIs estimated at two
 711 different transmit powers $P_T(\text{CE})$. Compared with the case of perfect CSI, the sum rate degrades
 712 significantly with the decrease of CSI estimation accuracy, i.e., with the decrease of the transmit
 713 power at the channel estimation stage. It can be clearly seen that due to the inter-user interference
 714 induced by CSI errors, the sum rates of the RZF and iterative WMMSE schemes barely increase
 715 with transmit power. Moreover, our proposed HBFRPD scheme exhibits a significant performance
 716 gain over the RZF and iterative WMMSE algorithms in the presence of CSI estimation errors. This
 717 result indicates that our proposed scheme can mitigate the interference caused by CSI errors and
 718 hence has better robustness to inaccurate CSI than the other schemes.

719 The cumulative distribution functions (CDFs) characterizing the sum rate performance achieved
 720 by different schemes are shown in Figure 14(b). Here, we consider the transmit power $P_T = 44$ dBm
 721 at the data transmission stage. Figure 14(b) shows that when the transmit power is $P_T(\text{CE}) =$
 722 34 dBm at the channel estimation stage, the proposed HBFRPD network has a probability of about
 723 64.6% to achieve a sum rate exceeding 30 bps/Hz, while the other two schemes can only achieve

724 16.3%. When the transmit power is $P_T(\text{CE}) = 44 \text{ dBm}$ at the channel estimation stage, our
 725 HBFRPD network has a probability of about 68.8% to achieve a sum rate exceeding 40 bps/Hz,
 726 which is significantly higher than the other two schemes. This result again confirms the superior
 727 performance of our proposed HBFRPD network over existing conventional schemes.

728 5.4 Computational Complexity Analysis

729 The computational complexity analysis of different schemes at the inference stage is presented
 730 in Table 1. All the numerical results are obtained on a PC with Intel(R) Core(TM) i9-10980XE
 731 CPU @ 3.00GHz and an Nvidia GeForce RTX 3090 GPU. The DL-based methods and the existing
 732 solutions are implemented on the PyCharm framework. The details are further elaborated as follows.

733 1) Channel estimation schemes: In the SOMP algorithm [55], correlation operation imposes sig-
 734 nificant computational complexity, where I is the number of iterations. The MMV-AMP algorithm
 735 [56] mainly requires matrix multiplication operations, but a large number of iterations I increases
 736 its computational complexity. The MMV-LAMP algorithm [57] has a low computational complexity
 737 because DL reduces the required number of iterations. The Transformer-CEN [51] also has a low
 738 computational complexity, and the main sources of its computational complexity come from self-
 739 attention and MLP sublayers. In the CNN-CEtraN [28], convolutional layers introduce significant
 740 computational complexity. By contrast, the MLP-mixer layers provide the majority of the com-
 741 putational complexity in our proposed SFDCEtra network, which is much lower than that of the
 742 CNN-CEtraN. We further meticulously count the numbers of floating-point operations per second
 743 (FLOPs) and run times per sample on CPU for different schemes in Table 1. Observe that at the
 744 inference stage, the FLOPs and run time per sample of the proposed scheme are lower than most
 745 benchmarks. Specifically, our SFDCEtra network imposes the second lowest run time per sample,
 746 and only the MMV-LAMP and Transformer-CEN 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 [53], 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 Conclusions

756 In this paper, we have proposed a DL-based transmission architecture 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

Table 1: Computational Complexity of Different Schemes.

Channel estimation scheme	Complexity	FLOPs	Run time/s
SOMP	$\mathcal{O}(G_d K M I + G_d^2 K I)$	4.707 G	0.1130
MMV-AMP	$\mathcal{O}(M K N^R I)$	5.337 G	0.7482
MMV-LAMP	$\mathcal{O}(M G_d K I)$	0.341 G	0.0689
Transformer-CEN	$\mathcal{O}(L_T (K d_T^2 + K^2 d_T))$	0.362 G	0.0103
CNN-CEtraN	$\mathcal{O}(Z_5^2 N^R K C_{32}^2)$	10.16 G	0.6039
Proposed	$\mathcal{O}(L_M (N_p^2 d_M + N_p d_M^2))$	1.066 G	0.0248
Beamforming scheme	Complexity	FLOPs	Run time/s
RZF	$\mathcal{O}((2U(M^B)^2 + (M^B)^3)K)$	24.58 K	0.1389
WMMSE	$\mathcal{O}(IK(U^2(M^B)^2 + U(M^B)^3))$	8.192 M	0.9184
Proposed	$\mathcal{O}(L_B U(K d_B^2 + K^2 d_B))$	0.512 G	0.0616

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.
 765 Moreover, the results have also demonstrated that our proposed beamforming scheme is superior to
 766 the existing designs in terms of achievable sum rate performance and robustness to imperfect CSI.
 767 Potential future research directions based on the outcomes of this paper include the practical discrete
 768 phase shifter, the analysis of the complex near-field channel, and sensing-aided communications.

769 Acknowledgments

770 **General:** Thank Yifei Zhang, Ziwei Wan, and others for any contributions.

771 **Funding:** This work was supported by the Natural Science Foundation of China (NSFC) [Grant
 772 62071044]; in part by the Shandong Province Natural Science Foundation [Grant ZR2022YQ62];
 773 and in part by the Beijing Nova Program.

774 **Author Contributions:** Y. Wang, Z. Gao, and S. Chen contributed to the algorithms and anal-
 775 ysis. Y. Wang and Z. Gao prepared the manuscript; all the authors contributed to the manuscript
 776 editing. C. Hu and D. Zheng initiated the presented concept and supervised the research.

777 **Competing Interests:** The authors declare that there is no conflict of interest regarding the
 778 publication of this paper.

779 Data Availability

780 Data are available from the corresponding author on reasonable request.

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 Journal of the Communications Society*, vol. 1, pp. 1–32, 2020.
- [2] C. Lin and G. Y. Li, “Indoor Terahertz communications: How many antenna arrays are needed?” *IEEE Transactions on Wireless Communications*, vol. 14, no. 6, pp. 3097–3107, 2015.
- [3] F. Sotiraki and W. Yu, “Hybrid digital and analog beamforming design for large-scale antenna arrays,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 10, no. 3, pp. 501–513, 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 Journal on Selected Areas in Communications*, vol. 38, no. 11, pp. 2450–2525, 2020.
- [5] C. Huang *et al.*, “Holographic mimo surfaces for 6g wireless networks: Opportunities, challenges, and trends,” *IEEE Wireless Communications*, vol. 27, no. 5, pp. 118–125, 2020.
- [6] M. D. Renzo *et al.*, “Smart radio environments empowered by reconfigurable ai meta-surfaces: An idea whose time has come,” *EURASIP Journal on Wireless Communications and Networking*, vol. 2019, no. 1, pp. 1–20, 2019.
- [7] C. Huang, A. Zappone, G. C. Alexandropoulos, M. Debbah, and C. Yuen, “Reconfigurable intelligent surfaces for energy efficiency in wireless communication,” *IEEE Transactions on Wireless Communications*, vol. 18, no. 8, pp. 4157–4170, 2019.
- [8] M. Di Renzo *et al.*, “Smart radio environments empowered by reconfigurable intelligent surfaces: How it works, state of research, and the road ahead,” *IEEE Journal on Selected Areas in Communications*, vol. 38, no. 11, pp. 2450–2525, 2020.
- [9] G. C. Alexandropoulos, K. Stylianopoulos, C. Huang, C. Yuen, M. Bennis, and M. Debbah, “Pervasive machine learning for smart radio environments enabled by reconfigurable intelligent surfaces,” *Proceedings of the IEEE*, vol. 110, no. 9, pp. 1494–1525, 2022.
- [10] M. Wu, Z. Gao, Y. Huang, Z. Xiao, D. W. K. Ng, and Z. Zhang, “Deep learning-based rate-splitting multiple access for reconfigurable intelligent surface-aided terahertz massive mimo,” *IEEE Journal on Selected Areas in Communications*, vol. 41, no. 5, pp. 1431–1451, 2023.
- [11] K. T. Selvan and R. Janaswamy, “Fraunhofer and fresnel distances: Unified derivation for aperture antennas,” *IEEE Antennas and Propagation Magazine*, vol. 59, no. 4, pp. 12–15, 2017.
- [12] M. Cui and L. Dai, “Channel estimation for extremely large-scale MIMO: Far-field or near-field?” *IEEE Transactions on Communications*, vol. 70, no. 4, pp. 2663–2677, 2022.
- [13] 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 Transactions on Communications*, vol. 70, no. 2, pp. 1391–1406, 2022.

- 819 [14] D. Mishra and H. Johansson, "Channel estimation and low-complexity beamforming design
820 for passive intelligent surface assisted miso wireless energy transfer," in *ICASSP 2019 - 2019*
821 *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019,
822 pp. 4659–4663.
- 823 [15] P. Wang, J. Fang, H. Duan, and H. Li, "Compressed channel estimation for intelligent reflecting
824 surface-assisted millimeter wave systems," *IEEE Signal Processing Letters*, vol. 27, pp. 905–
825 909, 2020.
- 826 [16] L. Wei, C. Huang, G. C. Alexandropoulos, C. Yuen, Z. Zhang, and M. Debbah, "Channel
827 estimation for RIS-empowered multi-user MISO wireless communications," *IEEE Transactions*
828 *on Communications*, vol. 69, no. 6, pp. 4144–4157, 2021.
- 829 [17] A. M. Elbir, A. Papazafeiropoulos, P. Kourtessis, and S. Chatzinotas, "Deep channel learn-
830 ing for large intelligent surfaces aided mm-Wave massive MIMO systems," *IEEE Wireless*
831 *Communications Letters*, vol. 9, no. 9, pp. 1447–1451, 2020.
- 832 [18] S. Liu, Z. Gao, J. Zhang, M. D. Renzo, and M.-S. Alouini, "Deep denoising neural network
833 assisted compressive channel estimation for mmwave intelligent reflecting surfaces," *IEEE*
834 *Transactions on Vehicular Technology*, vol. 69, no. 8, pp. 9223–9228, 2020.
- 835 [19] J. An *et al.*, "Codebook-based solutions for reconfigurable intelligent surfaces and their open
836 challenges," *IEEE Wireless Communications*, 2022.
- 837 [20] L. Wei *et al.*, "Joint channel estimation and signal recovery for ris-empowered multiuser com-
838 munications," *IEEE Transactions on Communications*, vol. 70, no. 7, pp. 4640–4655, 2022.
- 839 [21] M. Liu, X. Li, B. Ning, C. Huang, S. Sun, and C. Yuen, "Deep learning-based channel esti-
840 mation for double-ris aided massive mimo system," *IEEE Wireless Communications Letters*,
841 vol. 12, no. 1, pp. 70–74, 2022.
- 842 [22] Z. Wan, Z. Gao, and M.-S. Alouini, "Broadband channel estimation for intelligent reflecting
843 surface aided mmwave massive mimo systems," in *ICC 2020-2020 IEEE International Con-*
844 *ference on Communications (ICC)*, IEEE, 2020, pp. 1–6.
- 845 [23] M. B. Mashhadi and D. Gündüz, "Pruning the pilots: Deep learning-based pilot design and
846 channel estimation for mimo-ofdm systems," *IEEE Transactions on Wireless Communications*,
847 vol. 20, no. 10, pp. 6315–6328, 2021.
- 848 [24] Z. Wan, Z. Gao, F. Gao, M. Di Renzo, and M.-S. Alouini, "Terahertz massive mimo with holo-
849 graphic reconfigurable intelligent surfaces," *IEEE Transactions on Communications*, vol. 69,
850 no. 7, pp. 4732–4750, 2021.
- 851 [25] S. Zhang, Y. Liu, F. Gao, C. Xing, J. An, and O. A. Dobre, "Deep learning based channel
852 extrapolation for large-scale antenna systems: Opportunities, challenges and solutions," *IEEE*
853 *Wireless Communications*, vol. 28, no. 6, pp. 160–167, 2021.
- 854 [26] B. Lin, F. Gao, S. Zhang, T. Zhou, and A. Alkhateeb, "Deep learning-based antenna selection
855 and CSI extrapolation in massive MIMO systems," *IEEE Transactions on Wireless Commu-*
856 *nications*, vol. 20, no. 11, pp. 7669–7681, 2021.

- 857 [27] M. Xu, S. Zhang, C. Zhong, J. Ma, and O. A. Dobre, “Ordinary differential equation-based
858 CNN for channel extrapolation over RIS-assisted communication,” *IEEE Communications*
859 *Letters*, vol. 25, no. 6, pp. 1921–1925, 2021.
- 860 [28] S. Zhang, S. Zhang, F. Gao, J. Ma, and O. A. Dobre, “Deep learning-based RIS channel
861 extrapolation with element-grouping,” *IEEE Wireless Communications Letters*, vol. 10, no. 12,
862 pp. 2644–2648, 2021.
- 863 [29] K. Ying, Z. Gao, S. Lyu, Y. Wu, H. Wang, and M.-S. Alouini, “GMD-based hybrid beam-
864 forming for large reconfigurable intelligent surface assisted millimeter-Wave massive MIMO,”
865 *IEEE Access*, vol. 8, pp. 19 530–19 539, 2020.
- 866 [30] B. Di, H. Zhang, L. Song, Y. Li, Z. Han, and H. V. Poor, “Hybrid beamforming for recon-
867 figurable intelligent surface based multi-user communications: Achievable rates with limited
868 discrete phase shifts,” *IEEE Journal on Selected Areas in Communications*, vol. 38, no. 8,
869 pp. 1809–1822, 2020.
- 870 [31] Y. Ahn and B. Shim, “Deep learning-based beamforming for intelligent reflecting surface-
871 assisted mmWave systems,” in *2021 International Conference on Information and Communi-*
872 *cation Technology Convergence (ICTC)*, 2021, pp. 1731–1734.
- 873 [32] C. Pradhan, A. Li, L. Song, B. Vucetic, and Y. Li, “Hybrid precoding design for reconfigurable
874 intelligent surface aided mmWave communication systems,” *IEEE Wireless Communications*
875 *Letters*, vol. 9, no. 7, pp. 1041–1045, 2020.
- 876 [33] S. Zhang, H. Zhang, B. Di, Y. Tan, Z. Han, and L. Song, “Beyond intelligent reflecting sur-
877 faces: Reflective-transmissive metasurface aided communications for full-dimensional coverage
878 extension,” *IEEE Transactions on Vehicular Technology*, vol. 69, no. 11, pp. 13 905–13 909,
879 2020.
- 880 [34] Y. Youn *et al.*, “Demo: Transparent intelligent surfaces for sub-6 GHz and mmWave B5G/6G
881 systems,” in *2022 IEEE International Conference on Communications Workshops (ICC Work-*
882 *shops)*, 2022, pp. 1–2.
- 883 [35] D. Kitayama, Y. Hama, K. Goto, K. Miyachi, T. Motegi, and O. Kagaya, “Transparent dy-
884 namic metasurface for a visually unaffected reconfigurable intelligent surface: Controlling trans-
885 mission/reflection and making a window into an RF lens,” *Optics Express*, vol. 29, no. 18,
886 pp. 29 292–29 307, 2021.
- 887 [36] Y. Chen, L. Yan, and C. Han, “Hybrid spherical- and planar-Wave modeling and DCNN-
888 powered estimation of Terahertz ultra-massive MIMO channels,” *IEEE Transactions on Com-*
889 *munications*, vol. 69, no. 10, pp. 7063–7076, 2021.
- 890 [37] X. Wang, Z. Lin, F. Lin, and L. Hanzo, “Joint hybrid 3D beamforming relying on sensor-based
891 training for reconfigurable intelligent surface aided TeraHertz-based multiuser massive MIMO
892 systems,” *IEEE Sensors Journal*, vol. 22, no. 14, pp. 14 540–14 552, 2022.

- 893 [38] S. Hong, C. Pan, H. Ren, K. Wang, K. K. Chai, and A. Nallanathan, “Robust transmission
894 design for intelligent reflecting surface-aided secure communication systems with imperfect
895 cascaded CSI,” *IEEE Transactions on Wireless Communications*, vol. 20, no. 4, pp. 2487–
896 2501, 2021.
- 897 [39] Z. Chen, J. Tang, X. Y. Zhang, Q. Wu, G. Chen, and K.-K. Wong, “Robust hybrid beamform-
898 ing design for multi-RIS assisted MIMO system with imperfect CSI,” *IEEE Transactions on*
899 *Wireless Communications*, pp. 1–1, 2022.
- 900 [40] W. Xu, L. Gan, and C. Huang, “A robust deep learning-based beamforming design for RIS-
901 assisted multiuser MISO communications with practical constraints,” *IEEE Transactions on*
902 *Cognitive Communications and Networking*, vol. 8, no. 2, pp. 694–706, 2022.
- 903 [41] P. Larsson, “Lattice array receiver and sender for spatially orthonormal MIMO communica-
904 tion,” in *2005 IEEE 61st Vehicular Technology Conference*, vol. 1, 2005, 192–196 Vol. 1.
- 905 [42] F. Bohagen, P. Orten, and G. E. Oien, “Optimal design of uniform planar antenna arrays
906 for strong line-of-sight MIMO channels,” in *2006 IEEE 7th Workshop on Signal Processing*
907 *Advances in Wireless Communications*, 2006, pp. 1–5.
- 908 [43] X. Song, W. Rave, N. Babu, S. Majhi, and G. Fettweis, “Two-level spatial multiplexing using
909 hybrid beamforming for millimeter-Wave backhaul,” *IEEE Transactions on Wireless Commu-*
910 *nications*, vol. 17, no. 7, pp. 4830–4844, 2018.
- 911 [44] L. Yan, C. Han, and J. Yuan, “Energy-efficient dynamic-subarray with fixed true-time-delay
912 design for Terahertz wideband hybrid beamforming,” *IEEE Journal on Selected Areas in Com-*
913 *munications*, vol. 40, no. 10, pp. 2840–2854, 2022.
- 914 [45] C. Han, A. O. Bicen, and I. F. Akyildiz, “Multi-ray channel modeling and wideband character-
915 ization for wireless communications in the Terahertz band,” *IEEE Transactions on Wireless*
916 *Communications*, vol. 14, no. 5, pp. 2402–2412, 2015.
- 917 [46] Y. Wu, J. Kokkonen, C. Han, and M. Juntti, “Interference and coverage analysis for Tera-
918 hertz networks with indoor blockage effects and line-of-sight access point association,” *IEEE*
919 *Transactions on Wireless Communications*, vol. 20, no. 3, pp. 1472–1486, 2021.
- 920 [47] J. M. Jornet and I. F. Akyildiz, “Channel modeling and capacity analysis for electromagnetic
921 wireless nanonetworks in the Terahertz band,” *IEEE Transactions on Wireless Communica-*
922 *tions*, vol. 10, no. 10, pp. 3211–3221, 2011.
- 923 [48] B. Wang, F. Gao, S. Jin, H. Lin, and G. Y. Li, “Spatial- and frequency-wideband effects in
924 millimeter-Wave massive MIMO systems,” *IEEE Transactions on Signal Processing*, vol. 66,
925 no. 13, pp. 3393–3406, 2018.
- 926 [49] C.-K. Wen, W.-T. Shih, and S. Jin, “Deep learning for massive MIMO CSI feedback,” *IEEE*
927 *Wireless Communications Letters*, vol. 7, no. 5, pp. 748–751, 2018.

- 928 [50] A. Vaswani *et al.*, “Attention is all you need,” in *Advances in Neural Information Process-*
929 *ing Systems*, I. Guyon *et al.*, Eds., vol. 30, Curran Associates, Inc., 2017. [Online]. Available:
930 [https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-](https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf)
931 [Paper.pdf](https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- 932 [51] Y. Wang, Z. Gao, D. Zheng, S. Chen, D. Gunduz, and H. V. Poor, “Transformer-empowered
933 6g intelligent networks: From massive MIMO processing to semantic communication,” *IEEE*
934 *Wireless Communications*, pp. 1–9, 2022.
- 935 [52] I. O. Tolstikhin *et al.*, “MLP-mixer: An all-MLP architecture for vision,” *Advances in Neural*
936 *Information Processing Systems*, vol. 34, pp. 24 261–24 272, 2021.
- 937 [53] Q. Shi, M. Razaviyayn, Z.-Q. Luo, and C. He, “An iteratively weighted MMSE approach
938 to distributed sum-utility maximization for a MIMO interfering broadcast channel,” *IEEE*
939 *Transactions on Signal Processing*, vol. 59, no. 9, pp. 4331–4340, 2011.
- 940 [54] R. J. Williams and D. Zipser, “A learning algorithm for continually running fully recurrent
941 neural networks,” *Neural Computation*, vol. 1, no. 2, pp. 270–280, 1989.
- 942 [55] C.-R. Tsai, Y.-H. Liu, and A.-Y. Wu, “Efficient compressive channel estimation for millimeter-
943 Wave large-scale antenna systems,” *IEEE Transactions on Signal Processing*, vol. 66, no. 9,
944 pp. 2414–2428, 2018.
- 945 [56] M. Ke, Z. Gao, Y. Wu, X. Gao, and R. Schober, “Compressive sensing-based adaptive active
946 user detection and channel estimation: Massive access meets massive MIMO,” *IEEE Transac-*
947 *tions on Signal Processing*, vol. 68, pp. 764–779, 2020.
- 948 [57] X. Ma, Z. Gao, F. Gao, and M. Di Renzo, “Model-driven deep learning based channel esti-
949 mation and feedback for millimeter-Wave massive hybrid MIMO systems,” *IEEE Journal on*
950 *Selected Areas in Communications*, vol. 39, no. 8, pp. 2388–2406, 2021.