# Deep Multi-Stage CSI Acquisition for Reconfigurable Intelligent Surface Aided MIMO Systems

Shen Gao, Peihao Dong, *Member, IEEE*, Zhiwen Pan, *Member, IEEE*, and Geoffrey Ye Li, *Fellow, IEEE*

*Abstract*—This article aims to reduce huge pilot overhead when estimating the reconfigurable intelligent surface (RIS) relayed wireless channel. Motivated by the compelling grasp of deep learning in tackling nonlinear mapping problems, the proposed approach only activates a part of RIS elements and utilizes the corresponding cascaded channel estimate to predict another part. Through a synthetic deep neural network (DNN), the direct channel and active cascaded channel are first estimated sequentially, followed by the channel prediction for the inactive RIS elements. A three-stage training strategy is developed for this synthetic DNN. From simulation results, the proposed deep learning based approach is effective in reducing the pilot overhead and guaranteeing the reliable estimation accuracy.

*Index Terms*—RIS-aided MIMO systems, CSI acquisition, deep neural network, multi-stage training.

## I. INTRODUCTION

**T**HE recent advent of reconfigurable intelligent surface (RIS) has stirred up a plethora of research activities since it has potential to boost the network performance and reduce the cost [1]–[3]. The general RIS consists of an inexpensive smart surface, usually made of either tiny antenna elements or metamaterials, and some low power circuits. Through real-time reflecting adaption by an external controller, RIS bears the ability of manipulating the phase and amplitude of the impinging signal in order to focus the signal energy at the receiver as well as to mitigate the interference and security threats [4]–[6].

The above-mentioned brightening advantages of RIS build on the reliable channel state information (CSI). However, the CSI acquisition is challenging since the RIS-relayed link is a cascaded channel whose structure differs from the single-hop channel, which makes most channel estimation approaches exploiting single-hop channel statistics malfunction. In addition, the huge pilot overhead incurred by the RIS is a vital limit thwarting the dramatic performance improvement. In [7], a simple on/off operation mode for the RIS is proposed to sequentially estimate the cascaded channel associated with each active RIS element. By exploiting the array gain of the RIS, an efficient least-square (LS) based approach is developed in [8] and [9], which adapts the reflection coefficient of each RIS element as per a discrete Fourier transformation (DFT) matrix to improve the estimation accuracy. In [10],

S. Gao and Z. Pan are with the National Mobile Communications Research Laboratory, Southeast University, Nanjing 210096, China, and also with the Purple Mountain Laboratories, Nanjing 211100, China (e-mail: gaoshen@seu.edu.cn; pzw@seu.edu.cn).

P. Dong is with the College of Electronic and Information Engineering, Nanjing University of Aeronautics and Astronautics, Nanjing 211106, China (e-mail: phdong@nuaa.edu.cn).

G. Y. Li is with Department of Electrical and Electronic Engineering, Imperial College London, London SW7 2AZ, U.K. (e-mail: Geoffrey.Li@imperial.ac.uk).

compressed sensing (CS) is exploited to separately estimate the sparse channel of each hop by endowing the RIS with signal processing ability. In [11], a novel sparse form of the cascaded channel is uncovered to enable the sparse channel recovery with the reduced pilot overhead provided that the sparsity of each hop is known.

As one of the key technologies underlaying the pathway to smart radio, deep learning (DL) based signal processing has sparked a revolution in wireless communications [12]–[15], and thus inspires some attempts to address the channel estimation problem for RIS-aided systems. In [16] and [17], deep convolutional neural network (CNN) is applied to refine the LS channel estimate to further improve the accuracy. A complex-valued denoising CNN is proposed in [18] to enhance the CS-based estimate for the broadband user-RIS channel.

It can be seen that the CS-based channel estimation approaches highly rely on channel statistics as the prior knowledge and will suffer from performance degradation in the practical complicated scenarios. Although the LS-based approaches get rid of this dependence, they still have not well traded off between the estimation accuracy and the pilot overhead. The success of DL in estimating the single-hop channels [14], [15] has implied its potential to address the above-mentioned problems and thus enlightens us to figure out a solution along with this line. The main novelty and contribution of this article can be summarized as follows:

1) With partially activated RIS elements to reduce the pilot overhead, we propose a three-stage CSI acquisition framework successively including estimation of the direct channel, estimation of the cascaded channel for active RIS elements, and prediction of the cascaded channel for inactive RIS elements, and develop a synthetic deep neural network (DNN) to realize it.

2) For the prediction stage, we discover a simple yet efficient mapping relationship from the perspective of each base station (BS) antenna, which facilitates the DNN training under the accumulative estimation errors propagated from the previous stages and thus achieves the superior estimation accuracy.

*Notations*: In this article, we use upper and lower case boldface letters to denote matrices and vectors, respectively. $(\cdot)^T$, $(\cdot)^H$, $\|\cdot\|_F$, and $\mathbb{E}\{\cdot\}$ represent the transpose, conjugate transpose, Frobenius norm, and expectation, respectively. $\|\cdot\|$ denotes the Euclidian norm of a vector. $\mathrm{diag}(\mathbf{x})$ transforms vector $\mathbf{x}$ to a diagonal matrix. $\mathbf{0}_N$ denotes an all-zero column vector with $N$ entries. $\mathbf{1}_N$ denotes a column vector with each of $N$ entries equal to 1. $\mathbf{I}_N$ denotes an $N \times N$ identity matrix. $\mathcal{CN}(0, \sigma^2)$ represents a circular symmetric complex Gaussian distribution with variance $\sigma^2$. $|\mathcal{X}|$ denotes the cardinality of set $\mathcal{X}$.
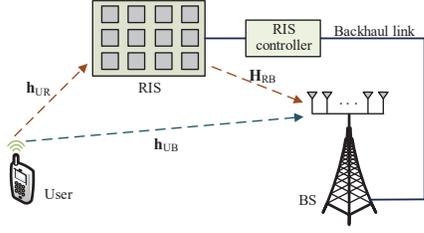
Fig. 1. A RIS-aided MIMO uplink system.

## II. System Model

We consider a RIS-aided multiple-input multiple-output (MIMO) uplink system shown in Fig. 1, where a single-antenna user transmits signals to a BS with $M$ antennas and a RIS with $N$ reflecting elements is deployed to enhance the communication quality by using extremely low power and cost. The reflecting coefficient of each RIS element can be adapted by the BS via a controller based on the real-time CSI in order to smartly rotate the phase of the incident signal. In Fig. 1, the MIMO system is operated in the time division duplex mode, where the BS acquires the CSI resorting to the pilot signals transmitted by the user.

At the $i$th time instant of a coherence interval, the received pilot signal at the BS is expressed as

$$\mathbf{y}_i = \sqrt{P}(\mathbf{h}_{\mathrm{UB}} + \mathbf{H}_{\mathrm{RB}}\mathrm{diag}(\boldsymbol{\phi}_i)\mathbf{h}_{\mathrm{UR}})x_i + \mathbf{z}_i, \qquad (1)$$

where $P$, $x_i$, $\boldsymbol{\phi}_i$, and $\mathbf{z}_i \sim \mathcal{CN}(0, \sigma_0^2\mathbf{I}_M)$ denote the user transmit power, the transmitted pilot signal, the reflecting coefficient vector of the RIS, and the additive white Gaussian noise (AWGN) at the BS, respectively. From [8], the $n$th entry of $\boldsymbol{\phi}_i \in \mathbb{C}^{N\times 1}$ corresponds to the reflecting coefficient of the $n$th RIS element and is written as $\phi_{ni} = \beta_{ni}e^{j\varphi_{ni}}$ with $\beta_{ni} \in [0, 1]$ and $\varphi_{ni} \in (0, 2\pi]$ accounting for the reflecting amplitude and phase rotation, respectively. In this article, we set $\beta_{ni} = 1$ for the active RIS elements to avoid the energy loss and simplify the RIS hardware structure. $\mathbf{h}_{\mathrm{UB}} \in \mathbb{C}^{M\times 1}$, $\mathbf{h}_{\mathrm{UR}} \in \mathbb{C}^{N\times 1}$, and $\mathbf{H}_{\mathrm{RB}} \in \mathbb{C}^{M\times N}$ denote the channel from the user to the BS, the channel from the user to the RIS, and the channel from the RIS to the BS, respectively.

To facilitate further signal processing, $\mathbf{y}_i$ is rewritten as

$$\begin{aligned}\mathbf{y}_i &= \sqrt{P}(\mathbf{h}_{\mathrm{UB}} + \mathbf{H}_{\mathrm{RB}}\mathrm{diag}(\mathbf{h}_{\mathrm{UR}})\boldsymbol{\phi}_i)x_i + \mathbf{z}_i \\ &\triangleq \sqrt{P}(\mathbf{h}_{\mathrm{UB}} + \mathbf{G}\boldsymbol{\phi}_i)x_i + \mathbf{z}_i.\end{aligned} \qquad (2)$$

Note that it is more tractable to estimate the equivalent cascaded channel $\mathbf{G}$ instead of estimating $\mathbf{H}_{\mathrm{RB}}$ and $\mathbf{h}_{\mathrm{UR}}$ separately and $\mathbf{G}$ can be directly used for the beamforming design during the data transmission. Hence, we will focus on the acquisition of $\mathbf{h}_{\mathrm{UB}}$ and $\mathbf{G}$ hereinafter.

Then the received pilot signals at the BS during $\tau + 1$ time instants are given in matrix form as

$$\mathbf{Y} = \sqrt{P}(\mathbf{h}_{\mathrm{UB}}\mathbf{1}_{\tau+1}^T + \mathbf{G}\boldsymbol{\Psi})\mathbf{X} + \mathbf{Z}, \qquad (3)$$

where $\boldsymbol{\Psi} = [\boldsymbol{\phi}_1, \ldots, \boldsymbol{\phi}_{\tau+1}]$ with $\boldsymbol{\phi}_1 = \mathbf{0}_N$, $\mathbf{X} = \mathrm{diag}([x_1, \ldots, x_{\tau+1}])$, and $\mathbf{Z} = [\mathbf{z}_1, \ldots, \mathbf{z}_{\tau+1}]$. Without loss of generality, we assume $\mathbf{X} = \mathbf{I}_{\tau+1}$ for simplicity, which yields

$$\mathbf{Y} = \sqrt{P}(\mathbf{h}_{\mathrm{UB}}\mathbf{1}_{\tau+1}^T + \mathbf{G}\boldsymbol{\Psi}) + \mathbf{Z}. \qquad (4)$$

Based on the general pilot transmission model sketched above, we investigate how to acquire reliable $\mathbf{h}_{\mathrm{UB}}$ and $\mathbf{G}$ with low pilot overhead via DL in the following.

## III. DL-Based Three-Stage CSI Acquisition

In this section, a novel three-stage CSI acquisition framework is proposed based on DL. We will first shed light on the basic idea of the framework and then design a synthetic DNN as its backbone for CSI acquisition. Finally, the online testing procedure of the framework is described.

### A. Basic Idea

In the proposed framework, $\mathbf{h}_{\mathrm{UB}}$ is first estimated via a DNN while all RIS elements are turned off. After then, a part of RIS elements are activated with their reflection coefficients adapted as per the DFT matrix [8], [9] so that the BS estimates the corresponding equivalent cascaded channel by using another DNN[1]. Finally, the DNN-based channel prediction is conducted to retrieve the equivalent cascaded channel associated with those inactive RIS elements. These three DNNs hook up in a sequential manner and compose a synthetic DNN. To achieve a good overall estimation accuracy, we design and train the DNNs separately in three stages[2].

### B. Three-Stage DNN Design

At the beginning of this subsection, we outright present the overall structure of the proposed CSI acquisition framework in Fig. 2 to facilitate the elaboration on the synthetic DNN design in each stage.

*1) Stage 1:* At the first time instant during pilot transmission, all RIS elements are turned off so that the BS can estimate $\mathbf{h}_{\mathrm{UB}}$ from the received pilot signal, which is the first column of $\mathbf{Y}$ in (4) and is written as

$$\mathbf{y}_1 = \sqrt{P}\mathbf{h}_{\mathrm{UB}} + \mathbf{z}_1. \qquad (5)$$

Then the LS estimate of $\mathbf{h}_{\mathrm{UB}}$ can be obtained as $\hat{\mathbf{h}}_{\mathrm{UB,ls}} = \frac{\mathbf{y}_1}{\sqrt{P}}$.

As shown in Fig. 2, $\hat{\mathbf{h}}_{\mathrm{UB,ls}}$ is then input into a **D**irect Channel **E**stimation DNN (DE-DNN) in attempts to approximate the true channel $\mathbf{h}_{\mathrm{UB}}$. Thus DE-DNN is trained with the sample tuple, $\langle\hat{\mathbf{h}}_{\mathrm{UB,ls}}, \mathbf{h}_{\mathrm{UB}}\rangle$, and the loss function,

$$\mathcal{L}_{\mathrm{DE}} = \frac{1}{N_{\mathrm{tr}}}\sum_{n=1}^{N_{\mathrm{tr}}}\|\mathbf{h}_{\mathrm{UB}}^{(n)} - \hat{\mathbf{h}}_{\mathrm{UB,dnn}}^{(n)}\|^2, \qquad (6)$$

where $N_{\mathrm{tr}}$ denotes the number of training samples, $\hat{\mathbf{h}}_{\mathrm{UB,dnn}}^{(n)}$ denotes the channel approximated by DE-DNN, and the superscript $(n)$ indicates the $n$th sample. In more detail, DE-DNN is fully-connected (FC) and includes three hidden layers, each of which applies rectified linear unit (ReLU) activation function and batch normalization (BN) to avoid gradient vanishing and overfitting. The output layer does not apply any activation function so that the label values need not to be tailored to fit the activation function.

*2) Stage 2:* In this stage, only a part of RIS elements, denoted by set $\mathcal{A}$ with $N_1 = |\mathcal{A}|$, are activated to reflect the pilot signals transmitted by the user to the BS. Considering the efficient DFT matrix based reflection mode for the active

---

[1]By turning on only a part of RIS elements, relatively accurate LS channel estimate for these active RIS elements can be obtained even with reduced pilot overhead, which can be further exploited by the elaborated DNNs to construct the complete channel.

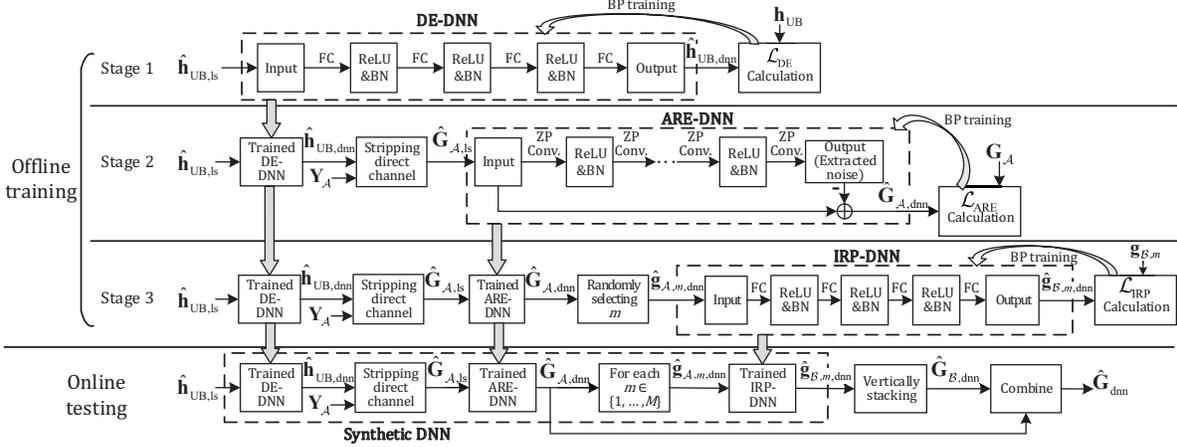[2]The DNNs are trained offline and thus the computational cost of training is relatively trivial.

Fig. 2. DL-based three-stage CSI acquisition framework.

elements [8], [9], the received pilot signals at the BS from the second to the $(N_1 + 1)$th time instant can be expressed as

$$\mathbf{Y}_{\mathcal{A}} = \sqrt{P}(\mathbf{h}_{\text{UB}}\mathbf{1}_{N_1}^T + \mathbf{G}_{\mathcal{A}}\boldsymbol{\Phi}_{N_1}) + \mathbf{Z}_{\mathcal{A}}, \qquad (7)$$

where $\mathbf{G}_{\mathcal{A}} \in \mathbb{C}^{M \times N_1}$ denotes the equivalent cascaded channel associated with the active RIS elements whose columns are fetched from $\mathbf{G}$ according to the indexes in $\mathcal{A}$, the DFT matrix, $\boldsymbol{\Phi}_{N_1} \in \mathbb{C}^{N_1 \times N_1}$, indicates the reflection coefficients, and $\mathbf{Z}_{\mathcal{A}}$ is the corresponding AWGN.

As $\mathbf{h}_{\text{UB}}$ has been estimated by the DE-DNN, its estimate $\hat{\mathbf{h}}_{\text{UB,dnn}}$ will be subtracted from $\mathbf{Y}_{\mathcal{A}}$, which yields the LS estimate of $\mathbf{G}_{\mathcal{A}}$, i.e.,

$$\hat{\mathbf{G}}_{\mathcal{A},\text{ls}} = \left(\frac{\mathbf{Y}_{\mathcal{A}}}{\sqrt{P}} - \hat{\mathbf{h}}_{\text{UB,dnn}}\mathbf{1}_{N_1}^T\right)\frac{\boldsymbol{\Phi}_{N_1}^H}{N_1}$$
$$= \mathbf{G}_{\mathcal{A}} + \left((\mathbf{h}_{\text{UB}} - \hat{\mathbf{h}}_{\text{UB,dnn}})\mathbf{1}_{N_1}^T + \frac{\mathbf{Z}_{\mathcal{A}}}{\sqrt{P}}\right)\frac{\boldsymbol{\Phi}_{N_1}^H}{N_1}. \quad (8)$$

From Fig. 2, $\hat{\mathbf{G}}_{\mathcal{A},\text{ls}}$ will be refined by an **A**ctive **R**IS Channel **E**stimation DNN (ARE-DNN) to output a version closer to the true channel, $\mathbf{G}_{\mathcal{A}}$. Similarly, the loss function for ARE-DNN training is given by

$$\mathcal{L}_{\text{ARE}} = \frac{1}{N_{\text{tr}}}\sum_{n=1}^{N_{\text{tr}}}\|\mathbf{G}_{\mathcal{A}}^{(n)} - \hat{\mathbf{G}}_{\mathcal{A},\text{dnn}}^{(n)}\|_F^2, \qquad (9)$$

where $\hat{\mathbf{G}}_{\mathcal{A},\text{dnn}}^{(n)}$ denotes the output of ARE-DNN.

Since $\hat{\mathbf{G}}_{\mathcal{A},\text{ls}}$ is corrupted by the residual estimation error of Stage 1 in addition to AWGN, we invoke the more efficient residual network structure to design ARE-DNN. As shown in Fig. 2, the input, $\hat{\mathbf{G}}_{\mathcal{A},\text{ls}}$, flows in ARE-DNN through two ways, one of which includes several convolutional layers to successively distill the aggregated noise from $\hat{\mathbf{G}}_{\mathcal{A},\text{ls}}$ while another one is a shortcut representing the identity mapping. The two ways intersect by subtracting the extracted noise from the input, $\hat{\mathbf{G}}_{\mathcal{A},\text{ls}}$, and then a more purified channel, $\hat{\mathbf{G}}_{\mathcal{A},\text{dnn}}$, can be obtained. Specifically, the first way consists of eight zero padding (ZP) convolutional layers for feature extraction [14]. Each of the first seven layers applies 64 $3 \times 3$ kernels, ReLU activation function and BN while the last layer only applies 2 $3 \times 3$ kernels and directly output the filtered results. Simulation trails show that further increasing the number of convolutional layers will not improve the performance. In addition, using

the same or even a bit larger number of layers, stacking multiple residual network units cannot beat the current one-unit structure, which indicates that the uninterrupted layer structure is more efficient to extract the high-order features in this case.

*3) Stage 3:* Denote $\mathcal{B}$ as the index set of the inactive RIS elements with $N_2 = |\mathcal{B}|$. The aim of this stage is to infer the equivalent cascaded channel associated with $\mathcal{B}$, $\mathbf{G}_{\mathcal{B}}$, from $\hat{\mathbf{G}}_{\mathcal{A},\text{dnn}}$. A straightforward way to carry out this task is inputting $\hat{\mathbf{G}}_{\mathcal{A},\text{dnn}}$ into a CNN to approximate $\mathbf{G}_{\mathcal{B}}$. However, it is challenging to uncover this matrix mapping relationship since the structure of the equivalent cascaded channel is quite different from that of the single-hop channel. This inspires us to dissect the channel structure and find the more suitable input that the neural network can digest well.

Since $\mathbf{G} = \mathbf{H}_{\text{RB}}\text{diag}(\mathbf{h}_{\text{UR}})$, we start with analyzing the structure of $\mathbf{H}_{\text{RB}}$. According to Saleh-Valenzuela (SV) model, $\mathbf{H}_{\text{RB}}$ is given by

$$\mathbf{H}_{\text{RB}} = \sqrt{\frac{MN}{L_{\text{RB}}}}\sum_{l=1}^{L_{\text{RB}}}\alpha_{\text{RB},l}\mathbf{a}_{\text{B}}(\theta_l)\mathbf{a}_{\text{R}}^H(\varphi_l), \qquad (10)$$

where $L_{\text{RB}}$, $\alpha_{\text{RB},l}$, $\theta_l$, and $\varphi_l$ denote number of main paths, the complex gain of the $l$th path, the azimuth angles of arrival and departure (AoA/AoD) at the BS and the RIS, respectively. The response vector $\mathbf{a}_{\text{B}}(\theta_l)$ can be further expressed as $\mathbf{a}_{\text{B}}(\theta_l) = \frac{1}{\sqrt{M}}\left[1, e^{-j2\pi\frac{d}{\lambda}\sin(\theta_l)}, \ldots, e^{-j2\pi\frac{d}{\lambda}(M-1)\sin(\theta_l)}\right]^T$ with $d$ and $\lambda$ denoting the space between the adjacent antennas at the BS and the wavelength of the carrier frequency, respectively. Then the $m$th row of $\mathbf{H}_{\text{RB}}$ can be written as

$$\mathbf{h}_{\text{RB},m} = \sqrt{\frac{N}{L_{\text{RB}}}}\sum_{l=1}^{L_{\text{RB}}}\alpha_{\text{RB},l}e^{-j2\pi\frac{d}{\lambda}(m-1)\sin(\theta_l)}\mathbf{a}_{\text{R}}^H(\varphi_l). (11)$$

It is obvious that each row of $\mathbf{H}_{\text{RB}}$ exhibits a unified form containing all the channel information between the RIS and the BS. Consequently, the $m$th row of $\mathbf{G}$, which is given by $\mathbf{g}_m = \mathbf{h}_{\text{RB},m}\text{diag}(\mathbf{h}_{\text{UR}})$, $\forall m \in \{1, \ldots, M\}$, inherits this property with all the channel information of the RIS-relayed channel included.

Denote $\mathbf{g}_{\mathcal{A},m}$ and $\mathbf{g}_{\mathcal{B},m}$ as the parts of $\mathbf{g}_m$ corresponding to $\mathcal{A}$ and $\mathcal{B}$, respectively. Based on the perspective view of the

channel mentioned above, we can focus on $\mathbf{g}_{\mathcal{A},m}$ and $\mathbf{g}_{\mathcal{B},m}$, instead of $\mathbf{G}_{\mathcal{A}}$ and $\mathbf{G}_{\mathcal{B}}$, to extract the inherence underlaying the mapping relationship, which facilitates the DNN design and training with significantly improved prediction accuracy. Specifically, as shown in Fig. 2, we design an **I**nactive **R**IS Channel **P**rediction DNN (IRP-DNN) to approximate $\mathbf{g}_{\mathcal{B},m}$ by using $\hat{\mathbf{g}}_{\mathcal{A},m,\mathrm{dnn}}$, i.e., the $m$th row of $\hat{\mathbf{G}}_{\mathcal{A},\mathrm{dnn}}$, as the input with $m$ randomly selected from $\{1,\ldots,M\}$. The loss function for IRP-DNN training is expressed as

$$\mathcal{L}_{\mathrm{IRP}} = \frac{1}{N_{\mathrm{tr}}} \sum_{n=1}^{N_{\mathrm{tr}}} \|\mathbf{g}_{\mathcal{B},m}^{(n)} - \hat{\mathbf{g}}_{\mathcal{B},m,\mathrm{dnn}}^{(n)}\|^2, \qquad (12)$$

where $\hat{\mathbf{g}}_{\mathcal{B},m,\mathrm{dnn}}^{(n)}$ denotes the output of IRP-DNN. Since we focus on the level of vector mapping instead of matrix mapping, IRP-DNN is designed in a FC structure. There are three hidden layers applying ReLU activation function and BN while no activation function is applied in the output layer.

### C. Online Testing

After offline training, the CSI acquisition framework will be deployed for online testing. From Fig. 2, $\hat{\mathbf{h}}_{\mathrm{UB,ls}}$, which is obtained from $\mathbf{y}_1$, is first refined by DE-DNN to output $\hat{\mathbf{h}}_{\mathrm{UB,dnn}}$, based on which the component of the direct channel will be stripped from $\mathbf{Y}_{\mathcal{A}}$. Then $\hat{\mathbf{G}}_{\mathcal{A},\mathrm{ls}}$ is obtained and will be further purified via ARE-DNN to yield $\hat{\mathbf{G}}_{\mathcal{A},\mathrm{dnn}}$. Afterward, each row of $\hat{\mathbf{G}}_{\mathcal{A},\mathrm{dnn}}$ is input into IRP-DNN sequentially to predict the corresponding $\hat{\mathbf{g}}_{\mathcal{B},m,\mathrm{dnn}}$. Vertically stacking $\hat{\mathbf{g}}_{\mathcal{B},m,\mathrm{dnn}}$ ($\forall m \in \{1,\ldots,M\}$) constructs $\hat{\mathbf{G}}_{\mathcal{B},\mathrm{dnn}}$ and combining $\hat{\mathbf{G}}_{\mathcal{A},\mathrm{dnn}}$ and $\hat{\mathbf{G}}_{\mathcal{B},\mathrm{dnn}}$ finally obtains $\hat{\mathbf{G}}_{\mathrm{dnn}}$. It can be seen that the proposed CSI acquisition framework is able to estimate $\mathbf{h}_{\mathrm{UB}}$ and $\mathbf{G}$ with high accuracy while the pilot overhead can be reduced by the ratio of $\frac{N_1+1}{N+1}$.

In addition, the computational complexities of DE-DNN, ARE-DNN, and IRP-DNN are $\mathcal{O}(\sum_{i=2}^{L_{\mathrm{DE}}} N_{i-1}^{\mathrm{DE}} N_i^{\mathrm{DE}})$, $\mathcal{O}(MN_1^2 + MN_1 \sum_{i=1}^{L_{\mathrm{ARE}}} K_i^2 F_{i-1} F_i)$, and $\mathcal{O}(M \sum_{i=2}^{L_{\mathrm{IRP}}} N_{i-1}^{\mathrm{IRP}} N_i^{\mathrm{IRP}})$, respectively, where $L_{\mathrm{DE}}$, $L_{\mathrm{ARE}}$, and $L_{\mathrm{IRP}}$ denote the numbers of (convolutional) layers of DE-DNN, ARE-DNN, and IRP-DNN, $N_i^{\mathrm{DE}}$ and $N_i^{\mathrm{IRP}}$ denotes the corresponding numbers of neurons of the $i$th layer, $K_i$ is the side length of the filters used by the $i$th convolutional layer, $F_{i-1}$ and $F_i$ denote the numbers of input and output feature maps of the $i$th convolutional layer.

## IV. Simulation Results

In this section, numerical results are presented to validate the proposed DL-based CSI acquisition framework. The baseline schemes for comparison include the LS estimator [8], [9], orthogonal matching pursuit (OMP) [11] and ChannelNet [16]. The numbers of BS antennas and RIS elements are set as $M = 16$ and $N = 128$, respectively. The numbers of paths of $\mathbf{h}_{\mathrm{UB}}$, $\mathbf{h}_{\mathrm{UR}}$, and $\mathbf{H}_{\mathrm{RB}}$ are set as 3. The average noise power, $\sigma_0^2$, is normalized to 1 and the transmit power, $P$, is set as a relative value with respect to $\sigma_0^2$. The signal-to-noise ratio (SNR) in Figs. 3 and 4 is defined as $\mathrm{SNR} = 10 \lg \frac{P}{\sigma_0^2}$ (dB). Without loss of generality, assume that the pathloss has been absorbed into the transmit power for simplicity. Then the average power gain of each path of $\mathbf{h}_{\mathrm{UB}}$, $\mathbf{h}_{\mathrm{UR}}$, and $\mathbf{H}_{\mathrm{RB}}$ is set as 1 [19]. For the proposed framework, the training, validation, and testing sets

TABLE I
STRUCTURES OF DE-DNN, ARE-DNN, AND IRP-DNN

| | Layer type | Tensor size | Kernel size | Activation function |
|---|---|---|---|---|
| DE-DNN | Input | $2M$ | - | - |
| | Dense | 64 | - | ReLU |
| | Dense | 128 | - | ReLU |
| | Dense | 64 | - | ReLU |
| | Output | $2M$ | - | - |
| ARE-DNN | Input | $M \times N_1 \times 2$ | - | - |
| | ZP Conv. (7 layers) | $M \times N_1 \times 64$ | $3 \times 3$ | ReLU |
| | Output | $M \times N_1 \times 2$ | $3 \times 3$ | - |
| IRP-DNN | Input | $2N_1$ | - | - |
| | Dense | 128 | - | ReLU |
| | Dense | 256 | - | ReLU |
| | Dense | 256 | - | ReLU |
| | Output | $2N_2$ | - | - |

contain $90,000$, $10,000$, and $10,000$ samples, respectively. Adam is applied as the optimizer and the batch size is set as 128. The training of each DNN therein lasts 300 epochs with the initial learning rates $1 \times 10^{-3}$ and $1 \times 10^{-4}$ for the first 200 epochs and the remaining 100 epochs, respectively. The detailed DNN structures are listed in Table I. For offline training, the label of DE-DNN, $\mathbf{h}_{\mathrm{UB}}$, is generated according to the SV channel model and the input, $\hat{\mathbf{h}}_{\mathrm{UB,ls}}$, is generated by $\hat{\mathbf{h}}_{\mathrm{UB,ls}} = \frac{\mathbf{y}_1}{\sqrt{P}}$. Then the label of ARE-DNN, $\mathbf{G}_{\mathcal{A}}$, is obtained by fetching the corresponding columns from $\mathbf{G}$ as per the indexes in $\mathcal{A}$. The input, $\hat{\mathbf{G}}_{\mathcal{A},\mathrm{ls}}$, is calculated as (8), where $\mathbf{h}_{\mathrm{UB}}$ is the training label of DE-DNN and $\hat{\mathbf{h}}_{\mathrm{UB,dnn}}$ is the corresponding output. The label of IRP-DNN, $\mathbf{g}_{\mathcal{B},m}$, is the $m$th row of $\mathbf{G}_{\mathcal{B}}$, where $\mathbf{G}_{\mathcal{B}}$ is the complementary matrix of the training label, $\mathbf{G}_{\mathcal{A}}$, and $m$ is arbitrarily selected from $\{1,\ldots,M\}$. The input, $\hat{\mathbf{g}}_{\mathcal{A},m,\mathrm{dnn}}$, is the $m$th row of $\hat{\mathbf{G}}_{\mathcal{A},\mathrm{dnn}}$, which is the output of ARE-DNN when the input is $\hat{\mathbf{G}}_{\mathcal{A},\mathrm{ls}}$. The testing samples are generated similarly to training samples, but different path gains and AoAs/AoDs in the SV model are used. The normalized mean-squared error (NMSE) is used to evaluate the CSI acquisition performance and can be expressed as $\mathsf{NMSE} = \mathbb{E}\left\{\frac{\|\mathbf{h}_{\mathrm{UB}} - \hat{\mathbf{h}}_{\mathrm{UB,dnn}}\|^2}{\|\mathbf{h}_{\mathrm{UB}}\|^2}\right\}$ for $\mathbf{h}_{\mathrm{UB}}$ and $\mathsf{NMSE} = \mathbb{E}\left\{\frac{\|\mathbf{G} - \hat{\mathbf{G}}_{\mathrm{dnn}}\|_F^2}{\|\mathbf{G}\|_F^2}\right\}$ for $\mathbf{G}$. The pilot overhead ratio for estimating $\mathbf{G}$ is defined as $r = \frac{N_1}{N}$.

Fig. 3 plots the NMSE performance versus SNR for $\mathbf{h}_{\mathrm{UB}}$. From Fig. 3, both the ChannelNet and the proposed DE-DNN are effective to refine the LS estimate and further improve the accuracy. Compared with the ChannelNet, DE-DNN achieves almost same or even better performance, indicating that the DNN with FC structure is adequate to estimate the direct channel with high accuracy.

The estimation performance of $\mathbf{G}$ versus SNR is shown in Fig. 4 with $r = \frac{1}{4}$. The LS estimator keeps error floor over the considered SNR regime because the performance will be poor regardless of the SNR level once the pilot overhead $N_1$ is less than $N$. The ChannelNet is based on the LS estimation and thus also performs unsatisfactorily. The OMP approach outperforms the LS estimator and ChannelNet but cannot provide very accurate CSI due to the complication of the cascaded RIS channel. By contrast, the proposed framework divides the acquisition of $\mathbf{G}$ into estimation and prediction
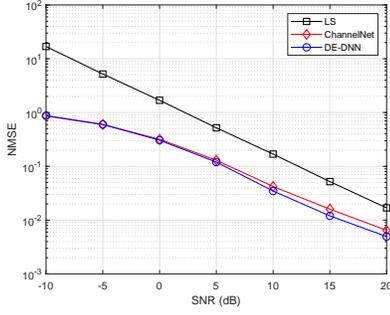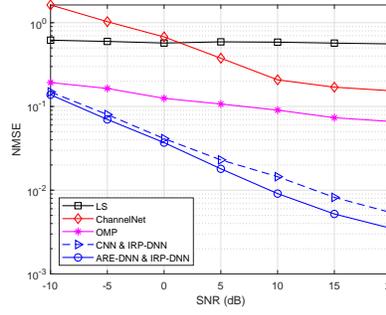
Fig. 3. NMSE versus SNR for the direct channel $\mathbf{h}_{\mathrm{UB}}$.

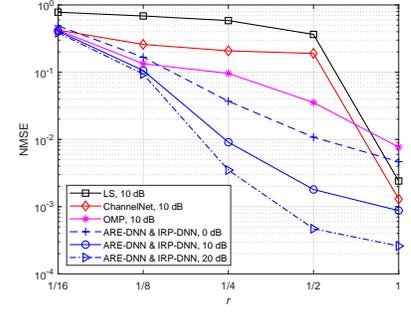Fig. 4. NMSE versus SNR for the equivalent cascaded channel $\mathbf{G}$.

Fig. 5. NMSE versus the pilot overhead ratio for $\mathbf{G}$.

stages with respective dedicated DNNs, i.e., ARE-DNN and IRP-DNN, which can learn the features of the RIS-relayed channel more comprehensively and thus achieve the superior performance. In addition, the performance will be a little worse if the residual network structure of ARE-DNN is replaced by a CNN with the same number of layers, which indicates the effectiveness of the residual network structure.

In Fig. 5, we further investigate the NMSE performance of $\mathbf{G}$ versus the pilot overhead ratio, $r$. The proposed framework (ARE-DNN and IRP-DNN) always outperforms the LS estimator, ChannelNet, and OMP at different values of $r$, which demonstrates that the proposed framework is more robust to the reduction of the pilot overhead. For the proposed framework, the performance gap between different SNRs increases with $r$ since the dominating factor turns into SNR from the pilot overhead.

Considering the time complexity for practical implementation, the runtimes of the proposed scheme and the ChannelNet are $1.37 \times 10^{-4}$ seconds and $2.75 \times 10^{-4}$ seconds, respectively, on the GTX 2080Ti GPU while the runtime of the OMP is $5.87 \times 10^{-3}$ seconds on the Intel(R) Core(TM) i7-3770 CPU. The proposed scheme consumes minimum runtime owing to elaborated design and efficient parallel computing.

## V. CONCLUSION

In this article, we develop a three-stage CSI acquisition framework for the RIS-aided MIMO uplink system based on an elaborated synthetic DNN. It includes three dedicated DNNs in charge of estimating the direct channel, estimating the cascaded channel for active RIS elements, and predicting the cascaded channel for inactive RIS elements, respectively. The three DNNs with specialized structures are trained and hook up sequentially. Simulation results show that the proposed CSI acquisition framework can achieve superior performance without relying on high pilot overhead and exact knowledge on channel statistics. In future work, the proposed framework can be extended to the broadband channel exploiting frequency correlation or common sparsity [14], [18].

## REFERENCES

[1] Q. Wu and R. Zhang, "Towards smart and reconfigurable environment: Intelligent reflecting surface aided wireless network," *IEEE Commun. Mag.*, vol. 58, no. 1, pp. 106–112, Jan. 2020.

[2] C. Huang, A. Zappone, G. C. Alexandropoulos, M. Debbah, and C. Yuen, "Reconfigurable intelligent surfaces for energy efficiency in wireless communication," *IEEE Trans. Wireless Commun.*, vol. 18, no. 8, pp. 4157–4170, Aug. 2019.

[3] C. Huang et al., "Holographic MIMO surfaces for 6G wireless networks: Opportunities, challenges, and trends," *IEEE Wireless Commun.*, vol. 27, no. 5, pp. 118–125, Oct. 2020.

[4] M. Di Renzo *et al.*, "Smart radio environments empowered by reconfigurable intelligent surfaces: How it works, state of research, and the road ahead," *IEEE J. Sel. Areas Commun.*, vol. 38, no. 11, pp. 2450–2525, Nov. 2020.

[5] C. Huang, R. Mo and C. Yuen, "Reconfigurable intelligent surface assisted multiuser MISO systems exploiting deep reinforcement learning," *IEEE J. Sel. Areas Commun.*, vol. 38, no. 8, pp. 1839–1850, Nov. 2020.

[6] Z. Yang, W. Xu, C. Huang, J. Shi, and M. Shikh-Bahaei, "Beamforming design for multiuser transmission through reconfigurable intelligent surface," *IEEE Trans. Commun.*, vol. 69, no. 1, pp. 589–601, Jan. 2021.

[7] D. Mishra and H. Johansson, "Channel estimation and low-complexity beamforming design for passive intelligent surface assisted MISO wireless energy transfer," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, Brighton, U.K., May 2019, pp. 4659–4663.

[8] B. Zheng and R. Zhang, "Intelligent reflecting surface-enhanced OFDM: Channel estimation and reflection optimization," *IEEE Wireless Commun. Lett.*, vol. 9, no. 4, pp. 518–522, Apr. 2020.

[9] T. L. Jensen and E. D. Carvalho, "On optimal channel estimation scheme for intelligent reflecting surfaces based on a minimum variance unbiased estimator," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, Barcelona, Spain, May 2020, pp. 5000–5004.

[10] A. Taha, M. Alrabeiah, and A. Alkhateeb, "Enabling large intelligent surfaces with compressive sensing and deep learning," *arXiv preprint arXiv:1904.10136*, Apr. 2019.

[11] P. Wang, J. Fang, H. Duan, and H. Li, "Compressed channel estimation for intelligent reflecting surface-assisted millimeter wave systems," *IEEE Signal Process. Lett.*, vol. 27, pp. 905–909, 2020.

[12] H. Ye, G. Y. Li, and B.-H. Juang, "Power of deep learning for channel estimation and signal detection in OFDM systems," *IEEE Wireless Commun. Lett.*, vol. 7, no. 1, pp. 114–117, Feb. 2018.

[13] Z.-J. Qin, H. Ye, G. Y. Li, and B.-H. Juang, "Deep learning in physical layer communications," *IEEE Wireless Commun.*, vol. 26, no. 2, pp. 93–99, Apr. 2019.

[14] P. Dong, H. Zhang, G. Y. Li, I. Gaspar, and N. NaderiAlizadeh, "Deep CNN-based channel estimation for mmWave massive MIMO systems," *IEEE J. Sel. Topics Signal Process.*, vol. 13, no. 5, pp. 989–1000, Sep. 2019.

[15] S. Gao, P. Dong, Z. Pan, and G. Y. Li, "Deep learning based channel estimation for massive MIMO with mixed-resolution ADCs," *IEEE Commun. Lett.*, vol. 23, no. 11, pp. 1989–1993, Nov. 2019.

[16] A. M. Elbir, A. Papazafeiropoulos, P. Kourtessis, and S. Chatzinotas, "Deep channel learning for large intelligent surfaces aided mm-Wave massive MIMO systems," *IEEE Wireless Commun. Lett.*, vol. 9, no. 9, pp. 1447–1451, Sep. 2020.

[17] N. K. Kundu and M. R. McKay, "A deep learning-based channel estimation approach for MISO communications with large intelligent surfaces," in *Proc. IEEE Int. Symp. Personal, Indoor and Mobile Radio Communications*, London, United Kingdom, Aug. 2020, pp. 1–6.

[18] S. Liu, Z. Gao, J. Zhang, M. D. Renzo, and M. -S. Alouini, "Deep denoising neural network assisted compressive channel estimation for mmWave intelligent reflecting surfaces," *IEEE Trans. Veh. Technol.*, vol. 69, no. 8, pp. 9223–9228, Aug. 2020.

[19] Z. He and X. Yuan, "Cascaded channel estimation for large intelligent metasurface assisted massive MIMO," *IEEE Wireless Commun. Lett.*, vol. 9, no. 2, pp. 210–214, Feb. 2020.