# Beamspace Channel Estimation for Wideband Millimeter-Wave MIMO: A Model-Driven Unsupervised Learning Approach

Hengtao He, Rui Wang, Weijie Jin, Shi Jin, *Senior Member, IEEE,* Chao-Kai Wen, *Senior Member, IEEE,* and Geoffrey Ye Li, *Fellow, IEEE*

**Abstract**

Millimeter-wave (mmWave) communications have been one of the promising technologies for future wireless networks that integrate a wide range of data-demanding applications. To compensate for the large channel attenuation in mmWave band and avoid high hardware cost, a lens-based beamspace massive multiple-input multiple-output (MIMO) system is considered. However, the beam squint effect in wideband mmWave systems makes channel estimation very challenging, especially when the receiver is equipped with a limited number of radio-frequency (RF) chains. Furthermore, the real channel data cannot be obtained before the mmWave system is used in a new environment, which makes it impossible to train a deep learning (DL)-based channel estimator using real data set beforehand. To solve the problem, we propose a model-driven unsupervised learning network, named learned denoising-based generalized expectation consistent (LDGEC) signal recovery network. By utilizing the Stein's unbiased risk estimator loss, the LDGEC network can be trained only with limited measurements corresponding to the pilot symbols, instead of the real channel data. Even if designed for unsupervised learning, the LDGEC network can be supervisingly trained with the real channel via the denoiser-by-denoiser way. The numerical results demonstrate that the LDGEC-based channel estimator significantly outperforms state-of-the-art compressive sensing-based algorithms when the receiver is equipped with a small number of RF chains and low-resolution ADCs.

H. He, R. Wang, W. Jin and S. Jin are with the National Mobile Communications Research Laboratory, Southeast University, Nanjing 210096, China (e-mail: hehengtao@seu.edu.cn, wang_rui@seu.edu.cn, jinweijie@seu.edu.cn, and jinshi@seu.edu.cn).

C.-K. Wen is with the Institute of Communications Engineering, National Sun Yat-sen University, Kaohsiung 804, Taiwan (e-mail: chaokai.wen@mail.nsysu.edu.tw).

G. Y. Li is with the Department of Electrical and Electronic Engineering, Imperial College London. (email: geoffrey.li@imperial.ac.uk).

**Index Terms**

mmWave communications, deep learning, model-driven, deep unfolding, beam squint, massive MIMO, beamspace channel estimation, unsupervised learning

## I. INTRODUCTION

MmWave communications have been considered as a promising technology to support the very high data rate in future wireless communications since it can provide a tenfold increase in the bandwidth [1]–[3]. However, as the carrier frequency increases, the mmWave signals suffer from much more severe attenuation, which becomes a vital issue in mmWave communications. Leveraging the large antenna arrays employed at the transmitter and receiver, massive multiple-input multiple-output (MIMO) can perform directional beamforming to achieve a high beamforming gain, which helps overcome large pathloss of mmWave signals and guarantees sufficient received signal-to-noise ratio (SNR). However, the hardware cost and power consumption both increase with the number of RF chains, which is sometimes unaffordable if a dedicated RF chain is used for each of a huge number of antennas. To reduce the number of required RF chains, we can resort to beamspace massive MIMO with a discrete lens array (DLA), which has been first proposed in [4] and successfully employed in millimeter-wave (mmWave) communications. However, the number of RF chains is much smaller than that of antennas, and we cannot directly observe the complete channel in the baseband [5], thus incurring challenges for beamspace channel estimation.

### A. Related Work

For beamspace channel estimation, several works utilize compressive sensing (CS) techniques [5]–[9] in mmWave band. The training-based scheme in [8] first scans all the beams and retains only a few strong ones. Then, the least-square (LS) algorithm is employed for estimating the reduced-dimensional beamspace channel. In [9], a modified version of [8] reduces the overhead of beam training by simultaneously scanning several beams with the help of power splitters at the BS. However, the aforementioned algorithms are not optimized for lens-based mmWave systems because the lens antenna array has energy-focusing capability, and the received signal matrix from the lens antenna array is characterized by sparsity and concentration. The support detection based scheme in [6] further reduces the pilot overhead, which directly estimates the channel support by exploiting the sparsity of the beamspace channel. In [7], the channel matrix

is regarded as a 2-dimensional (2D) natural image and is then estimated by the cosparse analysis approximate message passing (SCAMPI) algorithm derived from the image recovery field. The SCAMPI algorithm models the channel as a sparse generic $L$-term Gaussian mixture (GM) probability distribution and uses the expectation-maximization (EM) algorithm to learn the GM parameters from the current estimated data.

Previous works only address narrowband mmWave systems. For wideband mmWave massive MIMO systems, the physical propagation delays of electromagnetic waves traveling across the whole array will become large and comparable to the time-domain sample period. In such a case, different antenna elements will receive different time-domain symbols, which is known as the spatial-wideband effect [30] and causes beam squint in the frequency domain. As a result, the AoAs (AoDs) will become frequency-dependent; thereby, channel estimation becomes very challenging, especially in mmWave band. The successive support detection (SSD) technique proposed in [31] applies successive interference cancelation to estimate the channel. The main idea is that each sparse path component has frequency-dependent support determined by its spatial direction and is estimated using beamspace windows. However, some important characteristics of mmWave channels, such as sparsity and channel correlation between adjacent antennas and subcarriers, are not considered. These characteristics are significant for performance improvement in channel estimation, but are difficult to be characterized by traditional model-based method.

Recently, deep learning (DL) has been applied to physical layer communications [10]–[12], [14]–[18], such as channel state information (CSI) feedback [12], signal detection [13], [14], channel estimation [15]–[17], precoder design [18], traffic analysis [19]–[21], and end-to-end transceiver design [22], [23]. DL-based physical layer communications may be data-driven and model-driven [11]. By incorporating domain knowledge into network design, model-driven DL can reduce the demand for computing resources and training time, which is more attractive for wireless communications [15]. As a promising model-driven DL approach, deep unfolding has been first proposed in [24] and applied to sparse signal recovery [25] and image processing [26]. The main idea is unfolding the iterative algorithm into a deep neural network and adding learnable parameters, which has been successfully applied to physical layer communications [11], [15], [27]. For example, the deep unfolding-based channel estimator has been successfully applied for narrowband beamspace mmWave massive MIMO systems in [15]. It outperforms state-of-the-art CS-based algorithms and can achieve excellent performance even with a small number of RF chains. However, the existing DL-based channel estimator utilizes a supervised

way [15]–[17], thereby a large number of real channel data are required to train the network, which defeats the point of channel estimation in the first place. This is because we cannot obtain the true channel data for training the network when the system is equipped in a practical environment. Therefore, how to train the DL-based channel estimator without the true channel data is significantly important.

## B. Contributions

In this study, we develop a DL-based channel estimator for lens-based mmWave massive MIMO systems. Instead of considering supervised learning for narrowband beamspace channel estimation [15], we investigate unsupervised learning for a wideband system and take the beam squint effect into consideration. To the best of our knowledge, this paper is the first study applies model-driven unsupervised DL network into wideband mmWave beamspace massive MIMO systems and considers the beam squint effect. The main contributions are summarized as follows.

- We first formulate the wideband beamspace channel estimation problem as a compressed image recovery problem. By incorporating an advanced denoising convolutional neural network (DnCNN) into the generalized expectation consistent signal recovery (GEC-SR) algorithm [32], [33], we develop a model-driven DL network, named the learned denoising-based GEC (LDGEC) network. The LDGEC network uses the Steins unbiased risk estimator (SURE) as the loss function; thereby, it can be trained only with the received signals not the real channel data. By utilizing *layer-by-layer training*, the LDGEC-based estimator can significantly outperform state-of-the-art CS algorithms even without the real channel data.
- Even if designed for unsupervised learning, the LDGEC network can also be supervisingly trained with the real channel data, thereby further improve channel estimation performance with available channel data. In this case, we can train the DnCNN denoiser in the *denoiser-by-denoiser* way, where the DnCNN denoiser is trained independently without including the whole GEC algorithm, thereby reducing the training complexity significantly.
- To further reduce the cost and power consumption, we investigate the LDGEC-based channel estimator for systems with hardware-friendly low-resolution ADCs. Numerical results demonstrate that little performance loss is caused for LDGEC-based channel estimator when the mmWave beamspace system is with low-resolution ADCs and reduced number of RF chains.
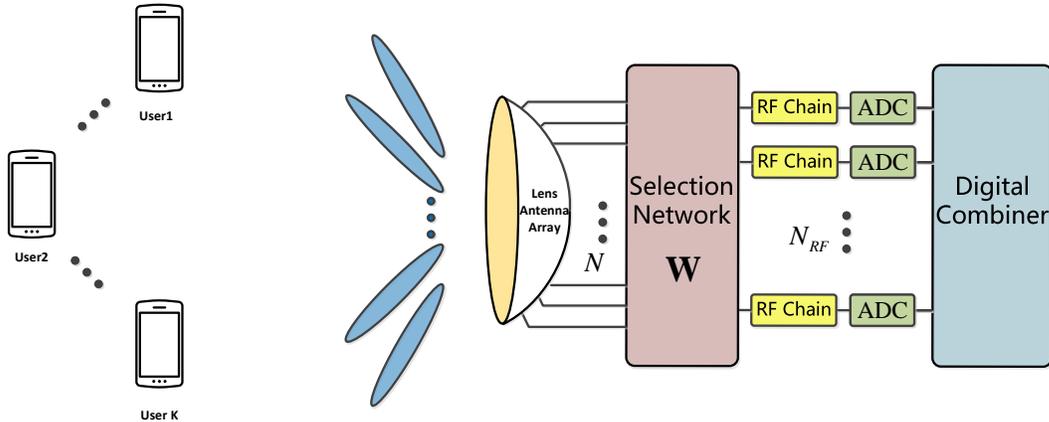
Fig. 1. The architecture of lens-based wideband beamspace mmWave MIMO-OFDM system.

*Notations*—For any matrix $\mathbf{A}$, $\mathbf{A}^T$ and $\mathrm{tr}(\mathbf{A})$ denote the transpose and the trace of $\mathbf{A}$, respectively. In addition, $\mathrm{Diag}(\mathbf{v})$ is the diagonal matrix with $\mathbf{v}$ on the diagonal, and $\mathbf{d}(\mathbf{Q})$ is the diagonalization operator, which returns a constant vector containing the average diagonal elements of $\mathbf{Q}$. Furthermore, $\mathbb{E}\{\cdot\}$ represents the expectation operator. A circular complex Gaussian with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Omega}$ can be described by the probability density function:

$$\mathcal{N}_{\mathbb{C}}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\Omega}) = \frac{1}{\det(\pi \boldsymbol{\Omega})} e^{-(\mathbf{z}-\boldsymbol{\mu})^H \boldsymbol{\Omega}^{-1}(\mathbf{z}-\boldsymbol{\mu})}.$$

We use $\mathrm{D}z$ to denote the real Gaussian integration measure

$$Dz = \phi(z)dz, \quad \phi(z) \triangleq \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}},$$

$\mathrm{D}z_c = \frac{e^{-|z|^2}}{\pi} dz$ to denote the complex Gaussian integration measure, $\Phi(x) \triangleq \int_{-\infty}^{x} \mathrm{D}z$ to denote the cumulative Gaussian distribution function.

The remaining part of this paper is organized as follows. Section II formulates the wideband beamspace channel estimation as a compressed image recovery problem. Next, a model-driven DL network is provided for beamspace channel estimation using SURE loss in Section III. Furthermore, the network can also be trained with the real channel data in Section IV. Then, numerical results are presented in Section V. Finally, Section VI concludes the paper.

## II. System Model and Problem Formulation

In this section, we first present the lens-based mmWave MIMO-OFDM systems. After introducing the beam squint effect, we formulate the wideband beamspace channel estimation as a compressed image recovery problem.

## A. Beamspace channel model

As illustrated in Fig.1, we consider an uplink wideband bemspace mmWave MIMO-OFDM system, where the BS employs an $N$-element lens antenna array and $N_{RF}$ RF chains to simultaneously serve $K$ single-antenna users. Applying the classical Saleh-Valenzuela channel model [34], the spatial channel $\mathbf{h}_m \in \mathbb{C}^{N \times 1}$ at sub-carrier $m$ is given by

$$\mathbf{h}_m = \sqrt{\frac{N}{L}} \sum_{l=1}^{L} \alpha_l e^{-j2\pi\tau_l f_m} \mathbf{a}(\phi_{l,m}), \tag{1}$$

for $m = 1, 2, \ldots, M$ where $L$ is the number of resolvable paths, $\alpha_l$ and $\tau_l$ are the complex gain and the time delay of the $l$-th path, respectively. Furthermore, $\mathbf{a}(\phi_{l,m})$ is the array response vector and $\phi_{l,m}$ is the spatial direction at sub-carrier $m$ defined as

$$\phi_{l,m} = \frac{f_m}{c} d \sin \theta_l, \tag{2}$$

where $f_m = f_c + \frac{f_s}{m}(m - 1 - \frac{M-1}{2})$ is the frequency of sub-carrier $m$ with $f_c$ and $f_s$ representing the carrier frequency and bandwidth, respectively. Furthermore, $c$ is the speed of light, $\theta_l$ is the physical direction, and $d$ is the antenna spacing, which is usually designed according to the carrier frequency as $d = c/2f_c$. Consider a uniform linear lens array in the BS, the array response vector $\mathbf{a}(\phi_{l,m})$ is given by,

$$\mathbf{a}(\phi_{l,m}) = \frac{1}{\sqrt{N}}[e^{-j2\pi\phi_{l,m}(-\frac{N-1}{2})}, e^{-j2\pi\phi_{l,m}(-\frac{N+1}{2})}, \ldots, \tag{3}$$
$$e^{-j2\pi\phi_{l,m}(\frac{N-1}{2})}]^T.$$

The conventional channel in the spatial domain in (1) can be transformed to the beamspace domain by employing a carefully designed lens antenna array, as shown in Fig. 1. Specifically, this lens antenna array plays the role of an $N$-element spatial discrete Fourier transform (DFT) matrix $\mathbf{F}$, which contains the array response vectors of $N$ orthogonal directions (beams) covering the entire space as

$$\mathbf{F} = [\mathbf{a}(\bar{\phi}_1), \mathbf{a}(\bar{\phi}_2), \ldots, \mathbf{a}(\bar{\phi}_N)], \tag{4}$$

where $\bar{\phi}_n = \frac{1}{N}(n - \frac{N+1}{2})$ for $n = 1, 2, \ldots, N$ are the spatial directions pre-defined by the lens antenna array. Accordingly, the wideband beamspace channel $\tilde{\mathbf{h}}_m$ at sub-carrier $m$ can be expressed as

$$\tilde{\mathbf{h}}_m = \mathbf{F}^H \mathbf{h}_m = \sqrt{\frac{N}{L}} \sum_{l=1}^{L} \alpha_l e^{-j2\pi\tau_l f_m} \tilde{\mathbf{c}}_{l,m}, \tag{5}$$
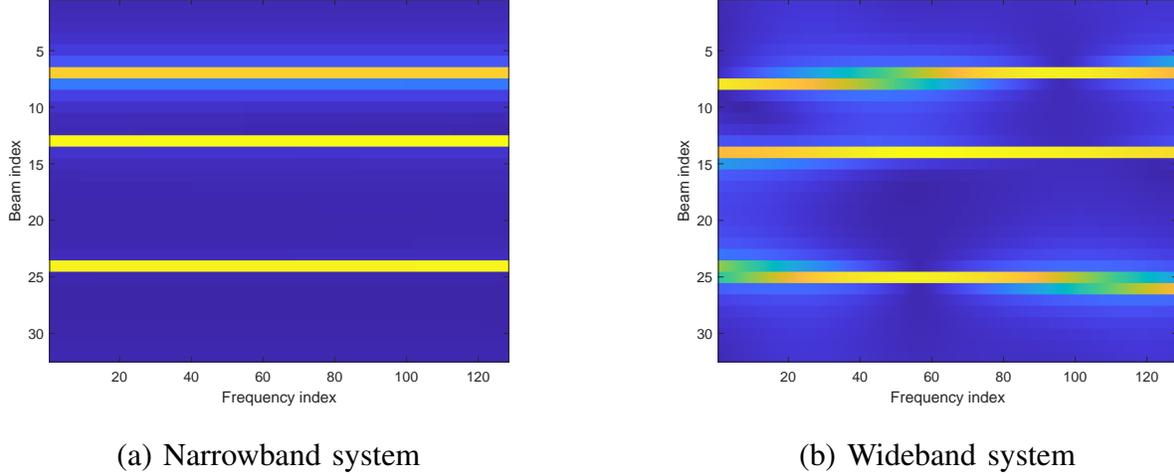
(a) Narrowband system

(b) Wideband system

Fig. 2. An illustration of a beam-frequency channel, where $L = 3$, $M = 128$, $N = 32$, $f_c = 28$ GHz, $f_s = 4$ GHz.

where $\tilde{\mathbf{c}}_{l,m}$ denotes the $l$-th path component at sub-carrier $m$ in the beamspace, and $\tilde{\mathbf{c}}_{l,m}$ is determined by $\phi_{l,m}$ as

$$\tilde{\mathbf{c}}_{l,m} = \mathbf{F}^H \mathbf{a}(\phi_{l,m}) \tag{6}$$

$$= [\Xi(\phi_{l,m} - \bar{\phi}_1), \Xi(\phi_{l,m} - \bar{\phi}_2), \ldots, \Xi(\phi_{l,m} - \bar{\phi}_N)]^T,$$

where $\Xi(x) = \frac{\sin N2\pi x}{\sin \pi x}$ is the Dirichlet sinc function.

*B. Beam Squint*

Before formulating the wideband beamspace channel estimation problem, we introduce the beam squint effect [30]. Based on the power-focusing capability of $\Xi(x)$, we know that most of the power of $\tilde{\mathbf{c}}_{l,m}$ is focused on only small number of elements. Additionally, due to the limited scattering in the mmWave systems, the number of resolvable paths, $L$, is generally small. However, the beam power distribution of the $l$-th path component will be different at different sub-carriers, i.e., $\tilde{\mathbf{c}}_{l,m_1} \neq \tilde{\mathbf{c}}_{l,m_2}$ for $m_1 \neq m_2$, since $\phi_{l,m}$ is frequency-dependent in wideband mmWave systems (i.e., $f_m \neq f_c$). This effect is termed as beam squint [30], which is a key difference between wideband and narrowband systems.

To show the beam squint effect, we present the energy diagram of the beam-frequency channel matrix in narrowband and wideband systems. We consider a beamspace system with $L = 3$, $M =$

128, $N = 32$, $f_c = 28$ GHz. Furthermore, we set $f_s = 4$ GHz in a wideband system and $f_s = 20$ MHz in the narrowband system. As illustrated in Fig. 2, the beam power distribution of the $l$-th path component is almost similar at different sub-carriers in the narrowband system. Therefore, the beamspace channel support (the index of non-zero elements) at different frequencies can be assumed to be the same.

Owing to the beam squint effect, the beam power distribution in wideband systems varies significantly over frequency. We denote the beam channel vector corresponding to the $l$-th path and the $m$-th frequency as $\tilde{\mathbf{h}}_{l,m}$. As shown is in Fig. 2, the index of the strongest element in $\tilde{\mathbf{h}}_{3,0}$ (i.e., the yellow bar at the bottom of the Fig. 2) is 24, while the index of the strongest element in $\tilde{\mathbf{h}}_{3,256}$ is 26. Thus, the beamspace channel supports at different frequencies are different. The characteristic of the beam-frequency matrix will bring a significant challenge for wideband beamspace channel estimation. Furthermore, the channel correlation between adjacent antennas and subcarriers is subtle, which is difficult to be characterized by the traditional approaches. Conversely, DL has the powerful capability to learn the correlation from the data, which is more promising for wideband channel estimation involved in the beam squint effect.

## C. Problem Formulation

In uplink channel estimation, the user devices transmit pilot sequences to the BS, and the channel is assumed to remain unchanged during this period. We use the orthogonal pilot sequence for different users. Therefore, the channel estimation can be performed for each user independently. Considering a specific user, the pilot at sub-carrier $m$ in instant $q$, $s_{m,q}$, is transmitted. The received signal vector $\mathbf{y}_{m,q} \in \mathbb{C}^{N \times 1}$ at the BS is given by

$$\mathbf{y}_{m,q} = \mathbf{W}_q \tilde{\mathbf{h}}_m s_{m,q} + \mathbf{W}_q \mathbf{n}_{m,q}, \tag{7}$$

for $m = 1, 2, \ldots, M$ and $q = 1, 2, \ldots, Q$, where $\mathbf{n}_{m,q} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{0}, \sigma_n^2 \mathbf{I})$ represents a Gaussian noise vector. $\mathbf{W}_q \in \mathbb{C}^{N_{RF} \times N}$ is the adaptive selection network but fixed for different sub-carriers. We set $s_{m,q} = 1$ for convenience as the pilot signal is known at the receiver side. Thus, the received signal $\bar{\mathbf{y}}_m$ in $Q$ instants is given by

$$\bar{\mathbf{y}}_m = [\mathbf{y}_{m,1}^T, \ldots, \mathbf{y}_{m,Q}^T]^T = \bar{\mathbf{W}} \tilde{\mathbf{h}}_m + \mathbf{W} \mathbf{n}_m, \tag{8}$$

where $\bar{\mathbf{W}} = [\mathbf{W}_1^T, \mathbf{W}_2^T, \ldots, \mathbf{W}_Q^T]^T \in \mathbb{C}^{Q N_{RF} \times N}$ and $\mathbf{n}_m^{\text{eq}} = [\mathbf{n}_{m,1}^T, \ldots, \mathbf{n}_{n,Q}^T]^T$. In this paper, low-cost $one$-bit phase shifters are utilized in the adaptive selection network $\mathbf{W}_q$. Therefore, the elements of $\bar{\mathbf{W}}$ are randomly selected from the set $\frac{1}{\sqrt{Q N_{RF}}}\{-1, +1\}$ with equal probability.

From Fig. 2 and Section II-A, the beamspace channel vectors at different subcarriers, even if different due to beam squint, are correlated through antenna array response vector $\mathbf{a}(\phi_{l,m})$, which is highly similar to a 2D natural image. By stacking $M$ beamspace channel vectors into a matrix, we have the following signal recovery model,

$$[\bar{\mathbf{y}}_1, \bar{\mathbf{y}}_2, \ldots, \bar{\mathbf{y}}_M] = \bar{\mathbf{W}}[\tilde{\mathbf{h}}_1, \tilde{\mathbf{h}}_2, \ldots, \tilde{\mathbf{h}}_M] + [\mathbf{n}_1^{\mathrm{eq}}, \mathbf{n}_2^{\mathrm{eq}}, \ldots, \mathbf{n}_M^{\mathrm{eq}}]. \tag{9}$$

If we regard beam-frequency matrix $[\tilde{\mathbf{h}}_1, \tilde{\mathbf{h}}_2, \ldots, \tilde{\mathbf{h}}_M]$ as a 2D natural image, many compressed image recovery method can be borrowed here for beamspace channel estimation, which enables us to develop a model-driven-DL-based channel estimation network.

Before introducing the DL-based channel estimation network, we first obtain a linear transformation model by stacking the $\mathbf{y}_m$, $\tilde{\mathbf{h}}_m$ and $\mathbf{n}_m^{\mathrm{eq}}$ into

$$\mathbf{y} = \mathbf{A}\mathbf{h} + \mathbf{n}, \tag{10}$$

where $\mathbf{y} = [\bar{\mathbf{y}}_1^T, \bar{\mathbf{y}}_2^T, \ldots, \bar{\mathbf{y}}_M^T]$, $\mathbf{h} = [\tilde{\mathbf{h}}_1^T, \tilde{\mathbf{h}}_2^T, \ldots, \tilde{\mathbf{h}}_M^T]^T$, $\mathbf{n} = [(\mathbf{n}_1^{\mathrm{eq}})^T, (\mathbf{n}_2^{\mathrm{eq}})^T, \ldots, (\mathbf{n}_M^{\mathrm{eq}})^T]^T$, and $\mathbf{A} = (\mathbf{I} \otimes \bar{\mathbf{W}})$. We denote $\otimes$ as the matrix Kronecker product. The linear transformation model (10) will be utilized in the subsequent discussion.

## III. Unsupervised Learning for Beamspace Channel Estimation

In this section, we propose a model-driven unsupervised DL network for wideband beamspace channel estimation, named LDGEC-based channel estimator. As in [11], the network is specially designed by unfolding an iterative algorithm, GEC algorithm, with the DL-based denoiser. After introducing the network architecture and DnCNN denoiser, we elaborate the SURE loss and the layer-by-layer training approach, which are critical to implementing LDGEC network with unsupervised learning.

### A. LDGEC-based channel estimator

As illustrated in Fig. 3, the input of the LDGEC network is the received signal vector, $\tilde{\mathbf{y}}$, and the linear transform matrix, $\mathbf{A}$, while the final output is $\hat{\mathbf{h}}^{\mathrm{out}}$, the estimated channel vector. The LDGEC network consists of $T$ layers connected in cascade. We replace the posterior mean estimator in the GEC algorithm with the DnCNN denoiser and deep unfold the GEC algorithm into neural network. Each iteration of the GEC algorithm can be interpreted as each layer of the LDGEC network. As each layer of LDGEC has the same structure except the learnable parameters in DnCNN denoiser, we omit the layer index $t$ in Fig. 3 and Algorithm 1.
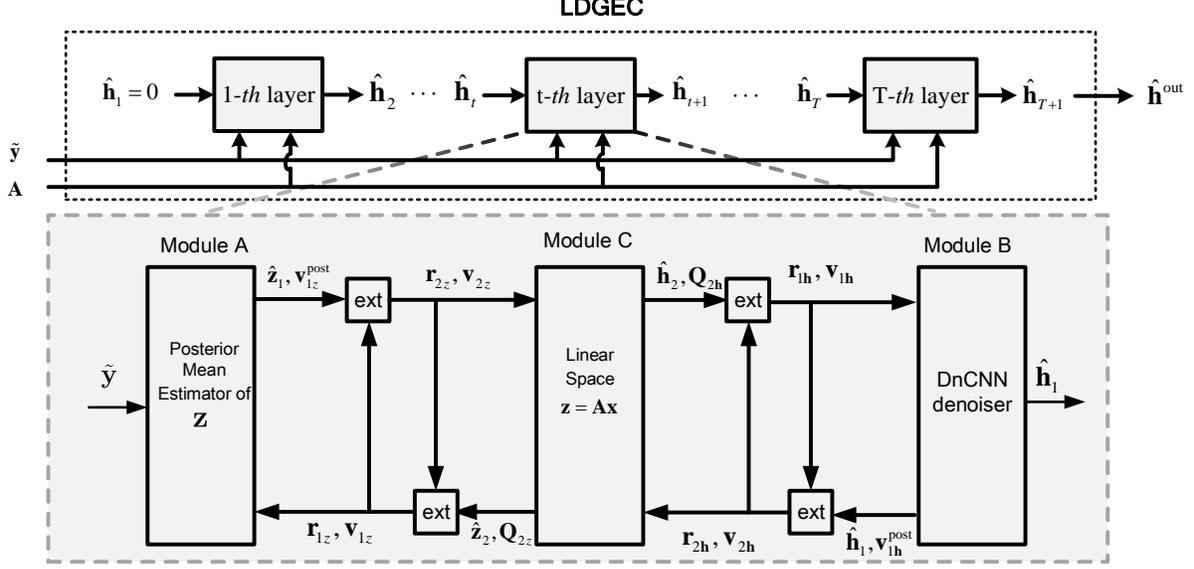
Fig. 3. The network structure of LDGEC-based channel estimator.

As illustrated in the figure, each layer of the LDGEC network has three modules. Specifically, Module A computes the posterior mean and variance of $\mathbf{z} = \mathbf{A}\mathbf{h}$, module B performs denoising from the noisy signal, $\mathbf{r}_{1h}$, by using the advanced DnCNN denoiser, and module C provides the framework that constrains the estimation problem into the linear space $\mathbf{z} = \mathbf{A}\mathbf{h}$. Modules A, B, and C are executed iteratively, as in the figure. In addition, each module uses the turbo principle as in iterative decoding; that is, each module passes the extrinsic messages to its next module. The three modules are executed iteratively until convergence or terminated by fixed number of layers.

Before introducing the principle of the LDGEC network, we define two auxiliary variables,

$$P_h = 1 \quad \text{and} P_z = P_h \cdot \text{tr}(\mathbf{A}^H\mathbf{A})/MN_{RF}Q, \tag{11}$$

which are interpreted as the powers of $h_n$ and $z_n$, respectively. $h_n$ and $z_n$ denote the $n$-th element in $\mathbf{h}$ and $\mathbf{z}$, respectively. The $P_h$ and $P_z$ are important for network initialization. The algorithm for the LDGEC-based channel estimator is listed in Algorithm 1.

To better understand the LDGEC network, we provide detailed explanations. Lines 1–2 compute the posterior mean and variance of $z_n$ from quantized measurements $\tilde{y}_n$, and the expectation

---

**Algorithm 1:** LDGEC-based channel estimator

---

**Input:** Received signals $\tilde{\mathbf{y}}$, linear transform matrix $\mathbf{A}$, likelihood $\mathsf{P}(\tilde{\mathbf{y}}|\mathbf{z})$

**Output:** Recovered signal $\hat{\mathbf{h}}^{\mathrm{out}}$.

**Initialize:** $t \leftarrow 1$, $\mathbf{r}_{1\mathbf{z}} \leftarrow \mathbf{0}$, $\mathbf{r}_{2\mathbf{h}} \leftarrow \mathbf{0}$, $\mathbf{v}_{1\mathbf{z}} \leftarrow P_z \mathbf{1}$, and $\mathbf{v}_{2\mathbf{h}} \leftarrow P_h \mathbf{1}$.

**for** $t = 1, \cdots, T$ **do**

    **Module A:**

    (1) Compute the posterior mean and covariance of $\mathbf{z}$

1    $\hat{\mathbf{z}}_1 = \mathrm{E}\left\{\mathbf{z}|\mathbf{r}_{1\mathbf{z}}, \mathbf{v}_{1\mathbf{z}}\right\}$,

2    $\mathbf{v}_{1\mathbf{z}}^{\mathrm{post}} = \mathrm{Var}\left\{\mathbf{z}|\mathbf{r}_{1\mathbf{z}}, \mathbf{v}_{1\mathbf{z}}\right\}$.

    (2) Compute the extrinsic information of $\mathbf{z}$

3    $\mathbf{v}_{2\mathbf{z}} = \mathbf{1} \oslash \left(\mathbf{1} \oslash \mathbf{v}_{1\mathbf{z}}^{\mathrm{post}} - \mathbf{1} \oslash \mathbf{v}_{1\mathbf{z}}\right)$,

4    $\mathbf{r}_{2\mathbf{z}} = \mathbf{v}_{2\mathbf{z}} \odot \left(\hat{\mathbf{z}}_1 \oslash \mathbf{v}_{1\mathbf{z}}^{\mathrm{post}} - \mathbf{r}_{1\mathbf{z}} \oslash \mathbf{v}_{1\mathbf{z}}\right)$.

    **Module C:**

    (3) Compute the mean and covariance of $\mathbf{h}$ from the linear space

5    $\mathbf{Q}_{2\mathbf{h}} = \left(\mathrm{Diag}(\mathbf{1} \oslash \mathbf{v}_{2\mathbf{h}}) + \mathbf{A}^H \mathrm{Diag}(\mathbf{1} \oslash \mathbf{v}_{2\mathbf{z}})\mathbf{A}\right)^{-1}$,

6    $\hat{\mathbf{h}}_2 = \mathbf{Q}_{2\mathbf{h}} \left(\mathbf{r}_{2\mathbf{h}} \oslash \mathbf{v}_{2\mathbf{h}} + \mathbf{A}^H \mathbf{r}_{2\mathbf{z}} \oslash \mathbf{v}_{2\mathbf{z}}\right)$.

    (4) Compute the extrinsic information of $\mathbf{h}$

7    $\mathbf{v}_{1\mathbf{h}} = \mathbf{1} \oslash \left(\mathbf{1} \oslash \mathbf{d}(\mathbf{Q}_{2\mathbf{h}}) - \mathbf{1} \oslash \mathbf{v}_{2\mathbf{h}}\right)$,

8    $\mathbf{r}_{1\mathbf{h}} = \mathbf{v}_{1\mathbf{h}} \odot \left(\hat{\mathbf{h}}_2 \oslash \mathbf{d}(\mathbf{Q}_{2\mathbf{h}}) - \mathbf{r}_{2\mathbf{h}} \oslash \mathbf{v}_{2\mathbf{h}}\right)$.

    **Module B:**

    (5) Compute the mean and covariance of $\mathbf{h}$

9    $\hat{\mathbf{h}}_1 = D_{\hat{\sigma}}(\mathbf{r}_{1\mathbf{h}}, \mathbf{v}_{1\mathbf{h}})$

10    $\mathbf{v}_{1\mathbf{h}}^{\mathrm{post}} = \frac{1}{MN}\mathrm{div}D_{\hat{\sigma}}(\mathbf{r}_{1\mathbf{h}})\mathrm{avg}(\mathbf{v}_{1\mathbf{h}})$

    (6) Compute the extrinsic information of $\mathbf{h}$

11    $\mathbf{v}_{2\mathbf{h}} = \mathbf{1} \oslash \left(\mathbf{1} \oslash \mathbf{v}_{1\mathbf{h}}^{\mathrm{post}} - \mathbf{1} \oslash \mathbf{v}_{1\mathbf{h}}\right)$,

12    $\mathbf{r}_{2\mathbf{h}} = \mathbf{v}_{2\mathbf{h}} \odot \left(\hat{\mathbf{h}}_1 \oslash \mathbf{v}_{1\mathbf{h}}^{\mathrm{post}} - \mathbf{r}_{1\mathbf{h}} \oslash \mathbf{v}_{1\mathbf{h}}\right)$.

    **Module C:**

    (7) Compute the mean and covariance of $\mathbf{z}$ from the linear space

13    $\mathbf{Q}_{2\mathbf{h}} = \left(\mathrm{Diag}(\mathbf{1} \oslash \mathbf{v}_{2\mathbf{h}}) + \mathbf{A}^H \mathrm{Diag}(\mathbf{1} \oslash \mathbf{v}_{2\mathbf{z}})\mathbf{A}\right)^{-1}$,

14    $\hat{\mathbf{h}}_2 = \mathbf{Q}_{2\mathbf{h}} \left(\mathbf{r}_{2\mathbf{h}} \oslash \mathbf{v}_{2\mathbf{h}} + \mathbf{A}^H \mathbf{r}_{2\mathbf{z}} \oslash \mathbf{v}_{2\mathbf{z}}\right)$,

15    $\mathbf{Q}_{2\mathbf{z}} = \mathbf{A}\mathbf{Q}_{2\mathbf{h}}\mathbf{A}^H$,

16    $\hat{\mathbf{z}}_2 = \mathbf{A}\hat{\mathbf{h}}_2$.

    (8) Compute the extrinsic information of $\mathbf{z}$

17    $\mathbf{v}_{1\mathbf{z}} = \mathbf{1} \oslash \left(\mathbf{1} \oslash \mathbf{d}(\mathbf{Q}_{2\mathbf{z}}) - \mathbf{1} \oslash \mathbf{v}_{2\mathbf{z}}\right)$,

18    $\mathbf{r}_{1\mathbf{z}} = \mathbf{v}_{1\mathbf{z}} \odot \left(\hat{\mathbf{z}}_2 \oslash \mathbf{d}(\mathbf{Q}_{2\mathbf{z}}) - \mathbf{r}_{2\mathbf{z}} \oslash \mathbf{v}_{2\mathbf{z}}\right)$.

w.r.t. the posterior

$$\mathsf{P}_Z(z_n|\tilde{y}_n) = \frac{\mathsf{P}_{\text{out}}(\tilde{y}_n|z_n)\mathsf{P}_Z(z_n)}{\int \mathsf{P}_{\text{out}}(\tilde{y}_n|z_n)\mathsf{P}_Z(z_n)dz_n}, \tag{12}$$

where $\mathsf{P}_Z(z_n)$ is assumed to be $\mathcal{N}_{\mathbb{C}}(z_n; r_{1z,n}, v_{1z,n})$. To clearly understand Lines 1 and 2 in Algorithm 1, we take the quantized and unquantized channels as two examples.

**Unquantized channel**: If with infinite-resolution ADCs, the received signal at the BS, $\tilde{\mathbf{y}} = \mathbf{y}$ and the posterior probability $\mathsf{P}_{\text{out}}(\tilde{y}_n|z_n)$ is given by

$$\mathsf{P}_{\text{out}}(\tilde{y}_n|z_n) = \frac{1}{\pi\sigma_n^2}e^{|\tilde{y}_n - z_n|/\sigma_n^2}. \tag{13}$$

Thus, the explicit expressions of the posterior mean and variance will be

$$\hat{z}_1 = r_{1z} + \frac{v_{1z}}{v_{1z} + \sigma_n^2}(\tilde{y} - r_{1z}), \tag{14}$$

$$v_{1z}^{\text{post}} = v_{1z} - \frac{v_{1z}^2}{v_{1z} + \sigma_n^2}, \tag{15}$$

**Quantized channel**: If the low-resolution ADCs are used in the BS, the received signal $\tilde{\mathbf{y}} = \mathcal{Q}_c(\mathbf{y})$, where $\mathcal{Q}_c$ is the complex-valued quantizer. Then, the explicit expressions of the posterior mean and variance can be derived similar to [35, Appendix A] as

$$\hat{z}_1 = r_{1z} + \frac{\text{sign}(\tilde{y})v_{1z}}{\sqrt{2(\sigma_n^2 + v_{1z})}}\left(\frac{\phi(\eta_1) - \phi(\eta_2)}{\Phi(\eta_1) - \Phi(\eta_2)}\right), \tag{16}$$

$$v_{1z}^{\text{post}} = \frac{v_{1z}}{2} - \frac{(v_{1z})^2}{2(\sigma_n^2 + v_{1z})} \times$$
$$\left(\frac{\eta_1\phi(\eta_1) - \eta_2\phi(\eta_2)}{\Phi(\eta_1) - \Phi(\eta_2)} + \left(\frac{\phi(\eta_1) - \phi(\eta_2)}{\Phi(\eta_1) - \Phi(\eta_2)}\right)^2\right), \tag{17}$$

where

$$\eta_1 = \frac{\text{sign}(\tilde{y})r_{1z} - \min\{|r^{\text{low}}|, |r^{\text{up}}|\}}{\sqrt{(\sigma_n^2 + v_{1z})/2}}, \tag{18a}$$

$$\eta_2 = \frac{\text{sign}(\tilde{y})r_{1z} - \max\{|r^{\text{low}}|, |r^{\text{up}}|\}}{\sqrt{(\sigma_n^2 + v_{1z})/2}}, \tag{18b}$$

where $r^{\text{low}}$ and $r^{\text{up}}$ are the lower and upper thresholds associated with $\tilde{y}_n$, respectively. For notational convenience, we omit index $n$ and have

$$r^{\text{low}} = \begin{cases} \tilde{y} - \frac{\Delta}{2}, & \text{for } \tilde{y} \geq -\left(\frac{2^{\kappa}}{2} - 1\right)\Delta, \\ -\infty, & \text{otherwise}, \end{cases} \tag{19a}$$

and

$$r^{\text{up}} = \begin{cases} \tilde{y} + \frac{\Delta}{2}, & \text{for } \tilde{y} \leq \left(\frac{2^\kappa}{2} - 1\right)\Delta, \\ \infty, & \text{otherwise.} \end{cases} \tag{19b}$$

In this paper, we mainly focus on a typical uniform midrise quantizer with quantization step size $\Delta$. It maps a real-valued input into the nearest value in

$$\mathcal{R}_\kappa \triangleq \left\{ \left(-\frac{1}{2} + b\right)\Delta; \; b = -\frac{2^\kappa}{2} + 1, \cdots, \frac{2^\kappa}{2} \right\}, \tag{20}$$

where $\kappa$ is the quantization bits.

The real and imaginary parts are quantized separately, and each complex-valued channel can be decoupled into two real-valued channels. Expressions (16) and (17) are the estimators only for the real part of $\hat{z}_1$. To facilitate notation, we have abused $\tilde{y}$ and $\hat{z}_1$ in (16) and (17) to denote $\text{Re}(\tilde{y})$ and $\text{Re}(\hat{z}_1)$, respectively, and we omit index $n$ in the aforementioned expression. The estimator for the imaginary part $\text{Im}(\hat{z}_1)$ can be obtained similarly as (16) and (17) while $\tilde{y}$ and $b$ should be replaced by $\text{Im}(\tilde{y})$ and $b'$, respectively.

Lines 3–4 compute the extrinsic information of $\mathbf{z}$ using the turbo principle. Lines 5–6 perform the linear minimum mean-squared error (LMMSE) estimate of $\mathbf{h}$ under the following assumption,

$$\mathbf{r_{2z}} = \mathbf{z}_2 + \mathbf{w_{2z}}, \tag{21}$$

where $\mathbf{w_{2z}} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{0}, \text{Diag}(\mathbf{v_{2z}}))$, $\mathbf{z}_2 = \mathbf{A}\mathbf{h}_2$, and $\mathbf{h}_2 \sim \mathcal{N}_{\mathbb{C}}(\mathbf{h}_2; \mathbf{r_{2h}}, \text{Diag}(\mathbf{v_{2h}}))$. Lines 7–8 compute the extrinsic information of $\mathbf{h}$ and pass it to module B as a prior information. Lines 9–10 estimate the mean, $\hat{\mathbf{h}}_1$, and variance, $\mathbf{v}_{1\mathbf{h}}^{\text{post}}$ by considering the true prior $\mathsf{P}(\mathbf{h})$, which is assumed to estimate $\mathbf{h}$ from several AWGN observations, that is,

$$\mathbf{r_{1h}} = \mathbf{h} + \mathbf{w_{1h}}, \tag{22}$$

where $\mathbf{w_{1h}} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{0}, \text{Diag}(\mathbf{v_{1h}}))$. As channel $\mathbf{h}$ can be regarded as a 2D natural image, we utilize the advanced DnCNN denoiser [36] in Lines 9–10 to recover channel $\mathbf{h}$ from equivalent noisy observations $\mathbf{r_{1h}}$. Lines 11–12 compute the extrinsic information of $\mathbf{h}$ using the turbo principle, and lines 13–16 constrain the estimated problem into a linear space $\mathbf{z} = \mathbf{A}\mathbf{h}$ which performs the same procedure as Lines 5–6. Lines 17–18 compute the extrinsic information of $\mathbf{z}$ and pass to module A as the prior information.

Posterior variance $\mathbf{v}_{1h}^{\text{post}}$ is determined by $\text{div}D_{\hat{\sigma}^t}(\mathbf{r}_{1h})\text{avg}(\mathbf{v}_{1h})$, where the divergence $\text{div}D_{\hat{\sigma}^t}(\mathbf{r}_{1h})$ is simply the sum of the partial derivates with respect to each element of $\mathbf{r}_{1h}$. It can be expressed by

$$\text{div}D_{\hat{\sigma}^t}(\mathbf{r}_{1h}) = \sum_{i=1}^{n} \frac{\partial D_{\hat{\sigma}^t}(\mathbf{r}_{1h})}{\partial r_{1h,i}}, \tag{23}$$

where $r_{1h,i}$ is the $i$-th element of $\mathbf{r}_{1h}$. Although simple denoisers often yield a closed form for their divergence, high-performance denoisers are often data-dependent; making it very difficult to characterize their input-output relationship explicitly for most DL-based denoisers. Therefore, we calculate a good approximation for the divergence.

We use the following Monte-Carlo approximation to compute divergence $\text{div}D_{\hat{\sigma}^t}(\cdot)$. Using an independent and identically distributed (i.i.d.) random vector $\mathbf{b} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, we can estimate the divergence with

$$\text{div}D_{\hat{\sigma}^t} = \lim_{\epsilon \to 0} \mathbb{E}_{\mathbf{b}} \left\{ \mathbf{b}^T \left( \frac{D_{\hat{\sigma}^t}(\mathbf{r}_{1h} + \epsilon\mathbf{b}) - D_{\hat{\sigma}^t}(\mathbf{r}_{1h})}{\epsilon} \right) \right\} \tag{24}$$

$$\approx \frac{1}{\epsilon} \mathbf{b}^T (D_{\hat{\sigma}^t}(\mathbf{r}_{1h} + \epsilon\mathbf{b}) - D_{\hat{\sigma}^t}(\mathbf{r}_{1h})), \tag{25}$$

where $\epsilon = \|\mathbf{r}_{1h}\|_{\infty}/1000$ is an arbitrary small number. Equation (24) is originated from the law of large numbers. The expectation can be approximated with Monte Carlo sampling and a single sample can well approximate the expectation.

To improve robustness, we use an auto-regressive filter to smooth the update of $(\mathbf{v}_{1z}, \mathbf{r}_{1z})$ by

$$\mathbf{v}_{1z}^{t+1} = \beta \cdot \mathbf{1} \oslash \left( \mathbf{1} \oslash \mathbf{d}(\mathbf{Q}_{2z}^{t+1}) - \mathbf{1} \oslash \mathbf{v}_{2z}^{t+1} \right) + (1 - \beta)\mathbf{v}_{1z}^t, \tag{26}$$

$$\mathbf{r}_{1z}^{t+1} = \beta \cdot \mathbf{v}_{1z}^{t+1} \oslash \left( \mathbf{1} \oslash \mathbf{d}(\mathbf{Q}_{2z}^{t+1}) - \mathbf{1} \oslash \mathbf{v}_{2z}^{t+1} \right) + (1 - \beta)\mathbf{r}_{1z}^t, \tag{27}$$

where a small $\beta$ is the damping factor. Furthermore, a small constant threshold $\epsilon_1 = 5 \times e^{-7}$ should be set to restrict the minimum variance allowed per iteration and avoid numerical instabilities, that is, $\mathbf{v}_{1z} = \max(\epsilon_1, \mathbf{v}_{1z})$ and $\mathbf{v}_{1h}^{\text{post}} = \max(\epsilon_1, \mathbf{v}_{1h}^{\text{post}})$.

### B. DnCNN denoiser

The denoiser used in the LDGEC network plays a key role in channel estimation. We consider the state-of-the-art DnCNN denoiser[1]. The DnCNN is first proposed in [36] to handle the

---

[1] Although several new denoisers have been proposed for image denoising problem recently [37], they have similar performance with DnCNN denoiser.
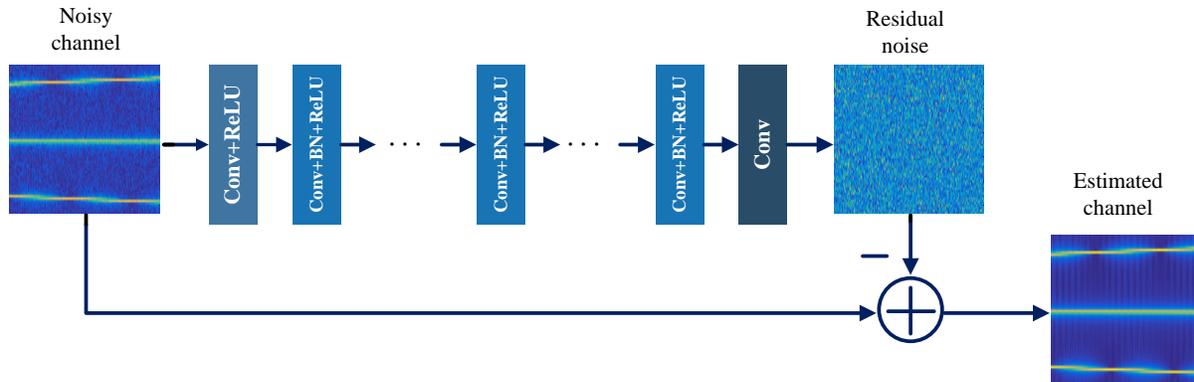
Fig. 4. Network architecture of the DnCNN denoiser.

Gaussian denoising problem with an unknown noise level, which is more accurate and faster than competing techniques. Fig. 4 illustrates the network architecture of the DnCNN denoiser. It consists of 20 convolutional layers. The first convolutional layer uses 64 different $3 \times 3 \times 1$ filters and is followed by a rectified linear unit (ReLU). Each of the succeeding 18 convolutional layers uses 64 different $3 \times 3 \times 64$ filters, each followed by batch-normalization and a ReLU. The final convolutional layer uses one separate $3 \times 3 \times 64$ filters to reconstruct the signal. Instead of learning a mapping directly from a noisy image to a denoised image, learning the residual noise is beneficial.

We plot three pseudo-color images of noisy channel, residual noise, and estimated channel in Fig. 4. The network is given the noisy observation $\mathbf{h} + \hat{\sigma}\mathbf{w}$ as an input, where $\mathbf{w}$ is the AWGN and noise variance $\hat{\sigma}$ is uniformly generated from a specific interval. The network produces residual noise $\hat{\mathbf{z}}$, rather than an estimated channel $\hat{\mathbf{h}}$, as an output. This method, known as residual learning [38], renders the network to remove the highly structured natural image rather than the unstructured noise. Consequently, residual learning improves both the training times and accuracy of a network. Furthermore, the DnCNN adopts the method of batch normalization, can speed up the training process, and boost the denoising performance [36].

*C. Stein's unbiased risk estimator*

Recently, many DL-based channel estimators have been proposed for different communication scenarios [15]. A common limitation of these works is the extensive real channel data should be obtained before training the network because the MSE function $\mathrm{MSE} = \mathbb{E}\|\hat{\mathbf{h}} - \mathbf{h}\|$ involves in the

real channel data $\mathbf{h}$. These requirements bring significant challenges when the real data cannot be obtained, e.g., the DL-based channel estimator is equipped in a new channel environment and only received signal $\tilde{\mathbf{y}}$ is obtained at the BS. Furthermore, the existing DL-based channel estimator utilizes supervised way and requires a large number of real channel data, which defeats the point of channel estimation. In this circumstance, how to train the network without the real channel data and only with measurements corresponding to the pilot symbols is significantly important. To solve the problem, we introduce the SURE loss function [39], which is a classical approach for learning images from noisy observations and has been extended to linear noisy measurements. It has been applied in medical imaging, microscopy, and astronomy, where the ground truth data is rarely available [40]. The SURE loss is an unbiased estimator of

$$\text{MSE} = \mathbb{E}_{\mathbf{w}}[\frac{1}{P}\|\mathbf{h} - f(\mathbf{y_w})\|^2], \tag{28}$$

where $\mathbf{y_w} = \mathbf{h} + \mathbf{w}$ and $f(\cdot)$ is an estimator of $\mathbf{h}$ from $\mathbf{y_w}$. Therefore, it can be used for training a DL-based denoiser to replace the MSE loss.

The goal of channel estimation is to reconstruct a channel $\mathbf{h}$ from noisy linear observations $\mathbf{y} = \mathbf{Ah} + \mathbf{n}$ with the known linear transform matrix $\mathbf{A}$. We are given training measurements $\mathbf{y}^1, \mathbf{y}^2, \ldots, \mathbf{y}^D$ but not the real channel $\mathbf{h}^1, \mathbf{h}^2, \ldots, \mathbf{h}^D$. Without access to $\mathbf{h}^1, \mathbf{h}^2, \ldots, \mathbf{h}^D$, the ground truth data, we cannot train the DnCNN denoiser by minimizing the traditional MSE loss function. Fortunately, we can use the SURE loss instead. SURE is a model selection technique first proposed by its namesake in [39]. It provides an unbiased estimate of the MSE for an estimator of the mean of a Gaussian distributed random vector with an unknown mean. Let $\mathbf{x}$ denote a vector we would like to estimate from noisy observations $\mathbf{r_{1h}} = \mathbf{h} + \mathbf{w_{1h}}$ where $\mathbf{w_{1h}} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{0}, \text{Diag}(\mathbf{w_{1h}}))$. We assume the DnCNN function $f_{\boldsymbol{\theta}}()$ is a weakly differentiable function parameterized by $\boldsymbol{\theta}$, which receives noisy observations $\mathbf{r_{1h}}$ as input and provides an estimate of $\mathbf{h}$ as output. Then, according to [39], [40], we can express the expectation of the MSE of the real channel $\mathbf{h}$ and equivalent noisy observations $\mathbf{r_{1h}}$ with respect to the random variable $\mathbf{w_{1h}}$ as follows,

$$\text{MSE} = \mathbb{E}_{\mathbf{w_{1h}}}[\frac{1}{P}\|\mathbf{h} - f_{\boldsymbol{\theta}}(\mathbf{r_{1h}})\|^2]. \tag{29}$$

Then, the MSE loss can be computed as follows,

$$\mathbb{E}_{\mathbf{w_{1h}}}[\frac{1}{P}\|\mathbf{h} - f_{\boldsymbol{\theta}}(\mathbf{r_{1h}})\|^2] = \mathbb{E}_{\mathbf{w_{1h}}}[\frac{1}{P}\|\mathbf{r_{1h}} - f_{\boldsymbol{\theta}}(\mathbf{r_{1h}})\|^2] - v_{1h}^2 \tag{30}$$
$$+ \frac{2v_{1h}^2}{P}\text{div}(f_{\boldsymbol{\theta}}(\mathbf{r_{1h}})),$$

where $P = MN$ and $\text{div}(\cdot)$ stands for divergence defined as (23). Note that two terms within the SURE loss depend on parameter $\boldsymbol{\theta}$. The first term, $\mathbb{E}_{\mathbf{w_{1h}}}[\frac{1}{P}\|\mathbf{r_{1h}} - f_{\boldsymbol{\theta}}(\mathbf{r_{1h}})\|^2]$, indicates the difference between estimate $f_{\boldsymbol{\theta}}(\mathbf{r_{1h}})$ and observation $\mathbf{r_{1h}}$ (bias). The second term, $\frac{2v_{1h}^2}{P}\text{div}(f_{\boldsymbol{\theta}}(\mathbf{r_{1h}}))$, penalizes the denoiser for varying as the input is changed. Thus, SURE is a natural way to control the trade-off between the bias and variance of a recovery algorithm.

The critical challenge for using SURE in practice is to compute divergence $\text{div}(f_{\boldsymbol{\theta}}(\mathbf{r_{1h}}))$. For the advanced DnCNN denoiser, the divergence is hard or even impossible to express analytically. Therefore, we cannot obtain a closed form for the divergence. Similar to (24), we can use a Montel Carlo method to estimate the divergence $\text{div}(f_{\boldsymbol{\theta}}(\mathbf{r_{1h}}))$. Combining the SURE loss in (30) and the estimate of divergence $\text{div}(f_{\boldsymbol{\theta}}(\mathbf{r_{1h}}))$, we can minimize the MSE loss function of a denoising problem without ground truth data. Note that minimizing SURE loss requires propagating gradients with respect to the Monte Carlo estimate of divergence (23). Although the gradients are challenging to compute by hand, we can resort to TensorFlow's auto-differentiation capabilities to propagate it.

### D. Layer-by-layer training

A significant reason for the LDGEC-based channel estimator trained with SURE loss is layer-by-layer training. From (30), the computational process of the SURE loss requires noisy observations $\mathbf{r_{1h}}$ and equivalent noise variance $v_{1h}$. This method takes advantage of the fact that each layer of the LDGEC-based channel estimator is to solve a denoising problem with known variance $v_{1h}$ and noisy observations $\mathbf{r_{1h}}$. As the LDGEC network can decouple the linear model in (10) into several equivalent AWGN models (22) in each layer and $\mathbf{v_{1h}}$ computed in line 7 in Algorithm 1 is accurate enough to describe the variance of $\mathbf{w_{1h}}$. Therefore, we can train the $t$-th layer LDGEC network with the SURE loss, estimated variance $\mathbf{v}_{\mathbf{1h}}^t$ and noisy observations $\mathbf{r}_{\mathbf{1h}}^t$ through layer-by-layer training.

In the $t$-th round of the layer-by-layer training, the loss function is given by

$$L_{\text{SURE}}^t(\boldsymbol{\theta}^t) = \arg\min_{\boldsymbol{\theta}^t} \sum_{d=1}^{D} [\frac{1}{P}\|\mathbf{r}_{\mathbf{1h}}^{t,d} - f_{\boldsymbol{\theta}}(\mathbf{r}_{\mathbf{1h}}^{t,d})\|^2] - v_{1h}^2 \tag{31}$$
$$+ \frac{2v_{1h}^2}{P}\text{div}(f_{\boldsymbol{\theta}}(\mathbf{r}_{\mathbf{1h}}^{t,d})),$$

where $D$ is the number of mini-batches in the $t$-th round and $\mathbf{r}_{\mathbf{1h}}^{t,d}$ is the corresponding noisy observations for sample $\mathbf{h}^d$ in the $l$-th LDGEC network. $\boldsymbol{\theta}^t$ is the required learnable variables in

TABLE I Complexity of different channel estimators

| Estimators | LDGEC | SSD | BEACHES | OMP |
|---|---|---|---|---|
| Complexity | $\mathcal{O}(MN^3)$ | $\mathcal{O}(MN_{RF}QL^2\Omega^2)$ | $\mathcal{O}(MNlog(N))$ | $\mathcal{O}(NMN_{RF}QL\Omega)$ |

the $t$-th layer. After training the first to $t$-th layers, a new $t+1$ layer is appended to the LDGEC network and the entire network is trained again for $D$ mini-batches. Although the objective function is changed, the values of the variables $\boldsymbol{\theta}^0$, ..., $\boldsymbol{\theta}^{t-1}$ of the previous round are taken as the initial ones in the optimization process for the new round. In summary, the layer-by-layer training updates variables $\boldsymbol{\theta}^t$ in a sequential manner from the first layer to the last layer.

*E. Complexity analysis*

The computational complexity required for the LDGEC network in each layer is dominant by matrix inverse in lines $5$ and $13$ in Algorithm 1. Generally, the computational complexity of matrix inverse is $\mathcal{O}((MN)^3)$, which cannot be acceptable. As the matrix $\mathbf{A}$ is a block diagonal matrix and can be expressed as

$$\mathbf{A} = (\mathbf{I} \otimes \bar{\mathbf{W}}) = \begin{bmatrix} \bar{\mathbf{W}} & & \\ & \ddots & \\ & & \bar{\mathbf{W}} \end{bmatrix}. \tag{32}$$

Therefore, the matrix inverse in lines $5$ and $13$ in algorithm 1 can be computed by matrix inverse with respect to matrix $\bar{\mathbf{W}}$. By inversion of the partitioned matrix, the total complexity can be reduced from $\mathcal{O}((MN)^3)$ to $\mathcal{O}(MN^3)$. We compare the complexity of LDGEC network with other CS-based algorithms. As shown in Table I, the computational complexity of LDGEC is $\mathcal{O}(MN^3)$, while SSD and BEACHES algorithms have the complexity of $\mathcal{O}(MN_{RF}QL^2\Omega^2)$ and $\mathcal{O}(MNlog(N))$, respectively. Furthermore, the computational complexity of OMP is $\mathcal{O}(MN_{RF}QL^2\Omega^2)$. Generally, $\Omega$ is the beamspace windows and assumed as $\Omega = 4$ when $N = 256$ [31], and is much smaller than $N$. Although LDGEC network has higher computational complexity than CS-based algorithms but can achieve better performance in simulation results.

Fig. 5. The LDGEC-based channel estimator with denoiser-by-denoiser training.

## IV. SUPERVISED LEARNING FOR BEAMSPACE CHANNEL ESTIMATION

In this section, we investigate the LDGEC-based channel estimator for the mmWave beamspace MIMO systems with supervised learning. After presenting three training methods for the LDGEC network, we briefly introduce the denoiser-by-denoiser training.

### A. Denoiser-by-denoiser training

The LDGEC-based channel estimator can also be trained with supervised learning. We introduce three supervised training methods for the LDGEC-based channel estimator with supervised learning as follows,

- *End-to-end training*: We can train all the weights of the $T$ layer LDGEC network simultaneously end-to-end. This is the standard method for training a neural network but with high training complexity.

- *Layer-by-layer training*: We can train the LDGEC with layer-by-layer way by utilizing the MSE loss in (29). The training process is introduced in Section III.

- *Denoiser-by-denoiser training*: We can decouple the denoisers from the rest of the network and train the AWGN denoising problems at different noise levels.

The principle of the *denoiser-by-denoiser* is to train the DnCNN denoiser for the denoising problem solely instead of including the whole GEC algorithm into network training, thereby reducing the computational time in the training stage. Note that the DnCNN denoiser in the LDGEC network is trained for a specific noise level interval. As equivalent noise variance

$\hat{\sigma}_t^2 = \text{avg}(\mathbf{v_{1h}})$ is different for each layer in the LDGEC, we need to deploy different DnCNN denoisers for different layers. To address the issue, we decouple the denoisers from the rest of the network and train each on an AWGN denoising problem at different noise levels. In particular, we scale noise level $\hat{\sigma}^2$ by multiplying 255 as $\bar{\sigma}^2 = 255\hat{\sigma}^2$ and divide $\bar{\sigma}^2$ into intervals [0,10), [10,20), [20,40), [40,60), [60,80), [80,100), [100,150), [150,300), [300,500). For each noise interval, we generate noise variance $\bar{\sigma}^2$ uniformly and train a corresponding DnCNN denoiser.

After training the DnCNN denoiser, we deploy the trained DnCNN denoiser into the LDGEC network to perform channel estimation. As illustrated in Fig. 5, we use a selector to choose the corresponding DnCNN denoiser according to the equivalent noise variance $\hat{\sigma}_t$ for each layer, e.g., we use the denoiser for noise standard deviations between 40 and 60, if $\hat{\sigma}_t^2 * 255 = 55$.

## B. MMSE optimal performance

In [41], the layer-by-layer and denoiser-by-denoiser training for LDAMP are proven to achieve MMSE optimal performance under the following three conditions are satisfied,

- The elements of matrix $\mathbf{A}$ are i.i.d. Gaussian (or sub-gaussian) with zero mean and standard deviation $1/M_1$, where $M_1$ is the number of rows of $\mathbf{A}$.
- The noise, $\mathbf{n}$, is also i.i.d. Gaussian.
- The denoisers, $D_{\hat{\sigma}^t}(\cdot)$, at each layer are *Lipschitz continuous*[2].

Even if the theoretical results have proved that the denoiser-by-denoiser training is optimal, the numerical results in [41] show that LDAMP trained with denoiser-by-denoiser performs slightly worse than the end-to-end and layer-by-layer trained networks due to the discretization of the noise levels ignored in our theory. This gap can be reduced by using a finer discretization of the noise levels or by using deeper denoiser networks to handle a range of noise levels. Although matrix $\mathbf{A}$ in system model (10) is a block diagonal matrix, rather than a Gaussian matrix, we try to use the denoiser-by-denoiser training method for the LDGEC network and show the numerical results in the following section.

---

[2]A denoiser is said to be $L$-Lipschitz continuous if for every $\mathbf{x}_1$ and $\mathbf{x}_2$ we have $\|D(\mathbf{x}_1) - D(\mathbf{x}_2)\| \leq L\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2$ and the Lipschitz continuity of the convolutional denoiser can be ensured by using weight clipping and gradient norm penalization method [43].

| Simulation parameters | Value |
|---|---|
| Number of Paths (L) | 3 |
| Number of antennas (N) | 32 |
| Number of RF chains ($N_{RF}$) | 8 |
| Carrier frequency ($f_c$) | 28 GHz |
| Bandwidth ($f_s$) | 4 GHz |
| Number of subcarriers (M) | 64, 128 |
| Complex gain ($\beta_l$) | $\mathcal{N}_{\mathbb{C}}(0, 1)$ |
| Angle ($\theta_l$) | $\mathcal{U}(-\pi/2, \pi/2)$ |
| Maximum Delay ($\tau_{\max}$) | 20 ns |
| Delay ($\tau_l$) | $\mathcal{U}(0, \tau_{\max})$ |

TABLE II Simulation parameters

## V. NUMERICAL RESULTS

In this section, we provide numerical results to show the performance of the proposed model-driven DL network for wideband beamspace channel estimation. First, we elaborate on the implementation details. Then, the performances of the LDGEC-based channel estimator trained with denoiser-by-denoiser and layer-by-layer are presented. Finally, we investigate the performance of the LDGEC-based channel estimator with a reduced number of RF chains and low-resolution ADCs.
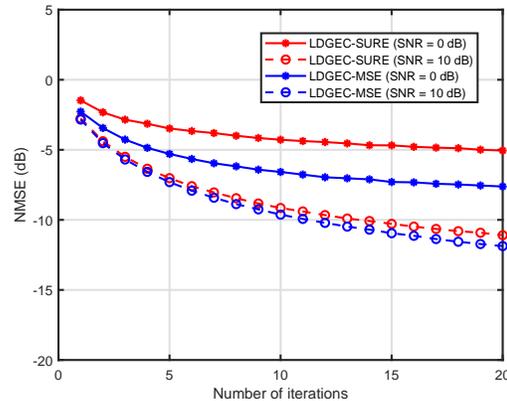
### A. Implementation details

The simulation parameters are listed in Table II. We use the normalized MSE (NMSE) to quantify the accuracy of channel estimation for each user, which is defined as

$$\text{NMSE} = \mathbb{E}\left\{ \|\hat{\mathbf{h}}^{\text{out}} - \mathbf{h}\|_2^2 / \|\mathbf{h}\|_2^2 \right\}, \tag{33}$$

In our simulation, the DL-based channel estimation network is implemented in Tensorflow by using a PC with GPU NVIDIA GeForce GTX 1080 Ti. The training, validation, and testing sets contain 19200, 6400, and 12800 samples, respectively, and are obtained from the Saleh-Valenzuela channel model in (1). The batch size is set to 16, and epoch is equal to 50. We generate the same adaptive selection network, $\mathbf{W}$, for the channel sample in each batch, which is generated independently for different batches. The LDGEC network is trained using the stochastic gradient descent method and the Adam optimizer. The training rate is set to be 0.001 initially

(a) Denoiser-by-denoiser training

(b) Layer-by-layer training

Fig. 6.   Convergence analysis of the LDGEC network.

and then dropped to $0.0001$. As existing deep learning APIs are mostly devoted to processing the real-valued data, we consider equivalent real-valued representation for the system model in (10). We set damping factor $\beta = 0.8$ except addition notes. The code will be available at https://github.com/hehengtao/LDGEC.

## B. Convergence analysis

*1) With real channel data*: Fig. 6(a) investigates the convergence of the LDGEC network with denoiser-by-denoiser training. SNR = $0$, $10$, and $15$ dB are considered. From the figure, the LDGEC network with denoiser-by-denoiser training converges within $6$ layers, and more layers are required when the SNR is increasing.

*2) Without real channel data*: Fig. 6(b) demonstrates the convergence of the LDGEC-based channel estimator with layer-by-layer training. Specifically, the LDGEC-MSE means training the LDGEC network with MSE loss function while LDGEC-SURE indicates training LDGEC network with SURE loss. From Fig. 6(b), the NMSE performance of LDGEC-SURE is close to that of LDGEC-MSE when SNR = $10$ dB. On the contrary, the performance gap is approximately $2.5$ dB when SNR = $0$ dB because the estimate of equivalent noise variance $v_{1h}$ and Monte-Carlo approximation of divergence (24) are not accurate enough in the low-SNR regime, thereby degrading the channel estimation performance.
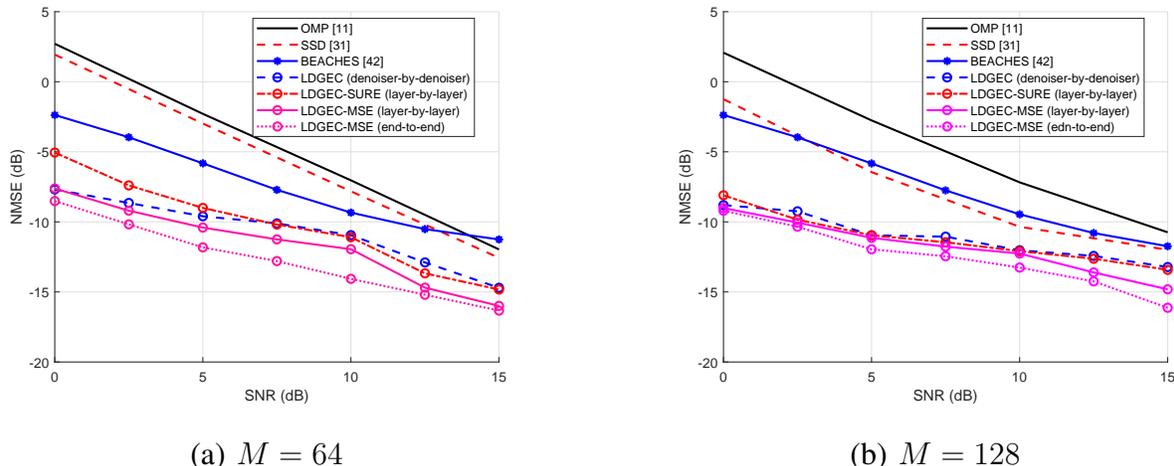
(a) $M = 64$  (b) $M = 128$

Fig. 7. NMSEs performance comparisons of the LDGEC network with other channel estimation algorithms.

## C. Performance comparison

Fig. 7 compares the performance of the LDGEC network with other channel estimation algorithms. For LDGEC network with denoiser-by-denoiser, we set the number of layers $T = 20$ for all SNR. For the LDGEC-SURE and LDGEC-MSE, we set the number of layers $T = 20$ for SNR $\leq 10$ dB while $T = 40$ for SNR $> 10$ dB, because the LDGEC with layer-by-layer training needs more layers to converge in higher SNR. $M = 64$ and $128$ are considered in the simulation. From the figure, the LDGEC-based channel estimator can outperform the traditional CS-based algorithms with different training methods, such as OMP [5], SSD [31], BEACHES [42]. Note that the LDGEC with layer-by-layer training can outperform that with denoiser-by-denoiser training if the MSE is considered as the loss function because we need to divide equivalent noise variance $\hat{\sigma}^2$ into several intervals and train one DnCNN denoiser for each interval in denoiser-by-denoiser training, respectively. Instead of using the coarse intervals, the layer-by-layer training employs the accurate equivalent noise variance estimate, $v_{1h}$, in each layer, thereby improves the denoising performance.

## D. Impact of measurement ratio

In Section II, the measurement ratio, defined by $\delta = QN_{RF}/N$ and involved in the number of RF chains, $N_{RF}$, and pilot length $Q$, influences the performance of channel estimation and

(a) Denoiser-by-denoiser training
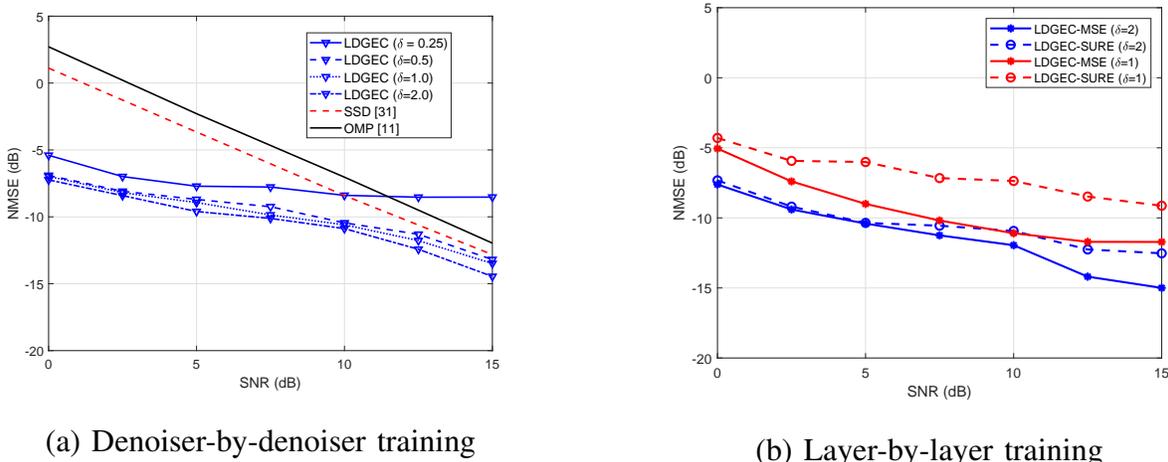
(b) Layer-by-layer training

Fig. 8. NMSE performance of LDGEC network with different measurement ratios for wideband beamspace mmWave MIMO systems.

related to the system overhead. Fig. 8(a) illustrates the performance of the LDGEC network with denoiser-by-denoiser training versus different measurement ratios. Since SSD algorithm cannot work when $\delta \leq 1$, we consider $\delta = 2$. From the figure, the performance of the LDGEC network improves as the measurement ratio increases. Interestingly, the LDGEC algorithm with $\delta = 1$ outperforms the SSD algorithm with $\delta = 2$. Furthermore, the performance of the LDGEC network with $\delta = 1$ is close to that with $\delta = 2$, which demonstrates the strong robustness to the reduced number of RF chains. As the measurement ratio is determined by $N_{RF}$ and $Q$, we can decrease the number of RF chains $N_{RF}$ by increasing the number of pilot length $Q$, which can reduce hardware cost and power consumption of the system significantly. Fig. 8(b) illustrates the performance of the LDGEC network with layer-by-layer training versus different measurement ratios. From the figure, we have similar conclusions to that of LDGEC with denoiser-by-denoiser training.

### E. Low-resolution ADC

The lens-based beamspace mmWave system can decrease the hardware cost by reducing the number of RF chains. However, a common limitation of the architectures is that the receiver RF chains include high-resolution ADCs, which are power-hungry devices, especially when large bandwidth is involved. The power consumption of a typical ADC roughly scales linearly with the bandwidth and grows exponentially with the quantization bits [44]. Many researchers have

studied the mmWave massive MIMO systems with low-resolution ADCs [32], [45], [46]. In this subsection, we investigate the LDGEC-based wideband beamspace channel estimation with low-resolution ADCs.



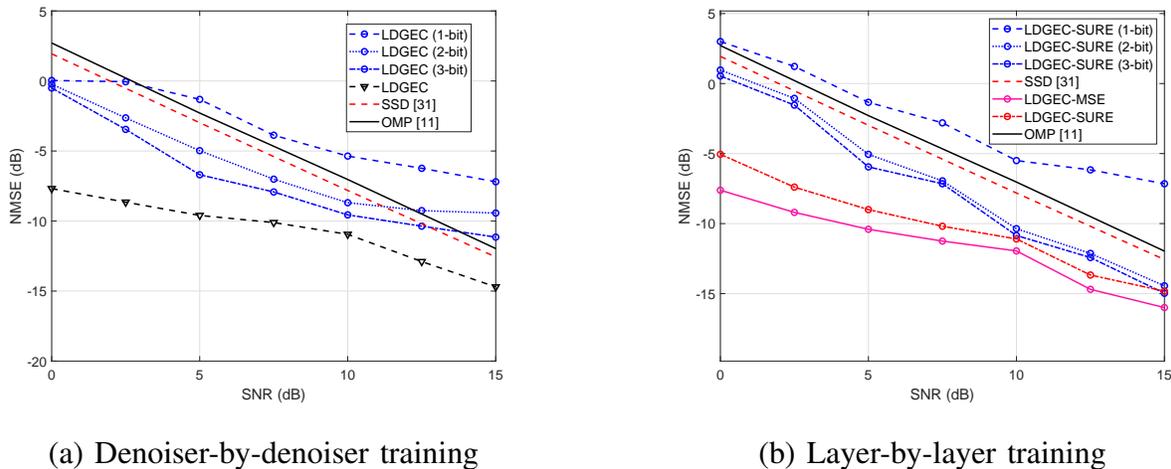(a) Denoiser-by-denoiser training



(b) Layer-by-layer training

Fig. 9. NMSE performance of LDGEC network for wideband beamspace mmWave MIMO systems with low-resolution ADC.

To improve the robustness of LDGEC network in quantized systems, we use the damping method presented in [33] where damping factor $\beta = 0.1^t$ is exponentially decreased. Fig. 9(a) compares the performance of the LDGEC channel estimator with low-resolution ADCs. From the figure, the LDGEC channel estimator with two-bit ADCs outperforms the SSD algorithm with infinite-bit ADCs when $\mathrm{SNR} < 10$ dB. Therefore, the LDGEC channel estimator can accurately estimate the channel from the quantized signal, thereby reducing the hardware cost of the systems.

Fig. 9(b) illustrates the performance of the LDGEC network with SURE loss and layer-by-layer training. From the figure, we have similar conclusions to that of LDGEC with denoiser-by-denoiser training. Furthermore, the LDGEC-SURE with two-bit ADCs outperforms the SSD algorithm with infinite-bit ADCs. Therefore, the LDGEC-SURE is a promising approach to perform beamspace channel estimation in mmWave systems with low-resolution ADC architecture even without real channel data.

## VI. CONCLUSION

We have developed a novel model-driven unsupervised DL network for wideband mmWave beamspace channel estimation. This network inherits the superiority of iterative signal recovery algorithms and the advanced DL-based denoiser, and thus presents excellent performance. The LDGEC network is easy to train and can be applied to a variety of selection networks. Furthermore, By utilizing the SURE loss, the LDGEC network can be trained without real channel data, enables the system to apply in a new environment. Simulation results demonstrate that the LDGEC-based channel estimator significantly outperforms state-of-the-art CS-based algorithms even for the receiver is equipped with a small number of RF chains and low-resolution ADCs. For future work, it will be interesting to apply the SURE technology to other DL-based wireless communication applications, such as CSI feedback and hybrid beamforming problems, to achieve unsupervised learning.

## REFERENCES

[1] A. L. Swindlehurst, E. Ayanoglu, P. Heydari, and F. Capolino, "Millimeter-wave massive MIMO: The next wireless revolution?" *IEEE Comm. Mag.*, vol. 52, no. 9, pp. 56–62, Sep. 2014.

[2] T. Bai and R. W. Heath, "Coverage and rate analysis for millimeter-wave cellular networks," *IEEE Trans. Wireless Comm.*, vol. 14, no. 2, pp. 1100–1114, Feb. 2015.

[3] R. W. Heath, N. Gonz¨¢lez-Prelcic, S. Rangan, W. Roh, and A. M. Sayeed, "An overview of signal processing techniques for millimeter wave mimo systems," *IEEE J. Sel. Topics Signal Process.*, vol. 10, no. 3, pp. 436–453, Apr. 2016.

[4] J. Brady, N. Behdad, and A. Sayeed, "Beamspace MIMO for millimeterwave communications: System architecture, modeling, analysis, and measurements," *IEEE Trans. Antennas Propag.*, vol. 61, no. 7, pp. 3814–3827, Jul. 2013.

[5] A. Alkhateeb, O. El Ayach, G. Leus, and R. W. Heath, "Channel estimation and hybrid precoding for millimeter wave cellular systems," *IEEE J. Sel. Topics Signal Process.*, vol. 8, no. 5, pp. 831–846, Oct. 2014.

[6] X. Gao, L. Dai, S. Han, C.-L. I, and X. Wang, "Reliable beamspace channel estimation for millimeter-wave massive MIMO systems with lens antenna array," *IEEE Trans. Wireless Commun.*, vol. 16, no. 9, pp. 6010–6021, Sept. 2017.

[7] J. Yang, C. Wen, S. Jin, and F. Gao, "Beamspace channel estimation in mmWave systems via cosparse image reconstruction technique," *IEEE Trans. Commun.*, vol. 66, no. 10, pp. 4767–4782, Oct. 2018.

[8] J. Hogan and A. Sayeed, "Beam selection for performance-complexity optimization in high-dimensional MIMO systems," in *Proc. Annu. Conf. Inf. Sci. Syst.*, Mar. 2016, pp. 337–342.

[9] L. Yang, Y. Zeng, and R. Zhang, "Efficient channel estimation for millimeter wave MIMO with limited RF chains," in *Proc. IEEE Int. Conf. Commun.*, May 2016, pp. 1–6.

[10] 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.

[11] H. He, S. Jin, C.-K. Wen, F. Gao, G. Y. Li, and Z. Xu, "Model-driven deep learning for physical layer communications," *IEEE Wireless Commun.*, vol. 26, no. 5, pp. 77–83, Oct. 2019.

[12] C.-K. Wen, W. T. Shih, and S. Jin, "Deep Learning for Massive MIMO CSI Feedback," *IEEE Wireless Commun. Lett.*, vol. 7, no. 5, pp. 748–751, Oct. 2018.

[13] H. Ye, G. Y. Li, and B.-H. F. 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.

[14] H. He, C.-K. Wen, S. Jin, and G. Y. Li, "Model-driven deep learning for MIMO detection," *IEEE Trans. Signal Process.*, vol. 68, pp. 1702–1715, Mar. 2020.

[15] H. He, C. K. Wen, S. Jin, and G. Y. Li, "Deep learning-based channel estimation for beamspace mmWave massive MIMO systems," *IEEE Wireless Commun. Lett.*, vol. 7, no. 5, pp. 852–855, Oct. 2018.

[16] C. Qi, P. Dong, W. Ma, H. Zhang, Z. Zhang, and G. Y. Li, "Acquisition of channel state information for mmWave massive MIMO: traditional and machine learning-based approaches," *arXiv preprint arXiv:2006.08894*, Jun. 2020.

[17] 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.

[18] W. Xia, G. Zheng, Y. Zhu, J. Zhang, J. Wang, and A. Petropulu, "A deep learning framework for optimization of MISO downlink beamforming," *IEEE Trans. Commun.*, vol. 68, no. 3, pp. 1866–1880, Mar. 2020.

[19] G. Aceto, D. Ciuonzo, A. Montieri, and A. Pescape, "Mobile encrypted traffic classification using deep learning: Experimental evaluation, lessons learned, and challenges," *IEEE Trans. Netw. Service Manage.*, vol. 16, no. 2, pp. 445–458, Jun. 2019.

[20] A. Rago, G. Piro, G. Boggia, and P. Dini, "Multi-task learning at the mobile edge: An effective way to combine traffic classification and prediction," IEEE Trans. Veh. Technol., vol. 69, no. 9, pp. 10362–10374, Sep. 2020.

[21] G. Aceto et al., "MIMETIC: Mobile encrypted traffic classification using multimodal deep learning," *Comput. Netw.*, vol. 165, Dec. 2019, Art. no. 106944.

[22] F. Aoudia and J. Hoydis, "End-to-End Learning of Communications Systems Without a Channel Model," *arXiv preprint arXiv:1804.02276*, Dec. 2018.

[23] H. Ye, L. Liang, G. Y. Li, and B.-H. F. Juang, "Deep learning based end-to-end wireless communication systems with GAN as unknown channel," *IEEE Trans. Wireless Commun.*,, vol. 19, no. 5, pp. 3133–3143, May. 2020.

[24] J. R. Hershey, J. Le Roux, and F. Weninger, "Deep unfolding: Model-based inspiration of novel deep architectures," *arXiv preprint arXiv:1409.2574*, Sep. 2014.

[25] D. Ito, S. Takabe, and T. Wadayama, "Trainable ISTA for sparse signal recovery," *IEEE Trans. Signal. Process.*, vol. 67, no. 12, pp. 3113-3125, Jun. 2019.

[26] V. Monga, Y. Li, and Y. C. Eldar, "Algorithm unrolling: Interpretable, efficient deep learning for signal and image processing," *arXiv preprint arXiv:1912.10557*, Dec. 2019.

[27] A. Balatsoukas-Stimming and C. Studer, "Deep unfolding for communications systems: A survey and some new directions," *arXiv preprint arXiv:1906.05774*, Jun. 2019.

[28] E. Balevi, A. Doshi, and J. G. Andrews, "Massive MIMO channel estimation with an untrained deep neural network," *IEEE Trans. Wireless Commun.*, vol. 19, no. 3, pp. 2079–2090, Mar. 2020.

[29] B. Wang, F. Gao, S. Jin, H. Lin, and G. Y. Li, "Spatial-and frequency wideband effects in millimeter-wave massive MIMO systems," *IEEE Trans. Signal Process.*, vol. 66, no. 13, pp. 3393–3406, Jul. 2018.

[30] B. Wang, F. Gao, S. Jin, H. Lin, G. Y. Li, S. Sun, and T. S. Rappaport, "Spatial-wideband effect in massive MIMO with application in mmWave systems," *IEEE Commun. Mag.*, vol. 56, no. 12, pp. 134–141, Dec. 2018.

[31] X. Gao, L. Dai, S. Zhou, A. M. Sayeed, and L. Hanzo, "Wideband beamspace channel estimation for millimeter-wave MIMO systems relying on lens antenna arrays," *IEEE Trans. Signal Process.*, vol. 67, no. 18, pp. 4809–4824, Sep. 2019.

[32] H. He, C.-K. Wen, and S. Jin, "Bayesian optimal data detector for hybrid mmwave MIMO-OFDM systems with low-resolution ADCs," *IEEE J. Sel. Topics Signal Process.*, vol. 12, no. 3, pp. 469–483, Jun. 2018.

[33] C. Wang, C.-K. Wen, S. Tsai, and S. Jin, "Decentralized expectation consistent signal recovery for phase retrieval," *IEEE Trans. Signal Process.*, vol. 68, pp. 1484–1499, 2020.

[34] A. Saleh and R. Valenzuela, "A statistical model for indoor multipath propagation," *IEEE J. Sel. Areas Commun.*, vol. 5, no. 2, pp. 128–137, May 1987.

[35] C.-K. Wen, C.-J. Wang, S. Jin, K.-K. Wong, and P. Ting, "Bayes-optimal joint channel-and-data estimation for massive MIMO with low-precision ADCs," *IEEE Trans. Signal Process.*, vol. 64, no. 10, pp. 2541–2556, May. 2016.

[36] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, "Beyond a Gaussian denoiser: Residual learning of deep CNN for image denoising," *IEEE Trans. Image Process.*, vol. 16, no. 9, pp. 3142–3155, Jul. 2017.

[37] C. Tian, L. Fei, W. Zheng, Y. Xu, W. Zuo, and C.-W. Lin, "Deep learning on image denoising: An overview," *arXiv preprint arXiv:1912.13171*, 2019.

[38] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2016, pp. 770–778.

[39] C. M. Stein, "Estimation of the mean of a multivariate normal distribution," *The Annals of Statistics*, pp. 1135–1151, 1981.

[40] C. A. Metzler, A. Mousavi, R. Heckel, and R. G. Baraniuk, "Unsupervised learning with stein's unbiased risk estimator," *arXiv preprint arXiv:1805.10531*, 2018.

[41] C. A. Metzler, Ali Mousavi, and R. G. Baraniuk, "Learned D-AMP: principled neural network based compressive image recovery," *arxiv preprint arXiv:1704.06625*, 2017.

[42] S. H. Mirfarshbafan, A. Gallyas-Sanhueza, R. Ghods, and C. Studer,"Beamspace channel estimation for massive MIMO mmWave systems: Algorithm and VLSI design," *arXiv preprint arXiv:1910.00756*, 2017.

[43] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. Courville, "Improved training of Wasserstein GANs," *arXiv preprint arXiv:1704.00028*, 2017.

[44] H.-S. Lee and C. G. Sodini, "Analog-to-digital converters: Digitizing the analog world," *Proc. IEEE*, vol. 96, no. 2, pp. 323–334, 2008.

[45] J. Zhang and L. Dai, Z. He, S. Jin, and X. Li, "Performance analysis of mixed-ADC massive MIMO systems over Rician fading channels," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 6, pp. 1327–1338, Jun. 2017.

[46] J. Mo, P. Schniter, and R. W. Heath, "Channel estimation in broadband millimeter wave MIMO systems with few-bit ADCs," *IEEE Trans. Signal Process.*, vol. 66, no. 5, pp. 1141–1154, Mar. 2018.