# Overview of Deep Learning-based CSI Feedback in Massive MIMO Systems

Jiajia Guo, *Graduate Student Member, IEEE*, Chao-Kai Wen, *Senior Member, IEEE*,
Shi Jin, *Senior Member, IEEE*, and Geoffrey Ye Li, *Fellow, IEEE*

*Abstract*—Many performance gains achieved by massive multiple-input and multiple-output depend on the accuracy of the downlink channel state information (CSI) at the transmitter (base station), which is usually obtained by estimating at the receiver (user terminal) and feeding back to the transmitter. The overhead of CSI feedback occupies substantial uplink bandwidth resources, especially when the number of the transmit antennas is large. Deep learning (DL)-based CSI feedback refers to CSI compression and reconstruction by a DL-based autoencoder and can greatly reduce feedback overhead. In this paper, a comprehensive overview of state-of-the-art research on this topic is provided, beginning with basic DL concepts widely used in CSI feedback and then categorizing and describing some existing DL-based feedback works. The focus is on novel neural network architectures and utilization of communication expert knowledge to improve CSI feedback accuracy. Works on bit-level CSI feedback and joint design of CSI feedback with other communication modules are also introduced, and some practical issues, including training dataset collection, online training, complexity, generalization, and standardization effect, are discussed. At the end of the paper, some challenges and potential research directions associated with DL-based CSI feedback in future wireless communication systems are identified.

*Index Terms*—CSI feedback, massive MIMO, deep learning, overview.

## I. INTRODUCTION

**T**HE 3rd Generation Partnership Project (3GPP) completed the first release of the fifth generation (5G) mobile communications, namely, Release 15, in 2018, laying the foundation for the global commercial deployment of 5G [1]. The three major usage scenarios of 5G networks include enhanced mobile broadband (eMBB) to ultra-reliable low-latency communications (URLLC) to massive machine type communications. To support the use cases, some novel techniques, including millimeter-wave transmission, network densification, and massive multiple-input and multiple-output (MIMO), have been introduced [2]. 3GPP has been working on 5G evolution in Releases 16 and 17 to enhance existing features further and enable new use cases [3], [4]. With Releases 17 specification work ongoing, 3GPP also started the plan for 5G-Advanced and recently approved the package including 27 work items in Release 18 [5]. In particular, the features of Release 18 work

Jiajia Guo and Shi Jin are with the National Mobile Communications Research Laboratory, Southeast University, Nanjing, 210096, P. R. China (email: jiajiaguo@seu.edu.cn, jinshi@seu.edu.cn).

Chao-Kai Wen is with the Institute of Communications Engineering, National Sun Yat-sen University, Kaohsiung 80424, Taiwan (e-mail: chaokai.wen@mail.nsysu.edu.tw).

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

(a) Autoencoder for image compression
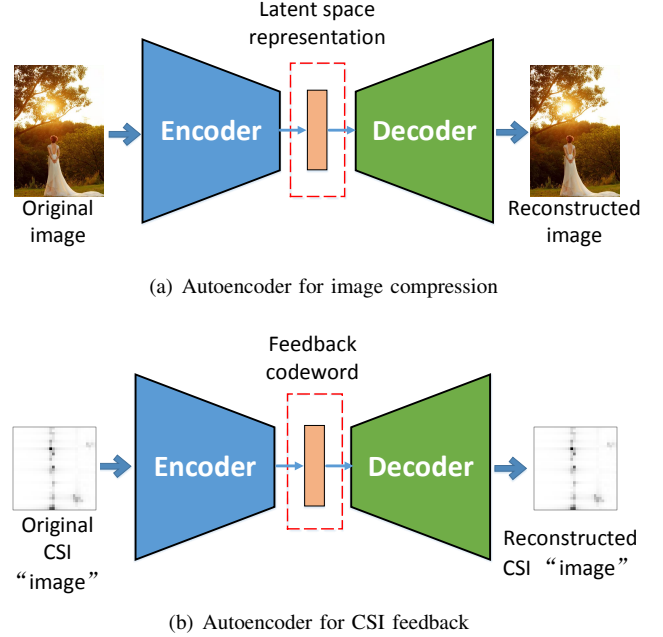


(b) Autoencoder for CSI feedback

Fig. 1. Illustration of autoencoder architectures. In image compression, the NN-based encoder compresses the original image into a low-dimensional representation and then the NN-based decoder reconstructs the image from the latent representation. The encoder and decoder are jointly trained. In the right sub-figure, the downlink CSI is regarded as a special type of "image".

on embracing artificial intelligence (AI) and machine learning (ML) technologies. Release 18 is expected to pave the way toward integrating AI and communications. MIMO evolution is one of the key features in 3GPP Release 18.

In the past 10 years, deep learning (DL) has achieved great success in many areas. Inspired by its success, DL has been introduced in wireless communications [6]–[9]. The DL technology can be used to enhance the conventional communication blocks. In [10] and [11], deep neural networks (NNs) are used to design a downlink beamforming (BF) matrix. In [12] and [13], DL is introduced for channel estimation and symbol detection, which has been validated by an over-the-air test in [14]. Furthermore, DL enables end-to-end communication systems [15], [16], where the transmitter and the receiver are represented by NNs in the form of autoencoder. The concept of end-to-end communication systems is validated in [17].

In massive MIMO, the base station (BS) is equipped with a large number (up to a few hundred) of active antennas and simultaneously serves multiple users. The knowledge of accurate channel state information (CSI) at the BS is essential to obtaining the performance gains of massive MIMO [18],
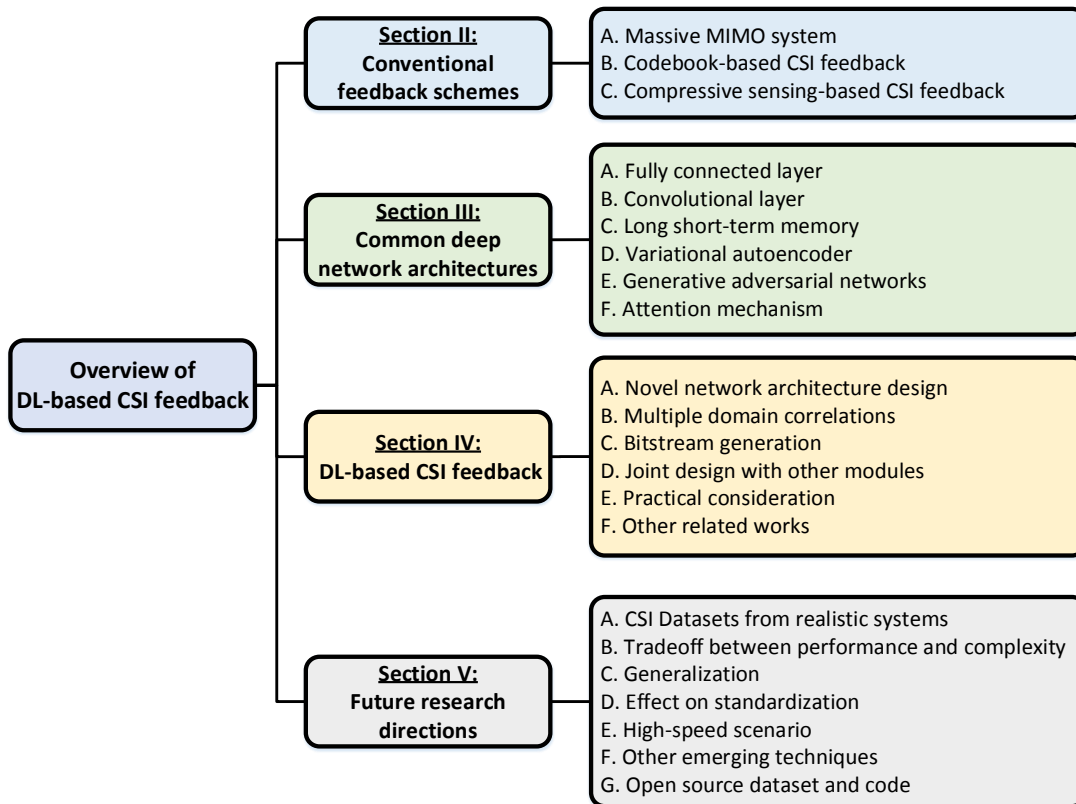
Fig. 2. Outline of article

[19]. Downlink CSI acquisition contains two main steps. First, the user estimates the downlink CSI utilizing the received pilot signals transmitted by the BS. Then, the user feeds the estimated downlink CSI back to the BS through the uplink control channel. In massive MIMO systems, a large number of antennas at the BS result in a vast CSI dimension and require a substantial feedback overhead. In addition, commercial deployments in 5G have observed that the user often experiences considerable performance loss due to the outdated CSI fed back by the user. Conventional CSI feedback methods are based on codebook [20] and compressive sensing (CS) [21], which cannot meet the requirement of low complexity and high accuracy. Therefore, potential CSI feedback enhancements are explored to improve the performance of massive MIMO systems. 3GPP Release 18 will study AI/ML for this use case [22].

CSI compression and feedback in massive MIMO can be also based on DL [23]. The common architecture adopted in DL-based feedback borrows the idea of the autoencoder used in image compression in Fig. 1. Fig. 1(a) shows that the encoder compresses the original image by the NNs to generate latent space representation. The dimension of this latent representation is much lower than that of the original image. Then, the NN-based decoder reconstructs the original image from the received latent representation. The NN-based encoder and decoder in data compression are trained in an end-to-end manner. The autoencoder-based image compression has substantially outperformed the conventional compression

techniques. Fig. 1(b) shows that the CsiNet framework in [23] regards downlink CSI as a special type of "image." The encoder at the user compresses the downlink CSI. The compressed CSI is then fed back to the transmitter. Upon obtaining the feedback information, the decoder at the BS reconstructs the CSI by NNs.

This paper provides an overview of DL-based CSI feedback in massive MIMO systems. It focuses on feedback performance, NN complexity, and the effect on other communication. First, the conventional CSI feedback methods, including codebook-based and CS-based feedback algorithms, are briefly introduced, and their main limitations are discussed. Then, a brief introduction to the basic concepts of the NNs and some representative NN architectures, which are widely used in the existing CSI feedback works, is provided[1]. Next, the existing DL-based feedback works are divided into six different categories, namely, the introduction of novel NN design, the utilization of multi-domain correlations, bitstream generation, joint design with other modules, practical considerations, and others. Finally, the main challenges of DL-based CSI feedback, especially in 3GPP standardization, are discussed, and the future research direction is identified.

The article is outlined in Fig. 2. Section II presents system models and some conventional feedback methods in massive MIMO systems. Section III describes basic NN concepts and representative architectures widely used in DL-based CSI

---

[1]This part can be skipped if the reader has a good understanding of basic DL concepts.

feedback. Section IV overviews existing works including the motivations, key ideas, and weaknesses. Section V illustrates the main challenges and the corresponding future research directions. Section VI concludes this paper.

*Notations:* In this paper, italic letters represent scalars. Bold-face lower-case and upper-case letters denote vectors and matrices, respectively. $\mathbb{C}^{m \times n}$ ($\mathbb{R}^{m \times n}$) denotes the space of $m \times n$ complex-valued (real-valued) matrices. $\mathbf{I}$ represents an identity matrix. The transpose, conjugate, Hermitian transpose, and inverse operations are represented by $(\cdot)^{\mathrm{T}}$, $(\cdot)^{*}$, $(\cdot)^{\mathrm{H}}$, and $(\cdot)^{-1}$, respectively. $\mathrm{Tr}(\cdot)$ and $E(\cdot)$ denote the trace and the expectation of a matrix, respectively. The Euclidean norm of a vector is written as $\| \cdot \|$. $\mathrm{round}(\cdot)$ represents the rounding operation. $[\mathbf{A}]_{i,j}$ represents the $(i,j)$-th element of matrix $\mathbf{A}$.

## II. CONVENTIONAL FEEDBACK SCHEMES

In this section, some representative conventional CSI feedback methods are presented. After introducing the fundamental signal model of massive MIMO systems, the basic ideas and the pros and cons of different methods are discussed.

### A. Massive MIMO System

For simplicity, a simple single-cell massive MIMO system operated in orthogonal frequency-division multiplexing mode with $N_{\mathrm{c}}$ subcarriers is considered. The BS is equipped with a uniform linear antenna array (ULA) with $N_{\mathrm{t}}$ ($\gg 1$) transmit antennas. $K$ users each have a single receive antenna. If the BS adopts a linear precoding algorithm, such as zero-forcing (ZF), the transmit signal from the BS at the $n$-th subcarrier will be

$$\mathbf{x}_n = \sum_{k=1}^{K} \mathbf{v}_{n,k} s_{n,k} = \mathbf{V}_n \mathbf{s}_n, \tag{1}$$

where $\mathbf{v}_{n,k} \in \mathbb{C}^{N_{\mathrm{t}} \times 1}$ and $s_k$ denote the linear precoding vector for the $k$-th user and the transmitted symbol of the $k$-th user, respectively, and $\mathbf{V}_n = [\mathbf{v}_{n,1}, \ldots, \mathbf{v}_{n,K}]$. The whole precoding matrix and the transmitted symbol should satisfy the power constraints as $\mathrm{Tr}(\mathbf{V}_n \mathbf{V}_n^{\mathrm{H}}) \leq P$ and $E(\mathbf{s}_n \mathbf{s}_n^{\mathrm{H}}) = \mathbf{I}$, respectively. The received signal at the $k$-th user over the $n$-th subcarrier can be expressed as:

$$y_{n,k} = \mathbf{h}_{n,k}^{\mathrm{T}} \mathbf{v}_{n,k} s_{n,k} + \sum_{i \neq k} \mathbf{h}_{n,k}^{\mathrm{T}} \mathbf{v}_{n,i} s_{n,i} + z_{n,k}, \tag{2}$$

where $\mathbf{h}_{n,k} \in \mathbb{C}^{N_{\mathrm{t}} \times 1}$ is the frequency response vector at the $n$-th subcarrier, and $z_{n,k} \sim \mathcal{CN}(0, \sigma^2)$ denotes the complex additive white Gaussian noise with zero mean and variance $\sigma^2$. The BS designs precoding matrix $\mathbf{V}_n$ for the $n$-th subcarrier using the entire downlink CSI matrix of all users, $\mathbf{H}_n = [\mathbf{h}_{n,1}, \ldots, \mathbf{h}_{n,K}]$. For example, the ZF precoding matrix can be expressed as

$$\mathbf{V}_n = c \mathbf{H}_n^{*} (\mathbf{H}_n^{\mathrm{T}} \mathbf{H}_n^{*})^{-1}, \tag{3}$$

where $c = \sqrt{P / \|(\mathbf{H}_n^{\mathrm{T}} \mathbf{H}_n^{*})^{-1}\|^2}$ is the power normalization factor.
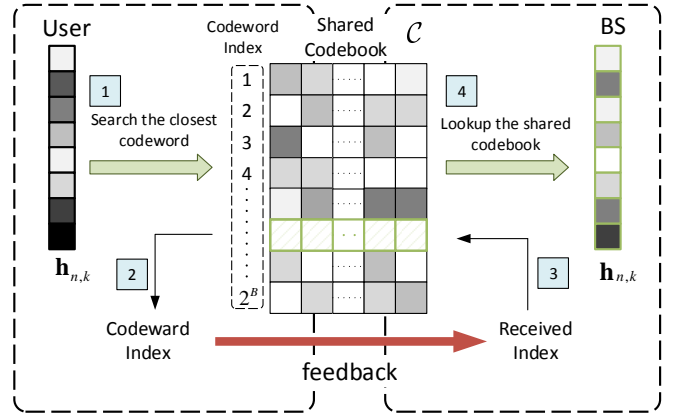


Fig. 3. Illustration of codebook-based CSI feedback. The codebook is known to the user and the BS. The user searches the codeword, which is the closest to the downlink CSI, and feeds back the corresponding index to the BS. Upon receiving the index, the BS can obtain the channel by looking up the shared codebook.

### B. Codebook-based CSI Feedback

In Fig. 3, random vector quantization (RVQ) [24] is used as an example to introduce the key idea of codebook-based feedback briefly. Assuming that the feedback bit number is $B$, the CSI codebook, $C$, shared by the BS and the user/user equipment (UE) contains $2^B$ $N_{\mathrm{t}}$-dimensional unit norm isotropic distributed vectors. The codeword for the downlink CSI $\mathbf{h}_{n,k}$ can be obtained by

$$\hat{\mathbf{h}}_{n,k} = \arg \max_{\mathbf{u} \in C} \|\mathbf{h}_{n,k}^{\mathrm{H}} \mathbf{u}\|. \tag{4}$$

The user feeds back the index of the selected codeword $\mathbf{u}$ to the BS via the uplink control channel. The BS obtains the corresponding codeword based on the received index.

The codebook-based CSI feedback has faced some challenges. Feedback accuracy is improved with codebook size $2^B$. For example, the TYPE II codebook in 5G new radio (NR) remarkably outperforms the TYPE I codebook at the expense of a substantial increase in feedback bit number. In addition, codeword search complexity increases with codebook size.

Although many algorithms, such as an adaptive codebook, have been proposed [20], the performance of feedback accuracy, complexity, and overhead of channel codebook needs to be improved further.

### C. CS-based CSI Feedback

The CSI matrix is sparse in certain domains, such as time, spatial, spatial-temporal, and spatial-frequency domains [21]. CS can be used to reduce the overhead of downlink CSI. Given that the number of scatter clusters is much smaller than that of the transmit antennas at the BS in massive MIMO systems, the CSI matrix can be represented by much fewer parameters, and the spatial domain turns into the sparse angular domain using discrete Fourier transform (DFT) as

$$\tilde{\mathbf{h}}_{n,k} = \mathbf{F} \mathbf{h}_{n,k}, \tag{5}$$

where $\mathbf{F} \in \mathbb{C}^{N_{\mathrm{t}} \times N_{\mathrm{t}}}$ stands for a DFT matrix. CSI compression at the user is implemented via a sensing matrix as

$$\mathbf{m} = \mathbf{\Phi} \tilde{\mathbf{h}}_{n,k}, \tag{6}$$

where $\mathbf{\Phi} \in \mathbb{C}^{M \times N_t}$ and $\mathbf{m}$ stand for the sensing matrix and the compressed measurement vector, respectively. To ensure a high-accuracy reconstruction at the BS, sensing matrix $\mathbf{\Phi}$ should satisfy the restricted isometry property.

Assuming that at most $p$ elements in the vector $\tilde{\mathbf{h}}_{n,k}$ are nonzero, the original high-dimension CSI vector can be accurately recovered from the measurement $\mathbf{m}$ when $N_t \gg M$ and $M > p$. The reconstruction problem can be formulated as the following minimization task

$$\min_{\tilde{\mathbf{h}}_{n,k}} \|\tilde{\mathbf{h}}_{n,k}\|_0, \quad \text{s.t.} \quad \mathbf{m} = \mathbf{\Phi}\tilde{\mathbf{h}}_{n,k}. \tag{7}$$

This optimization can be solved by the iterative algorithms, such as iterative shrinkage thresholding algorithm (ISTA) [25]. However, some challenges hinder the deployment of CS-based feedback: The CS-based CSI feedback is based on the sparsity assumption of CSI, which may not hold in practical systems. The complexity of the iterative reconstruction algorithms is too high to meet the real-time requirements.

## III. COMMON DEEP NETWORK ARCHITECTURES

In this section, common deep network architectures, including fully connected (FC) layers, convolutional layers, long short-term memory networks (LSTMs), variational autoencoder (VAE), generative adversarial networks (GAN), and attention mechanism, which are widely adopted in DL-based CSI feedback works, are briefly introduced.

### A. FC Layer

The FC layer is the vanilla NN layer, in which all input neurons are connected to all output neurons. This layer first multiplies the input vector by a weight matrix and then adds a bias vector, which can be formulated as

$$\mathbf{y}'_{\text{FC}} = \mathbf{W}_{\text{FC}}\mathbf{x}_{\text{FC}} + \mathbf{b}_{\text{FC}}, \tag{8}$$

where $\mathbf{x}_{\text{FC}} \in \mathbb{R}^{N_{\text{in}} \times 1}$ and $\mathbf{y}'_{\text{FC}} \in \mathbb{R}^{N_{\text{out}} \times 1}$ stand for the input and the output vectors, respectively, and $\mathbf{W}_{\text{FC}}$ and $\mathbf{b}_{\text{FC}}$ represent the weight matrix and the bias vector of the FC layer, respectively. This operation is linear and in real numbers. Following this operation, the activation function is implemented as

$$\mathbf{y}_{\text{FC}} = \sigma(\mathbf{y}'_{\text{FC}}) = \sigma(\mathbf{W}_{\text{FC}}\mathbf{x}_{\text{FC}} + \mathbf{b}_{\text{FC}}), \tag{9}$$

where $\sigma(\cdot)$ represents the non-linear activation function. In DL-based CSI feedback, the commonly used activation functions are Tanh, Sigmoid, ReLU, and Leaky ReLU functions, as follows:

$$Tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}, \tag{10}$$

$$Sigmoid(z) = \frac{1}{1 + e^{-z}}, \tag{11}$$

$$ReLU(z) = \max(0, z), \tag{12}$$

$$LeakyReLU(z) = \begin{cases} z & z \geq 0, \\ az & z < 0, \end{cases} \tag{13}$$

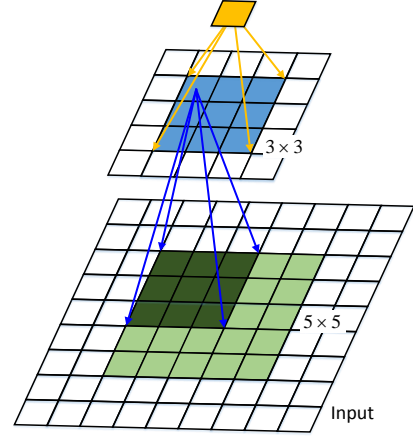where $a \in (0, 1)$ is a hyperparameter, namely, negative slope [26].



Fig. 4. Receptive field illustration of two stacked $3 \times 3$ convolutional layers. Stride $s$ is set as 1. The upper "pixel" is determined by the $3 \times 3$ square area in the middle. Each intermediate "pixel" is determined by the input $3 \times 3$ square area, which is overlapped with one another. Therefore, the upper "pixel" is determined by the $5 \times 5$ square area of the input.

The FC layer is widely used to extract features from the input vectors in computer vision. In addition to feature extraction, the FC layer can be used to change the dimension of the input vector. For example, in CsiNet [23], the last layer at the encoder and the first layer at the decoder are FC layers, which adjust the vector dimension by changing the neuron number of the FC output, that is, $N_{\text{out}}$.

NN complexity is evaluated by two metrics: the number of NN parameters and the number of floating point operations (FLOPs). The parameter of the FC layer can be calculated by

$$N_{\text{FC}} = N_{\text{in}} \times N_{\text{out}} + N_{\text{out}} \approx N_{\text{in}} \times N_{\text{out}}. \tag{14}$$

FLOP number [2] can be obtained by

$$F_{\text{FC}} = 2 \times N_{\text{in}} \times N_{\text{out}}. \tag{15}$$

From (14) and (15), the complexity of FC layers increases with the dimensions of the input and the output vectors.

### B. Convolutional Layer

The convolutional layer is composed of a linear convolution operation, which encompasses the multiplication of a set of weights with the input similar to sliding a filter across an input vector. The convolutional layer can learn features with invariance to shifts in the input [15]. The NN parameter number is substantially reduced in comparison with that of the FC layers.

Assuming that the convolutional layer is composed of $C_{\text{out}}$ filter weights $\mathbf{Q}_c \in \mathbb{R}^{a \times b}$ for $c = 1, \ldots, C_{\text{out}}$. $\mathbf{Q}^c$ is multiplied by input $\mathbf{x} \in \mathbb{R}^{w \times h}$ to generate a feature map $\mathbf{Y}^c$, which can be achieved by the convolution operation as

$$[\mathbf{Y}^c]_{i,j} = \sum_{m=0}^{a-1} \sum_{n=0}^{b-1} [\mathbf{Q}^c]_{a-m,b-n} [\mathbf{x}]_{1+s(i-1)-m,1+s(j-1)-n}, \tag{16}$$

where $s \geq 1$ is an integer hyperparameter named as stride. If input matrix $\mathbf{x}$ is padded with $\big((w + a - 1)(h + b - 1) - wh\big)$

---

[2]In this paper, the FLOPs number of the activation function is neglected.

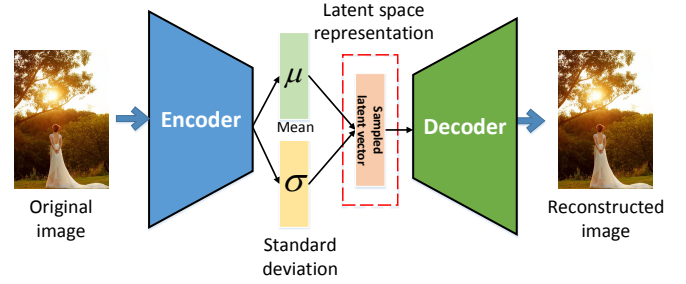Fig. 5. Illustration of an LSTM cell [30] that has feedback connections and allows information to persist



Fig. 6. Illustration of VAE architecture that encodes the input as a distribution over latent space instead of a single point like that in an autoencoder

zeros, that is, $[\mathbf{x}]_{i,j} = 0$ for all $i \notin [1, w]$ and $j \notin [1, h]$, the dimension of the feature map $\mathbf{Y}^c$ is $(1 + \lfloor \frac{w+a-2}{s} \rfloor) \times (1 + \lfloor \frac{l+b-2}{s} \rfloor)$. Assuming that the depth of the layer's input is $C_{\text{in}}$, the parameter number of this layer is

$$N_{\text{conv}} = (a \times b \times C_{\text{in}} + 1) \times C_{\text{out}}. \tag{17}$$

The number of FLOPs can be obtained by

$$F_{\text{conv}} = \left(2 \times C_{\text{in}} \times (a \times b)\right) \times w \times h \times C_{\text{out}}. \tag{18}$$

In convolutional NNs (CNNs), several and even hundreds of convolutional layers are stacked to extract the input features, such as ResNet-152 [27], increasing the receptive field of convolutional layers. The receptive field stands for the size of the region in the input, which produces the features. Fig. 4 shows an example with two stacked convolutional layers, whose filter sizes are $3 \times 3$ and strides are 1. The upper "pixel" is determined by the $3 \times 3$ square area in the middle. Each "pixel" in the middle is determined by "pixels" in the down $3 \times 3$ square area. The $3 \times 3$ filter slides across the whole matrix. Therefore, the upper "pixel" value is determined by the "pixels" in the $5 \times 5$ input area, and the receptive field of this NN block is $5 \times 5$.

The size of the receptive field is essential to the NN performance. As pointed out in [28], a logarithmic relationship exists between receptive field size and the accuracy of classification tasks. RepLKNet with $31 \times 31$ convolutional kernels [29] outperforms the most existing NNs in computer vision tasks. In DL-based feedback, many existing works focus on reducing feedback errors by adjusting the NN receptive field.

### C. LSTMs

The traditional NNs, such as FC and convolutional layers, make decisions only based on the current information without leveraging any previous information. However, in certain domains, such as consecutive frames in a video current, information is highly correlated with the previous one. Recurrent neural networks (RNNs) allow previous outputs to be used as inputs and can address this issue.

LSTM [31], which allows information to persist, is a widely used RNN architecture as shown in Fig. 5. The figure shows that an LSTM cell contains the input, forget, and output gates, namely, $\mathbf{i}_t$, $\mathbf{f}_t$, and $\mathbf{o}_t$, respectively, where $\mathbf{x}_t$, $\mathbf{C}_{t-1}$, and $\mathbf{h}_{t-1}$ represent the input to the LSTM cell, the state of the past cell, and the output of the past LSTM cell, respectively. The goal of the input gate $\mathbf{i}_t$ is to decide whether the input ( $\mathbf{x}_t$ and $\mathbf{h}_{t-1}$ ) is needed to be stored in the cell and drop the unwanted information. The forget gate, $\mathbf{f}_t$, determines whether to drop the previous state information $\mathbf{C}_{t-1}$ based on the input data. The state information of this cell, $\mathbf{C}_t$, is updated based on the forget and the input gates. The output of this cell, $\mathbf{h}_t$, is decided by the input data and the current state, $\mathbf{C}_t$. More details about LSTMs can be found in [30]. LSTM has many different variants and topologies, such as bidirectional LSTM (Bi-LSTM), and gated recurrent unit (GRU), referring to [32] for a tutorial.

### D. VAE

Fig. 1(a) shows that an autoencoder can realize a high dimension reduction with a high reconstruction accuracy. However, it cannot be used for the generative tasks because of the lack of structure among the latent vectors. To solve this problem, a variant of the autoencoder, namely, VAE, is introduced in [33]. VAE encodes the input as a distribution over the latent space instead of a single point like that in an autoencoder. The output of the encoder is the mean vector $\boldsymbol{\mu}$ and standard deviation vector $\boldsymbol{\delta}$. A point from latent space is sampled from a predefined distribution as

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\delta}\boldsymbol{\epsilon}, \tag{19}$$

where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Then, latent representation vector $\mathbf{z}$ is sent to the decoder to reconstruct the original data.

During the training of the autoencoder, mean-squared error (MSE) between the input and the output of the autoencoder is widely used as the loss function. However, the training of VAE has two main objectives: reconstructing the input and ensuring the latent vector $\mathbf{z}$ to be normally distributed. Therefore, its loss function is the sum of the reconstruction loss and the similarity loss. Reconstruction loss is the MSE loss used in autoencoder whereas similarity loss is the Kullback-Leibler divergence between standard Gaussian distribution and latent space distribution.

### E. GAN

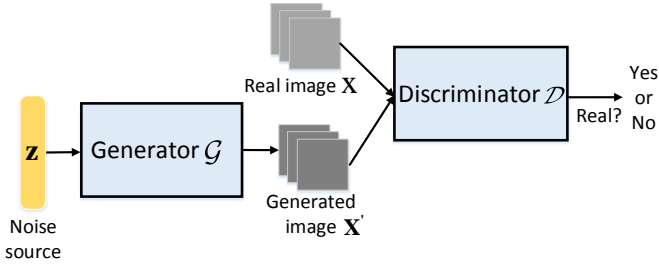The GAN framework [34], as shown in Fig. 7, is a class of DL-based generative models. The GAN framework consists

Fig. 7. Illustration of GAN architecture, including a generator, $\mathcal{G}$, and a discriminator, $\mathcal{D}$



Fig. 8. Illustration of an attention-based framework [39], including channel and spatial attention modules

of two NN-based submodels, namely, a generator $\mathcal{G}$ and a discriminator $\mathcal{D}$. The generator, $\mathcal{G}$, generates new plausible examples in the problem domain, whereas the discriminator, $\mathcal{D}$, classifies whether the input examples are real or generated by the generator. The two submodels compete with each other as in a game.

The training of GANs is based on a game-theoretic scenario, in which the generator, $\mathcal{G}$, competes against an adversary, that is, the discriminator, $\mathcal{D}$. Concretely, the two modules need to be trained jointly with two opposite goals at the same time:

- The training target of generator $\mathcal{G}$ is to fool discriminator $\mathcal{D}$, that is, maximizing the final classification error.
- The training target of discriminator $\mathcal{D}$ is to detect the fake examples generated by the generator $\mathcal{G}$, that is, minimizing the final classification error.

The opposite training goals force the two modules to try to beat each other, thereby simultaneously improving their performances. The final equilibrium state of the GAN training corresponds to the situation where the generator, $\mathcal{G}$, can generate data from the targeted distribution and the discriminator, $\mathcal{D}$, predicts "real" or "fake" with a probability 50% for all received examples.

*F. Attention Mechanism*

RNNs, LSTM, and GRU are typically established for processing sequential data, for example, language modeling and machine translation. Such networks consider computations along the symbolic positions of the input and output sequences. The attention mechanism allows dependencies to be modeled regardless of their distance in the input or output sequence, which has shown to achieve remarkable performance improvements [35]. The attention mechanism is first applied to the natural language processing domain to address the long sequence problem in the machine translation task [36] and has been extended to other applications, such as computer vision [37]. From a cognitive science perspective, humans only notice a portion of all visible information due to bottlenecks in information processing [38]. Bottlenecks inevitably exists in NNs' processing. Therefore, researchers propose a visual selective attention model, that is, attention mechanism, to simulate the visual perception of humans. The key problem of the attention mechanism is how to find the key features to capture long-range information interactions.

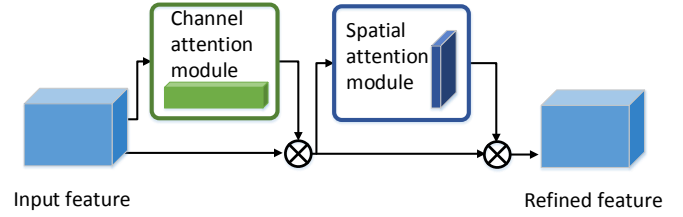In computer vision, the key features are identified by a mask, which quantifies the importance of each pixel or channel. The mask is not handcrafted but learned by another NN layer with new parameters. Fig. 8 shows an NN framework (called CBAM) [39], which consists of channel and spatial attention modules that are widely adopted in computer vision. Channel attention focuses on determining which feature map is meaningful when the input feature has tens or hundreds of channels. The spatial dimension of the input feature is squeezed by maximum- or average-pooling to generate a 1D vector, which is then forwarded to an NN module to produce a spatial attention map. Then, all input features maps are multiplied by the corresponding weight in the generated attention map. Spatial attention focuses on determining where the features are informative.

## IV. DL-BASED CSI FEEDBACK

The existing research directions in DL-based CSI feedback can be divided into six categories. The first three categories, which include novel NN architecture design, multi-domain correlation utilization, and bitstream generation, focus on improving the performance of DL-based CSI feedback. The remaining three categories, which include joint design with other modules, practical consideration, and other related works, focus on promoting the practical deployment of DL-based CSI feedback.

*A. Novel NN Architecture Design*

Conventional ML requires careful domain expertise to design a feature extractor. The main advantage of DL is that features can be learned from substantial training samples using an end-to-end approach, that is, manually designing feature extractors is not needed. However, NN architecture heavily affects the performance of DL-based algorithms and should be carefully designed. Table I shows the normalized MSE (NMSE) performance of the feedback NNs using the dataset published by [23]. Compression ratio (CR) represents the ratio of the codeword dimension to the original CSI dimension. If CR is 1/16 for outdoor channels, the performance gap between the CsiNet and the state-of-the-art NN is over 10 dB, which shows the great effect of the NN architecture on feedback accuracy.

The existing NN design works of CSI feedback can be divided into seven categories, as shown in Table II. Their key ideas are introduced, and a guideline for the future NN design of CSI feedback is provided.

TABLE I
NMSE (dB) PERFORMANCE OF THE NNs USING THE DATASET PUBLISHED BY [23].

| CR | 1/4 | | 1/8 | | 1/16 | | 1/32 | | 1/64 | |
|---|---|---|---|---|---|---|---|---|---|---|
| Scenarios | Indoor | Outdoor | Indoor | Outdoor | Indoor | Outdoor | Indoor | Outdoor | Indoor | Outdoor |
| CsiNet (Dec. 2017) [23] | −17.36 | −8.75 | \ | \ | −8.65 | −4.51 | −6.24 | −2.81 | −5.84 | −1.93 |
| ConvCsiNet (2018) [40] | −17.37 | −8.98 | \ | \ | −13.79 | −6.00 | −10.10 | −5.21 | **−7.72** | −4.48 |
| CsiNet+ (June 2019) [41] | −27.37 | −12.40 | −18.29 | −8.72 | −14.14 | −5.73 | −10.43 | −3.40 | −6.72 | −2.45 |
| Attention-CsiNet (Oct. 2019) [42] | −20.29 | −10.43 | \ | \ | −10.16 | −6.11 | −8.58 | −4.57 | −6.32 | −3.27 |
| CRNet (Oct. 2019) [43] | −26.99 | −12.71 | −16.01 | −8.04 | −11.35 | −5.44 | −8.93 | −3.51 | −6.49 | −2.22 |
| LSTM-Attention CsiNet (Jan. 2020) [44] | −22.00 | −10.20 | \ | \ | −11.00 | −5.80 | −8.80 | −3.70 | −7.20 | −2.40 |
| DS−NLCsiNet (Aug. 2020) [45] | −24.99 | −12.09 | −17.00 | −7.96 | −12.93 | −4.98 | −8.64 | −3.35 | \ | \ |
| DCGAN (Aug. 2020) [46] | −26.20 | −15.88 | \ | \ | −13.50 | −8.07 | −9.00 | −5.83 | −6.45 | −4.01 |
| PRVNet (Nov. 2020) [47] | −27.70 | −13.90 | \ | \ | −13.00 | −6.10 | −9.52 | −4.23 | −6.90 | −2.53 |
| CF-FCFNN (Jan. 2021) [48] | −20.07 | −11.61 | −15.14 | −10.08 | −12.35 | −9.12 | −8.86 | −8.42 | −6.60 | −7.25 |
| CLNet (Feb. 2021) [49] | −29.16 | −12.88 | −15.60 | −8.29 | −11.15 | −5.56 | −8.95 | −3.49 | −6.34 | −2.19 |
| ENet (May 2021) [50] | −26.00 | \ | \ | \ | −14.50 | \ | **−11.20** | \ | −7.50 | \ |
| DCRNet (June 2021) [51] | −30.61 | −13.72 | −19.92 | −10.17 | −14.02 | −6.35 | −9.88 | −3.95 | \ | \ |
| CsiNet+DNN (June 2021) [52] | \ | \ | \ | −17.88 | \ | −14.70 | \ | −14.42 | \ | −11.34 |
| MRFNet (July 2021) [53] | −25.76 | **−15.95** | \ | \ | −14.72 | −9.49 | −10.63 | −7.42 | −6.90 | −6.52 |
| ACCsiNet (Sep. 2021) [54] | \ | \ | \ | \ | −14.59 | −11.76 | −11.00 | −9.14 | −7.46 | −7.11 |
| DFECsiNet (Dec. 2021) [55] | −27.50 | −12.25 | −16.80 | −7.90 | −12.70 | −5.20 | −8.85 | −3.35 | −5.95 | −2.10 |
| TransNet (Feb. 2022) [56] | **−32.38** | −14.86 | **−22.91** | −9.99 | **−15.00** | −7.82 | −10.49 | −4.13 | −6.08 | −2.62 |
| CsiFormer (Feb. 2022) [57] | \ | \ | \ | \ | \ | \ | −9.32 | −3.51 | −6.85 | −2.25 |
| CVLNet (March 2022) [58] | \ | \ | \ | \ | −13.97 | −6.67 | −9.72 | −4.56 | \ | \ |

"\" means the performance is not reported. The methods are ordered by their publication time.

*1) Increasing Receptive Field:* The CsiNet architecture [23], shown in Fig. 9, is the first work that applies DL to CSI feedback. In this architecture, a convolutional layer is first employed at the encoder to extract the feature of the downlink CSI "images," and the filter size is set as $3 \times 3$, which is the smallest in the commonly used ones. The receptive field size is $3 \times 3$, too. At the decoder, the RefineNet, which consists of three series convolutional layers with $3 \times 3$ filters and adopts residual learning [71], is employed to refine the initially reconstructed CSI. The receptive field size of each RefineNet is $7 \times 7$, which is much smaller than the CSI "image" size, that is, $32 \times 32$, in [23].

As mentioned in Section III-B, the performance of CNN heavily depends on the size of the receptive field. Inspired by this observation, CsiNet+ [41] improves feedback performance by enlarging the receptive field size. From [41], the $3 \times 3$ receptive field, which is widely used to extract the edge information, is not suitable for the CSI feedback task. A convolutional filter with a small receptive field cannot make full use of the CSI sparsity in the angular-delay domain. Therefore, a convolutional layer with a much larger receptive field is adopted in CsiNet+ architecture. Two convolutional layers with $7 \times 7$ kernel sizes replace the original convolutional layer at the encoder, and the receptive field sizes of the first two convolutional layers in RefineNet are set as $7 \times 7$ and $5 \times 5$. This modification improves the performance of the original CsiNet. For example, NMSE is reduced from −17.36 dB to −20.80 dB when CR is 1/4. Based on [52], CsiNet+ does not work well when CR is low, such as 1/32 for outdoor channels. Therefore, two FC layers are embedded after the second convolutional layer in the RefineNet block, and more RefineNet blocks are employed. Moreover, the Swish function is used as the activation function. Inspired by [41], $5 \times 5$ and $8 \times 8$ convolution kernels are adopted at the encoder and decoder in [59].

In [60], seven convolutional layers with $3 \times 3$ filters are adopted at the decoder to expand the receptive field to 15, which is half of the CSI size. Compared with directly adopting a $15 \times 15$ convolutional operation, this can greatly reduce the complexity, including the numbers of NN parameters and FLOPs determined by (17) and (18). The receptive field is enhanced in [61] by stacking convolutional layers with $3 \times 3$ filters at the encoder to improve the quality of the features extracted from the input CSI.

*2) Multiple Resolutions:* The above works, such as CsiNet+ [41], focus on enlarging the NN receptive field. However, the CSI sparsity varies in different scenarios and even within different regions of a single CSI sample. As pointed out by [43], larger receptive field (or convolutional kernel) is preferred for sparser regions, and the convolution operation with a small kernel can extract finer features much better. Therefore, the CSI should be processed by the convolutions with different kernel sizes, namely, multiple resolutions.

The CRNet architecture proposed by [43] first introduces the multi-resolution architecture to CSI feedback. The encoder and decoder in CRNet adopt a multi-resolution architecture. Fig. 10 shows the encoder of CRNet and the main block at the decoder. The input CSI "image" passes through two parallel NN paths. The left one consists of three stacked convolutional layers, namely, with $3 \times 3$, $1 \times 9$, and $9 \times 1$ convolution kernels. The resolution (or receptive field size) of this path $11 \times 11$. The right path only consists of a convolutional layer with $3 \times 3$ convolution kernels. The resolution of this path is much smaller than that of the left one. Then, the outputs of two paths with different resolutions are concatenated and merged by a $1 \times 1$ convolution operation. Finally, an FC layer is adopted to reduce the dimension of the CSI. The decoder of CRNet is similar to that of CsiNet and only replaces the RefineNet block with the CRBlock, as shown in Fig. 10. The CRBlock is based on the encoder's NN architecture but uses larger convolution kernels and residual learning. The multiple resolution strategy greatly improves feedback performance. For example, when

TABLE II
NOVEL NN ARCHITECTURE DESIGN

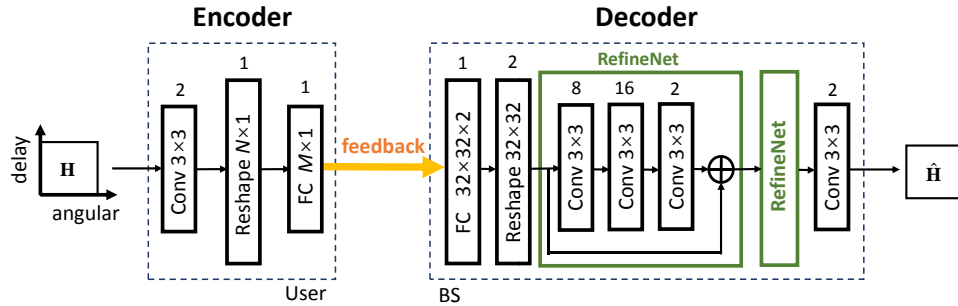| Key ideas | NN name | Main contributions in NN architectures |
|---|---|---|
| **Increasing Receptive Field** | CsiNet+ [41]<br><br>CsiNet+DNN [52]<br><br>MRNet [59]<br>CS-ReNet [60]<br>BCsiNet [61] | Replacing the original convolutional layer ($3 \times 3$) at the encoder by two convolutional layers ($7 \times 7$);<br>Setting the first two convolutional layers in RefineNet as $7 \times 7$ and $5 \times 5$;<br>Embedding two FC layers after the second convolutional layer in the RefineNet block;<br>Employing more RefineNet block;<br>Setting the convolutional sizes of the encoder and the decoder as $5 \times 5$ and $8 \times 8$;<br>Stacking seven convolutional layers with $3 \times 3$ filters at the decoder;<br>Stacking three $3 \times 3$ convolutional layers at the encoder to improve CSI feature quality; |
| **Multiple Resolutions** | CRNet [43]<br><br>DFECsiNet [55]<br>MSMCNet [62]<br><br>MRFNet [53] | The input CSI "image" passes through two parallel NN paths ($11 \times 11$ and $3 \times 3$) at the encoder;<br>RefineNet is replaced with CRBlock that uses larger convolution kernels and residual learning;<br>Two feature extraction paths ($7 \times 7$ and $3 \times 3$) are employed in parallel;<br>The MSMC block has three parallel convolution paths ($5 \times 5$, $7 \times 7$, and $9 \times 9$);<br>The MRFBlock has three parallel paths, with $5 \times 5$, $7 \times 7$, and $9 \times 9$ convolution kernels;<br>The reason for the success of the multiple resolution strategy is revealed via feature visualization; |
| **Fully Convolutional Layer** | ConvCsiNet [40]<br><br>DeepCMC [63]<br>ACCsiNet [54]<br>FullyConv [64] | The dimension reduction of CSI at the encoder is achieved by average pooling operations;<br>The dimension increase is realized by bilinear interpolation at the decoder;<br>The convolution kernels ($9 \times 9$ and $5 \times 5$) adopted are much larger than those in ConvCsiNet;<br>The asymmetric convolution block is adopted;<br>The dimension reduction is realized by a convolution operation by changing the stride; |
| **Attention Mechanism** | Attention-CsiNet [42]<br>SALDR [65]<br>CsiTransformer [66]<br>TransNet [56] | Channel attention modules are introduced to the decoder of the CSI feedback;<br>The encoder adopts a novel spatial attention mechanism, that is, patch-wise self-attention;<br>A single-layer transformer module replaces the traditional convolution operation;<br>A more powerful two-layer transformer architecture is adopted; |
| **GAN and VAE** | DCGAN [46]<br>PRVNet [47] | A GAN architecture is introduced to the training process;<br>VAE is introduced to CSI feedback; |
| **Well-designed Preprocessing** | CsiNet [23]<br><br><br>[67]<br>CLNet [49]<br>ENet [50]<br>P-SRNet [68] | The CSI matrix is first transformed to angular-delay domain;<br>Some rows, whose valuse are all close to zero, are removed;<br>The real and imaginary parts of truncated CSI are concatenated and then scaled in [0, 1];<br>Channel element magnitude is set as $A$ if it is over a predefined threshold $A$;<br>The real and imaginary parts are embedded in physical meaning by a $1 \times 1$ convolution;<br>The real and imaginary parts are fed back to the BS separately with the same autoencoder;<br>The rows full of near-zero elements are omitted; |
| **Others** | TiLISTA-Joint [69]<br><br>FISTA-Net [70]<br><br>CF-FCFNN [48] | The CSI matrix is compressed by an FC layer;<br>The ISTA algorithm is unfolded, and hyperparameter is learned by NNs;<br>The fast ISTA algorithm is unfolded;<br>The original CSI is represented by a basis and a residual part of the column space channel matrix;<br>The feedback NN architecture only consists of FC layers. |



Fig. 9. Illustration of the CsiNet architecture, in which the encoder compresses downlink CSI and the decoder reconstructs CSI from the feedback information. The encoder and decoder consist of convolutional and FC layers.

CR is 1/4 for the outdoor scenario, NMSE can be reduced from $-8.75$ dB to $-12.71$ dB with minimal increase in NN complexity.

An improved NN architecture, called DFECsiNet, is proposed in [55] for CSI feedback. The key component of DFEC-siNet is DFEBlock, in which two feature extraction paths are employed to extract the different resolution diversity of CSI in parallel. The two paths have different resolutions, $7 \times 7$ and $3 \times 3$. Similar to [43], residual learning is adopted in the DFEBlock at the decoder. The NN architecture in [62], called MSMCNet, adopts a novel multi-scale and multi-channel (MSMC) block based on the CRBlock [43]. The MSMC block has three parallel convolution paths with different receptive field sizes, $5 \times 5$, $7 \times 7$, and $9 \times 9$.

The NN architecture in [53], called MRFNet, reveals the reason for the success of the multiple resolution strategy via feature visualization. The encoder in MRFNet is the same as that in CsiNet, and the main modifications are employed to the decoder at the BS. The MRFBlock has three parallel paths, with $5 \times 5$, $7 \times 7$, and $9 \times 9$ convolution kernels. The feature number of each path is 64, which is much smaller than other works. Other architectures of MRFBlock are similar to that of the CRBlock. From feature visualization, different CSI features can be learned by the convolution operations with different kernel sizes. The convolution operation with a small kernel, such as $5 \times 5$, focuses on extracting the background or pattern information. The one with a large kernel, such as $9 \times 9$, is good at extracting values located in distinct regions.
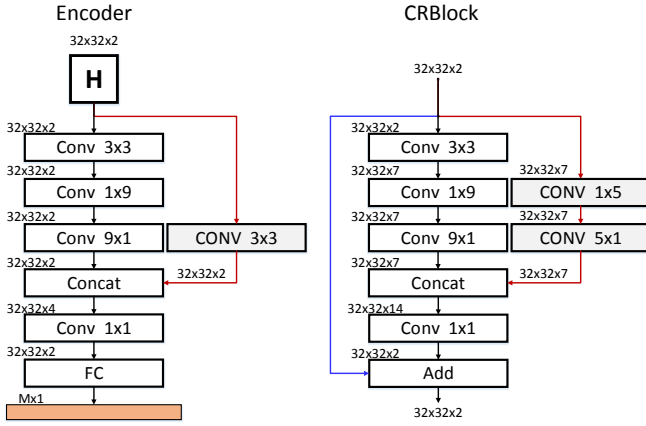
Fig. 10. Encoder and CRBlock of CRNet proposed in [43] with are two parallel paths, each with different resolutions
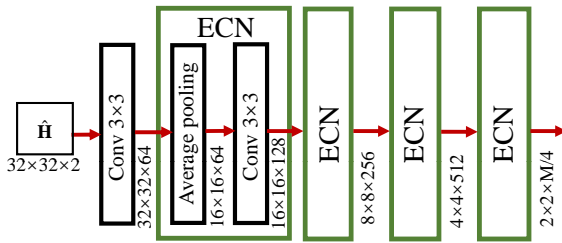


Fig. 11. Encoder architecture in ConvCsiNet [40], in which dimension reduction is achieved by average pooling rather than reducing the neuron number of FC layers

*3) Fully Convolutional Layer:* The existing works extract CSI's features by convolutional layers and compress and reconstruct the CSI by FC layers. Dimension reduction and increase of CSI are achieved by adjusting the neuron numbers of the last FC layer of encoder and the first FC layer of decoder, respectively.

Different from the existing works, ConvCsiNet proposed in [40] is based on convolutional layers without FC layers. Fig. 11 shows the encoder architecture of ConvCsiNet. The key component of the encoder is the encoded convolution network (ECN) block, each of which consists of an average pooling layer and a convolutional layer. The dimension reduction of CSI is achieved by average pooling operations, each of which reduces the CSI size by half. Therefore, four serial ECN reduces the CSI dimension from $32 \times 32$ to $2 \times 2$. The number of feature maps after each convolution operation is quite large, and the last one is $M/4$, where $M$ is the codeword length. The decoder of ConvCsiNet is also based on convolutional layers and the dimension increase is realized by bilinear interpolation. ConvCsiNet [40] can greatly improve CSI feedback accuracy when CR is low, such as $1/32$ and $1/16$. Moreover, fully convolutional architecture is flexible to the dimension of input CSI [63].

Similar to ConvCsiNet, DeepCMC [63] and ACCsiNet [54] compress and reconstruct the CSI by down sampling and up sampling layers, respectively. Moreover, the convolution kernels ($9 \times 9$ and $5 \times 5$) adopted in [63] are much larger than those adopted in [40]. The asymmetric convolution block
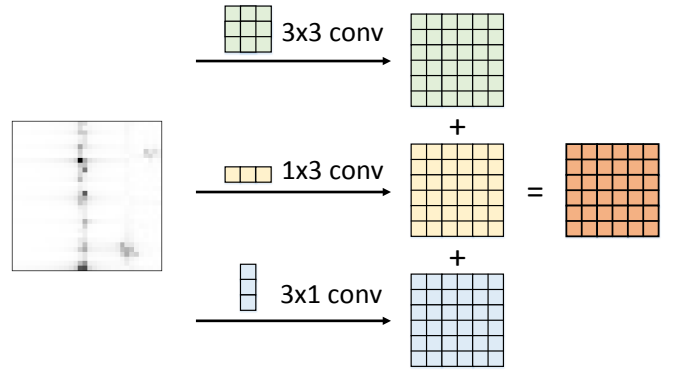


Fig. 12. Illustration of asymmetric convolution block in [54]

[72], as shown in Fig. 12, is introduced to enhance the CSI feature extraction of the convolution operation in [54]. The asymmetric convolution block consists of three parallel layers with $3 \times 3$, $1 \times 3$, and $3 \times 1$ convolution kernels, and the outputs are then summed up. This block enriches the feature space compared with the standard $3 \times 3$ convolution operation.

FullyConv in [64] is also based on fully convolutional layers. Different from [40], [54], [63], the dimension reduction in [64] is realized by convolution operation. In FullyConv, the convolution layer with the stride set as $4 \times 4 \times 4$ can achieve a CR of $1/64$. A transposed convolution (or deconvolution) layer with the same stride can restore the original size of the CSI at the decoder.

*4) Attention Mechanism:* The attention mechanism is first introduced to DL-based CSI feedback by Attention-CsiNet proposed in [42]. As mentioned in Section III-F, the importance of different feature maps is different. Therefore, NN performance can be improved if more attention is paid to the feature maps with more information. Based on this observation, channel attention modules are introduced to the decoder of the CSI feedback in [42], [73], and [74]. Fig. 13 shows the channel attention module and the Attention RefineNet block proposed in [42]. The goal of the attention module is to generate a vector to describe the importance of each feature map. First, a global average pooling is adopted to generate an $L \times 1$ vector. Then, two FC layers are employed to reconstruct the importance vector. The activation function of the last layer is Sigmoid to guarantee that all vector values are in the range (0, 1). Finally, the generated vector is multiplied by the input feature maps. NN can work better because more useful information can be highlighted and extracted with an attention vector. The Attention RefineNet block is similar to the original RefineNet used in CsiNet [23] but introduces an extra attention module.

The encoder in [65] adopts another kind of attention mechanism, that is, spatial attention. The original spatial attention generates attention weights for each "pixel" in the feature map. However, the correlation among the adjacent "pixels" is ignored. To solve this problem, a novel spatial attention mechanism, that is, patch-wise self-attention [75], is adopted in [65]. The key idea of this attention mechanism is to limit the scope of the original spatial attention to a local patch rather
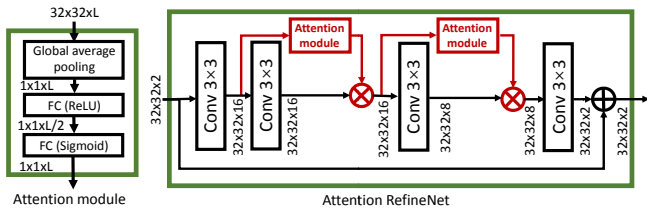
Fig. 13. Attention module and Attention RefineNet architecture proposed in [42]
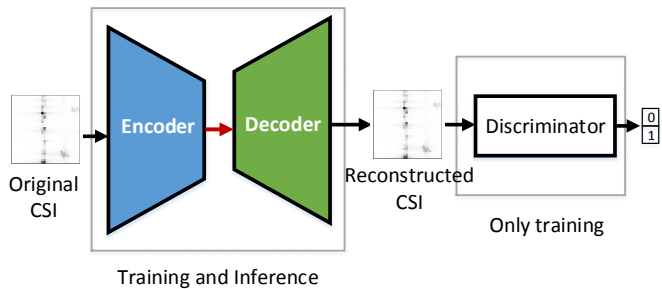


Fig. 14. Illustration of GAN architecture (called DCGAN) proposed in [46]. The autoencoder and discriminator $\mathcal{D}$ are jointly trained. However, the discriminator is not needed during inference.

than the entire feature map, thereby not only further improving feedback performance but also reducing the complexity of the attention module. Moreover, the decoder in [65] uses a novel RefineNet, named dense RefineNet, in which each NN layer passes its feature maps through all subsequent NN layers [76].

The attention mechanism in [42], [65] is based on CNN architecture. However, the transformer architecture [35] completely abandons the traditional CNN and RNN architectures, and only relies on the attention module to eschew recurrence. The transformer architecture is applied to CSI feedback by [66], in which the transformer module replaces the traditional convolution operation at the encoder and the RefineNet at the decoder. However, performance is slightly improved. The authors of [56] stated that CsiTransformer in [66] cannot fully utilize the transformer's power because only a single-layer transformer architecture is used. Therefore, a more powerful two-layer transformer architecture, namely, TransNet, is proposed. On the basis of fully excavating the power of the transformer, feedback performance is greatly improved. Convolutional transformer architecture is adopted by CsiFormer [57] to maintain long-range dependency of CSI.

*5) GAN and VAE:* The training and inference of the above works are based on the autoencoder architecture proposed by CsiNet [23]. Unlike these works, novel NN frameworks, namely, GAN and VAE, are introduced in [46] and [47] to DL-based CSI feedback.

Fig. 14 shows the GAN architecture (called DCGAN) in [46]. An extra discriminator $\mathcal{D}$ is added after the autoencoder. The method can help generate more plausible CSI compared with other DL techniques. During inference, the discriminator is not needed any more, and only the encoder and decoder are deployed to practical systems. This method can improve NN performance without changes to the original autoencoder-based inference architecture. Therefore, it can be easily introduced to the above works.

The novel NN architecture based on VAE, called PRVNet, is introduced in [47]. As in Section III-D, the loss function of the traditional VAE is defined as the sum of reconstruction loss and similarity loss. However, this loss function is not suitable for the DL-based CSI feedback problem. Therefore, reconstruction loss occupies a major position in the whole loss function. A parameter $\beta \in (0,1)$ is introduced in [47] to emphasize the importance of the reconstruction loss as

$$l = l_{\text{rec}} + \beta \cdot l_{\text{dis}}, \tag{20}$$

where $l_{\text{rec}}$ and $l_{\text{dis}}$ represent the reconstruction loss of CSI and the similarity loss of the distribution, respectively. This
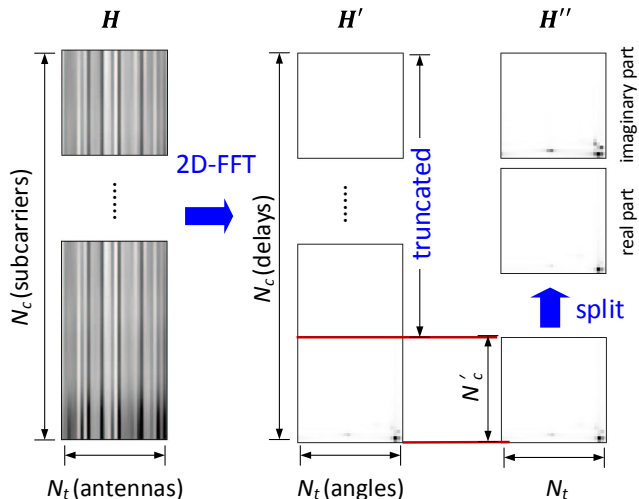


Fig. 15. CSI preprocessing workflow in [23]. Preprocessing consists of three steps: 2D-DFT, truncation, and splitting.

method introduces extra performance improvement compared with the traditional VAE-based CSI feedback.

*6) Well-designed Preprocessing:* Preprocessing is essential in data science. Preprocessing of CSI samples affects the performance of DL-based CSI feedback. Fig. 15 shows the preprocessing in [23], which has been adopted by most existing works. The estimated downlink CSI matrix is first transformed to angular-delay domain by a 2D-DFT operation. The sparsity characteristic of CSI holds in this domain. Given that the time delays between multi-path arrivals lie within a rather limited period, only the first $N_{\text{c}}^{'}$ rows contain values that are not close to zero. Therefore, only the first $N_{\text{c}}^{'}$ rows are retained, and the remaining are removed. DL libraries, such as TensorFlow [77] and PyTorch [78], only can build real-valued NNs. Thus, the real and imaginary parts of the complex truncated CSI are concatenated to formulate a real-valued 3D matrix. Last, each 3D matrix is scaled in [0, 1] by

$$\mathbf{H}_i^{\text{norm}} = \frac{1}{2}\left( \frac{\mathbf{H}_i^{''}}{max\left(abs\left([\mathbf{H}_0^{''}, \dots, \mathbf{H}_K^{''}]\right)\right)} + 1 \right). \tag{21}$$

The activation function of the last NN layer in CsiNet [23] is the Sigmoid defined in (11). Moreover, according to [23],
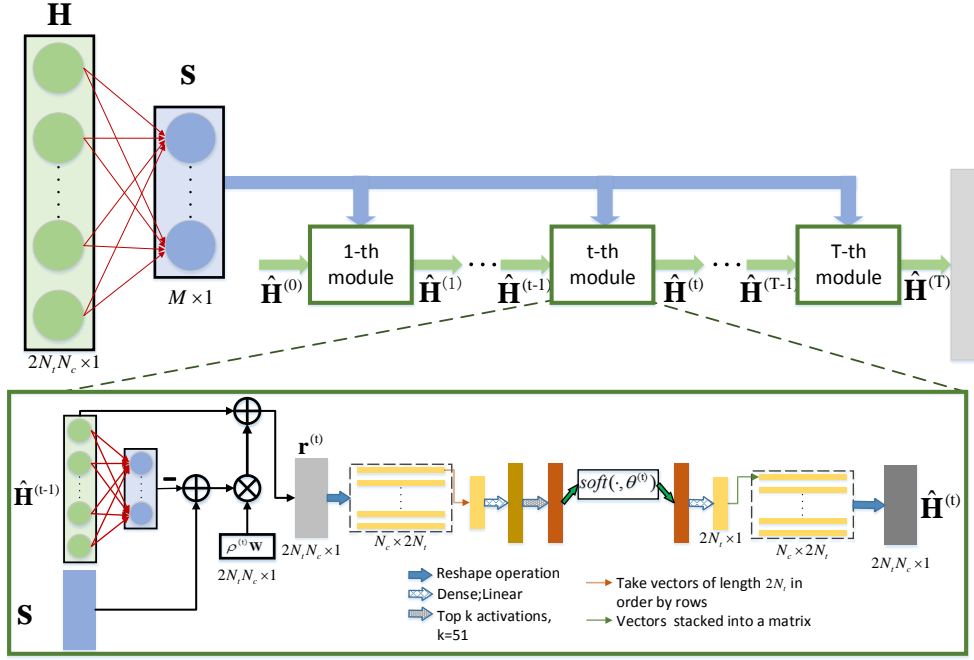
Fig. 16. Entire deep unfolding structure of TiLISTA-Joint in [69]

whether to transform CSI of the spatial domain into the angular domain has no effect on the CsiNet's performance.

Data normalization heavily affects the performance of DL, including accuracy and training complexity [79]. Therefore, a simple yet efficient data normalization method with clipping is proposed in [67] for DL-based CSI feedback. Some channel elements may have very high power, which affects the statistical operation of DL and can be regarded as outliers [67]. Hence, if a channel element has a magnitude over a threshold $A$, its magnitude is set as $A$, and its phase remains the same. Then, the clipped CSI matrix is scaled to [0, 1]. This kind of preprocessing can greatly improve NN performance and accelerate the convergence of NN training.

The CLNet proposed in [49] considers that the CSI is in the form of complex values with its physical meaning. The previous works overlook this problem and directly concatenate the real and imaginary parts of the CSI together, inevitably resulting in performance loss. The first layer in the previous works is a convolutional layer with $3 \times 3$ or larger kernels. Here, a $3 \times 3$ convolution operation is used as an example. $\mathbf{X} \in \mathbb{R}^{m \times m \times 2}$ is a 3D tensor that is extended from its 2D version by concatenation operation. $\mathbf{I} = [\mathbf{i}_1, \ldots, \mathbf{i}_C] \in \mathbb{R}^{m \times m \times C}$ represents the feature maps after convolution operation, and $C$ is the channel number of the feature maps. $a_n + b_n \mathrm{i}$ ($n = 1, 2, 3$) represents a $3 \times 3$ patch in $\mathbf{X}$, and $w_n$ denotes the weight of the convolution operation. The $3 \times 3$ convolution operation on this patch can be essentially formulated as the sum of two multiplication processes as[3]

$$
\begin{aligned}
i_1(1,1) &= [a_1, \ldots, a_9] \cdot [w_1, \ldots, w_9]^{\mathrm{T}} + [b_1, \ldots, b_9] \cdot [w_1, \ldots, w_9]^{\mathrm{T}} \\
&= [a_1 + b_1, \ldots, a_9 + b_9] \cdot [w_1, \ldots, w_9]^{\mathrm{T}}.
\end{aligned} \tag{23}
$$

[3]The bias terms of convolution operation are omitted in this part for simplicity.

From (23), the real and imaginary parts of the neighboring CSI elements are entangled, and nine complex CSI values are interpolated as a synthesized one, resulting in the loss of original physical meaning. CLNet overcomes this problem by replacing the original $3 \times 3$ convolution operation by a $1 \times 1$ convolution operation, in which the real and imaginary parts can be embedded in physical meaning as

$$
i_1(1,1) = [a_1] \cdot [w_1] + [b_1] \cdot [w_1], \tag{24}
$$

where the ratio between the real and imaginary parts ($a_1$ and $b_1$) is preserved; thus, phase information is preserved. An ablation study shows that a 1 dB reconstruction gain can be achieved when CR is 1/16 for the indoor scenario. Inspired by CLNet, CVLNet proposed in [58] introduces a complex-valued NN to realize the CSI compression and reconstruction and all operations in CLNet follow the law of complex number computation.

The ENet proposed in [50] feeds back the real and imaginary parts separately with the same autoencoder. From [50], the real and imaginary parts of the complex-valued CSI share the same distribution. Based on this observation, the autoencoder trained with the real parts of the CSI samples can be used to compress and reconstruct the imaginary parts of the CSI samples, which can greatly reduce the NN complexity.

P-SRNet [68] introduces a principal component mark (PCM) module before the encoder. Only partial rows in the truanted CSI, such as the CSI "images" in Figs. 14 and 15, have non-zero elements. Therefore, the rows full of near-zero elements are omitted. The square of the Euclidean norm of each CSI row is first calculated. The rows with Euclidean norms below a threshold are selected and do not need to be fed back. As shown in (14-18), NN complexity, including the numbers of weights and FLOPs, increases with input

TABLE III
MULTI-DOMAIN CORRELATION UTILIZATION

| Correlation types | NN name | Main contributions in correlation utilization |
|---|---|---|
| **Time Correlation** | CsiNet-LSTM [85] | The first CSI "frame" is compressed by a high-CR encoder with the highest quality; The other $T-1$ CSI "frames" are compressed by low-CR encoders; LSTM refines the reconstructed CSI with the information extracted from the former CSI; |
| | RecCsiNet [86] | The LSTM at the encoder compresses the CSI based on the current and previous CSI matrices; |
| | [87] | Feedback overhead is reduced by dynamically adjusting feedback interval of time varying channel; Feedback is not needed if prediction errors are tolerable; |
| **Partial Bidirectional Channel Correlation** | DualNet-MAG [88] | The CSI magnitude and phase is fed back separately; Uplink CSI magnitude is introduced into the downlink CSI magnitude reconstruction at the decoder.; |
| | UA-CsiConvLSTM [89] HyperRNN [90] | The initial recovered downlink CSI is concatenated with the entire uplink CSI for further reconstruction; The partial bidirectional correlation is utilized by adopting hypernetworks; |
| **Frequency Correlation** | Attention-CsiNet [42] SampleDL [91] | Bi-LSTM module is adopted to extract the subcarrier correlation to compress CSI; The original channel is uniformly sampled in the frequency domain before feedback; |
| **Correlation Among Nearby Users' CSI** | CoCsiNet [92] | Two nearby users cooperatively feed their CSI magnitudes back to the BS; The information contained in CSI magnitude is divided into individual and shared information; The final CSI is reconstructed from the recovered individual and shared information; |
| | Distributed DeepCMC [93] | The CSI magnitude and phase are fed back together; A joint feature decoder reconstructs the CSI of two users.; |

dimension. The drop of partial rows reduces the dimension of the input to NNs. Hence, this strategy can greatly reduce NN complexity. Additionally, a binary indicating vector needs to be sent to the BS to guarantee the decoder to reconstruct a CSI with the same dimension as the original CSI, that is, to indicate which rows have been omitted.

*7) Others:*

*a) Deep unfolding architecture:* Most of the existing works are based on an autoencoder architecture, which is data-driven and lacks enough theory explanations [80], [81]. However, the data-driven method performs better than the model-driven method, which introduces expert/domain knowledge into the model constraints. Deep unfolding [82] unfolds the inference iterations as NN layers and unties the model parameters via end-to-end learning. It is a combination of data-driven and model-driven methods, and has been regarded as a potential direction in wireless communications.

Deep unfolding for CSI feedback is first introduced in [69]. Fig. 16 shows the TiLISTA-Joint architecture proposed in [69]. In the tied ISTA algorithm, the reconstruction problem can be solved by the following iterative formulas [25], [83]

$$\mathbf{r}^{(t)} = \hat{\mathbf{H}}^{(t-1)} + \rho^{(t)}\mathbf{W}\big(\mathbf{s} - g(\hat{\mathbf{H}}^{(t-1)})\big), \qquad (25)$$

$$\hat{\mathbf{H}}^{(t)} = f_d\Big(soft\big(f_s(\mathbf{r}^{(t)}); \theta^{(t)}\big)\Big), \qquad (26)$$

where $\mathbf{s} \in \mathbb{R}^{m \times 1}$ and $\hat{\mathbf{H}}^{(t)} \in \mathbb{R}^{2N_tN_c \times 1}$ represent the codeword and the output of the $t$-th iteration[4], respectively; $f_s(\cdot)$ and $f_d(\cdot)$ denote the sparse transformation and the inverse transformation, respectively; $soft(\cdot; \cdot)$[5] represents the soft-thresholding function; $g(\cdot)$ stands for the NN-based compression operation that reduces CSI dimension from $\mathbb{R}^{2N_tN_c \times 1}$ to $M \times 1$; $\rho^{(t)}$ and $\mathbf{W}$ are the step size of the $t$-th iteration and a linear operator, respectively. The original signal in a sparse domain is reconstructed by the soft-thresholding function [84]. In traditional iterative algorithms, the hyperparameters of tied ISTA, namely, $\rho$, $\mathbf{W}$, $g(\cdot)$, $\theta$, are selected by a time-consuming grid search. In TiLISTA-Joint architecture, these parameters are learned by an end-to-end approach, that is,

gradient descent algorithm. Sparse transformation operations, $f_s(\cdot)$ and $f_d^{(\cdot)}$, are realized by two FC layers without bias terms. The TiLISTA-Joint architecture outperforms the traditional iterative algorithms and CsiNet by a large margin. The fast ISTA algorithm is unfolded in [70] with a two-stage low-rank feedback scheme, in which the original CSI is represented by a basis and a residual part of the column space channel matrix.

*b) Fully FC architecture:* As mentioned in Section IV-A3, some works try to realize CSI feedback fully by convolutional layers. By contrast, a feedback NN architecture (called CF-FCFNN) in [48] only consists of FC layers. The CF-FCFNN architecture based on FC layers can extract spatial features more sufficiently compared with CsiNet based on convolutional layers and greatly outperforms CsiNet, especially when the feedback difficulty is high. For example, when CR is 1/64 for the outdoor scenario, the NMSEs of CsiNet and CF-FCFNN are −2.02 and −7.25 dB, respectively. However, the NN parameter number of CF-FCFNN is rather large.

*B. Multi-domain Correlation Utilization*

In Section IV-A, some methods to improve CSI feedback by introducing novel NN architectures are discussed. Different from the images in computer vision, the CSI "images" contain rich information about geometrical wireless propagation, which can be exploited to further improve CSI feedback. Multi-domain correlations have been adopted by the conventional CSI feedback methods. For example, a distributed CSI acquisition framework is developed in [19] to utilize the joint sparsity structure in downlink CSI matrices owing to the shared local scatters of the physical propagation environment. The partial channel reciprocity-based CSI feedback codebook in [94] exploits that bidirectional channels have a similar angular-delay distribution. Inspired by these works, multi-domain correlations, including time correlation, partial bidirectional channel correlation, and correlation among nearby users' CSI, have been also introduced for DL-based feedback, as shown in Table III. In this part, how to embed these correlations into DL-based CSI feedback is briefly introduced.
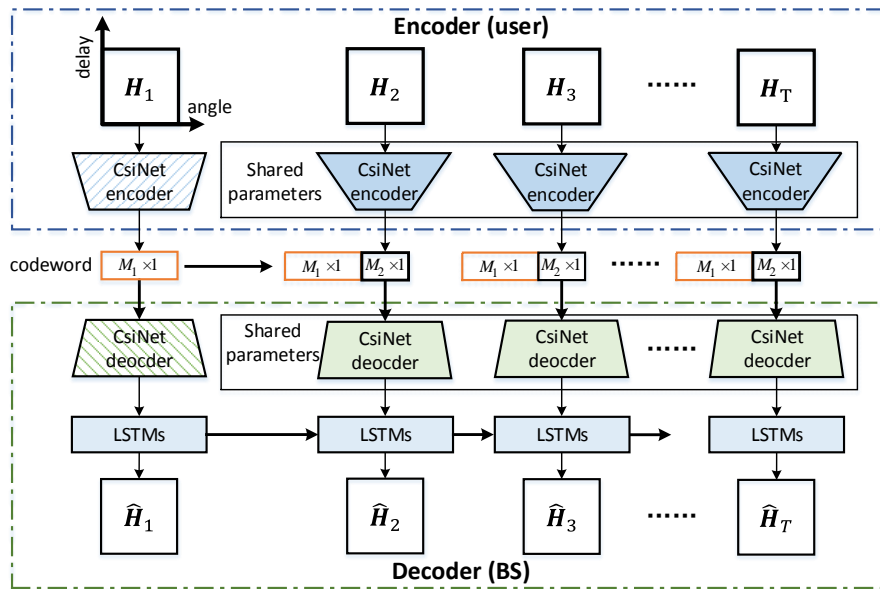
---

[4]The initial value $\hat{\mathbf{H}}^{(0)}$ is set as zero because of the CSI sparsity.

[5]The soft-thresholding function can be written as $soft(x; \theta) = \text{sign}(x)\max(0, |x| - \theta)$, where $\theta$ denotes the shrinkage threshold.

Fig. 17. Overall architecture of CsiNet-LSTM [85]. The first channel $\mathbf{H}_1$ and other $T - 1$ channels are compressed by high- and low-CR encoders, respectively. The output of the encoder, that is, codeword, is concatenated with the first one before being sent to the decoder. The initially reconstructed CSI is then refined by the LSTM modules.
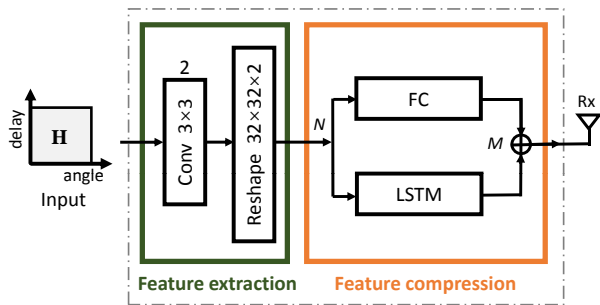


Fig. 18. Overall architecture of RecCsiNet encoder [86], which consists of feature extraction and compression modules

*1) Time Correlation:* In a time-varying scenario, the user location is not fixed. However, the user's moving distance within a short time (for example, feedback interval) is small. For a user with a moving speed of 360 km/h, the moving distance within 1 ms is only 0.1 m. Therefore, the environment around the user does not fully change. The channel is determined by the propagation environment; thus, the CSI at adjacent slots exhibits a high correlation. To model a CSI time evolution, the first-order Markov process can be adopted as [95], [96].

$$\mathbf{H}_t = \alpha \mathbf{H}_{t-1} + \sqrt{1 - \alpha^2} \mathbf{G}_t, \qquad (27)$$

where $\alpha \in [0, 1)$ represents the temporal correlation coefficient between the adjacent channels, and $\mathbf{G}_t$ denotes a zero-mean and unit-variance complex Gaussian matrix. $\alpha \to 1$ generates a time-invariant CSI matrix, and $\alpha = 0$ represents that the CSI has no time correlation. Considering the time correlation, the CSI of the time-varying scenario can be regarded as sequence data, such as video. However, time-varying CSI cannot be compressed fully the same as the video. The adjacent frames

of a video can be compressed together to save storage space. However, the user feeds back the estimated downlink CSI successively.

The novel NN architecture in [85], called CsiNet-LSTM, utilizes the time correlation to help CSI reconstruction at the decoder by LSTMs, as shown in Fig. 17. The basic encoders and decoders in CsiNet-LSTM are the same as those in CsiNet. Instead of feeding back a single CSI "image," the CsiNet-LSTM is designed for a sequence of CSI matrices. For a CSI sequence with length $T$[6], the first CSI "frame" is compressed by a high-CR encoder with the highest quality, and other $T - 1$ CSI "frames" are compressed by low-CR encoders. The low-CR encoders share the same NN parameters. The first CSI, $\mathbf{H}_1$, is used as a reference to help the reconstruction of the remaining CSI. Therefore, the last $T - 1$ codewords are concatenated with the first one. The CsiNet-based decoders recover the CSI from the codewords initially. Then, the initially reconstructed CSI is fed into a three-layer LSTM module. The LSTM module refines the reconstructed CSI with the information extracted from the former CSI. The simulation results show that CsiNet-LSTM has the least performance loss with the decrease of CRs compared with CsiNet.

CsiNet-LSTM only embeds LSTM into the decoder to exploit time correlation. However, no changes are introduced to the encoder of CsiNet-LSTM, that is, time correlation is ignored during compression. The RecCsiNet in [86] exploits LSTM to enhance the compression and reconstruction of a time-varying channel. Fig. 18 shows the encoder architecture of RecCsiNet. The encoder consists of two modules: feature extraction and compression. Similar to the CsiNet encoder, a convolutional layer with $3 \times 3$ kernels is first used to extract the feature of downlink CSI. The feature compression contains
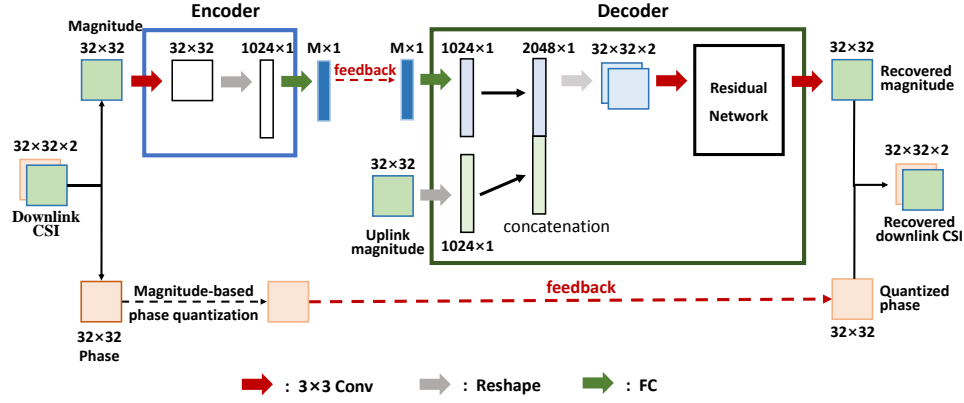
[6]The sequence length $T$ is set as ten in [85].

Fig. 19. Overall architecture of DualNet-MAG [88], which feeds back the CSI magnitude and phase. Uplink CSI magnitude is introduced into the reconstruction of the downlink CSI magnitude at the decoder.
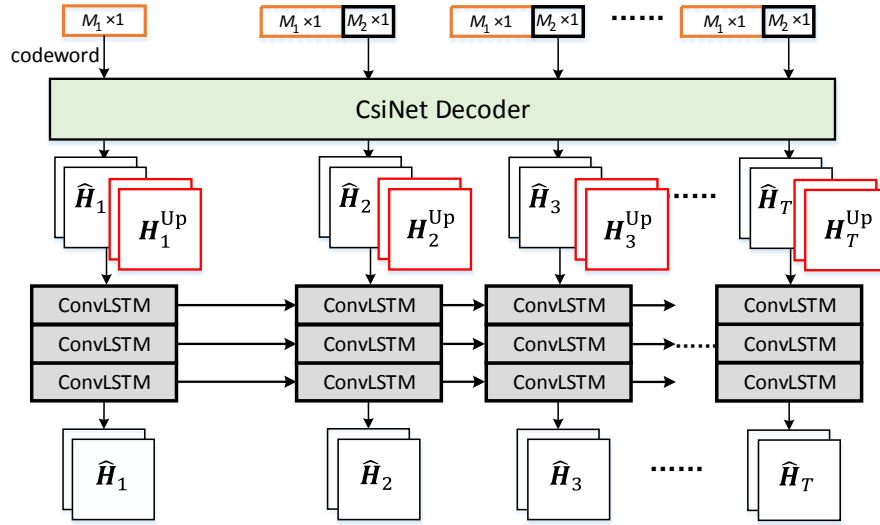


Fig. 20. Decoder architecture of UA-CsiConvLSTM [89], which introduces the time and partial bidirectional correlations to the CSI feedback

two parallel paths: a linear FC layer and an LSTM module. As described in (27), the current CSI contains the information of the previous one. The transmission bandwidth is wasted if the shared information is fed back repeatedly. Therefore, the LSTM module compresses the CSI based on the current and previous CSI matrices. The FC layer, which can be regarded as a jump connection, is used to accelerate the convergence of NN training. The decoder of RecCsiNet consists of two modules, feature uncompression and channel recovery. The feature uncompression module performs the inverse of the feature compression module. Then, the RefineNet proposed by [23] is used to improve reconstruction accuracy. Based on RecCsiNet architecture [86], a novel NN architecture, called ConvlstmCsiNet, in [97] replaces the feature extraction at the encoder and the RefineNet module at the decoder with more powerful NNs, thereby improving CSI feedback performance. An attention module is added after the LSTM of the feature compression/decompression block in RecCsiNet in [98].

Unlike [85], [86], [97], the feedback overhead is reduced in [87] by dynamically adjusting the feedback interval of the time-varying channel instead of reducing the overhead of each CSI feedback. A prediction NN in [87], which produces the current CSI based on the knowledge of the past CSI sequence, is shared by the user and the BS. If prediction errors are tolerable, the user does not need to feed back the current CSI, and the BS directly uses the CSI produced by the shared prediction NN. However, feedback is needed when the difference is over a predefined threshold. Numerical simulation shows that the MSE is reduced by 19.9% compared with the regular feedback strategy.

*2) Partial Bidirectional Channel Correlation:* Downlink CSI cannot be inferred from uplink CSI in frequency-division duplex (FDD) systems because the operating frequencies of the downlink and uplink are different. However, the signal propagation environment is the same for the downlink and uplink. Therefore, the bidirectional channels hold a partial correlation [94]. The accuracy of the reconstructed downlink CSI becomes better if uplink CSI is exploited.

The high correlation between the magnitudes of the bidirectional channels is exploited in [88]. Fig. 19 depicts the DualNet-MAG framework proposed by [88]. The quantized CSI phase is directly fed back via the uplink control channel. By contrast, the CSI magnitude is compressed by an NN-
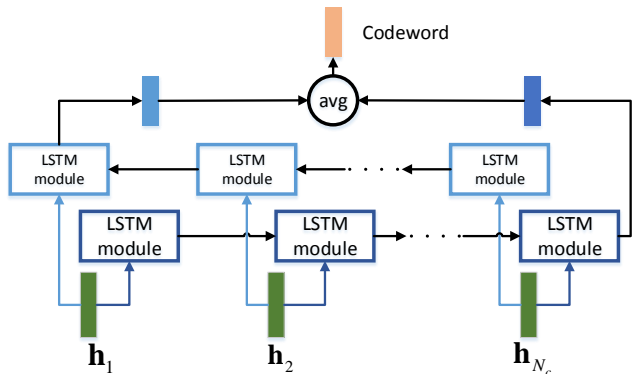
Fig. 21. Encoder architecture of Attention-CsiNet [42], which utilizes the correlation among adjacent subcarriers by Bi-LSTM modules
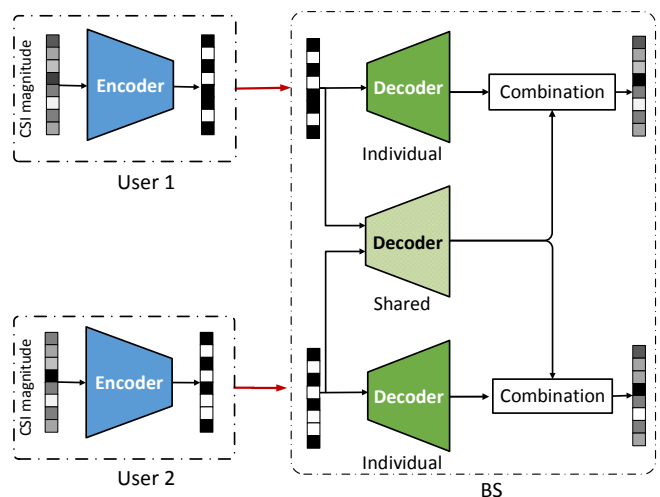


Fig. 22. Overall architecture of CoCsiNet [92], which consists of individual and shared decoders to recover the individual and shared information of the users' CSI magnitude

based encoder. Once receiving the feedback codeword, the BS concatenates the codeword with the corresponding uplink CSI magnitude and sends the concatenated vector to the decoder, which reconstructs the downlink CSI with the information not only from the feedback codeword but also from the uplink CSI magnitude. The simulation results show that the introduction of uplink CSI magnitude can greatly improve feedback accuracy. However, this feedback strategy [88] results in a bit-allocation between the magnitude and phase feedback. The original loss function in [88] is modified in [99] and [100] by directly introducing the phase and magnitude to the MSE function of CSI reconstruction to ensure an end-to-end training of CSI phase and magnitude feedback.

A novel feedback framework in [89], called UA-CsiConvLSTM, exploits time and partial bidirectional correlations, as shown in Fig. 20. The encoder part is the same as that of the CsiNet-LSTM [85]. Upon receiving the codeword, the BS initially reconstructs the downlink CSI with the decoder of the CsiNet [23]. Then, the recovered downlink CSI is concatenated with the uplink CSI. The concatenated vectors are sent to a three-layer ConvLSTM block, which utilizes time and partial bidirectional correlations to refine the initial downlink CSI. This strategy forces the NNs to learn and exploit the correlation automatically, thereby preventing the bit allocation between the magnitude and phase feedback in [88].

The HyperRNN in [90] utilizes the partial bidirectional correlation by adopting hypernetworks [101] instead of directly sending the uplink CSI to the decoder as [88], [89]. The key idea of hypernetworks is to use a single network (called as hypernetwork) to generate the weights for another NN. The estimated uplink CSI is sent to an FC layer to generate the weights of the NNs used to reconstruct the downlink CSI. The hypernetwork introduces uplink channel information into downlink channel reconstruction through these generated NN weights.

*3) Frequency Correlation:* The channels over adjacent subcarriers are highly correlated. Therefore, some works extract and utilize this correlation to further reduce feedback overhead.

The Attention-CsiNet proposed in [42] adds the LSTM modules to the encoder and decoder. As pointed out in [42],

CsiNet in [23] ignores the correlation between subcarriers. Therefore, the Bi-LSTM module is adopted to extract the subcarrier correlation to compress CSI as shown in Fig. 21. The final codeword is the average of the output of two LSTM modules. In CsiNet-LSTM [85], only a unidirectional LSTM is adopted because the current CSI is reconstructed with the help of the previous CSI but without the help of the next moment CSI. However, the channel feedback over a certain subcarrier can be enhanced by channels over all other subcarriers. Moreover, two LSTM modules in Bi-LSTM share the same NN weights, thereby dramatically reducing NN weight numbers.

For the novel compressive samples CSI feedback framework in [91], called SampleDL, the original channel is uniformly sampled in the frequency domain before feedback. Only channels over selected subcarriers are fed back to the BS. The autoencoder compresses and reconstructs the sampled CSI. Then, the reconstructed CSI is interpolated with 0 to recover its original dimension, and an extra NN is adopted to refine the interpolated CSI. This method can increase the feedback accuracy and reduce the NN complexity due to the reduction of NN input.

*4) Correlation Among Nearby Users' CSI:* The user number increases substantially in 6G. As pointed out in [102], user density may grow to hundreds per cubic meter, which poses a high requirement of spatial spectral efficiency. Based on practical measurements in [103], channel correlation is higher than 0.48 for all close-by users. The far-away users have an inter-user CSI correlation that is more than twice higher than that of the i.i.d. CSI, even when the distance of the users is over tens of wavelengths [7]. Based on this observation, the CSI correlation among nearby users can be utilized to improve feedback accuracy and reduce feedback overhead. Inspired by this, CoCsiNet [92] and distributed DeepCMC [93] introduce this correlation to DL-based CSI feedback.

---

[7] The frequency in [103] is 2.4 GHz, and the wavelength is 12.5 cm.

Fig. 22 shows the framework of the CoCsiNet proposed by [92]. Inspired by [88], CoCsiNet feeds back the CSI phase and magnitude, respectively. Two nearby users cooperatively feed their CSI magnitudes back to the BS due to the observation of correlation in the CSI magnitude domain. The information contained in CSI magnitude is divided into two kinds in [92]: individual and shared information. The encoders of the nearby users compress and quantize the CSI to generate a bitstream. Then, two different decoders are adopted to reconstruct the CSI from the feedback bitstream. The first decoder focuses on recovering the individual information in the user CSI. The shared decoder can recover the information shared by two users. The final CSI is reconstructed from the recovered individual and shared information. The information shared by nearby users does not need to be fed back any more, thereby reducing feedback overhead. Two magnitude-dependent methods introduce instant and statistical information of the CSI magnitude into the feedback of the CSI phase to feedback the CSI phase efficiently. Visualization of the encoder parameters in [92] shows that the nearby users can cooperatively feedback the shared path information after an end-to-end NN training.

Similar to [89], the distributed DeepCMC in [93] does not separately feedback the CSI magnitude and phase when exploiting the correlation among the nearby users. The encoder part of the distributed DeepCMC is the same as that of CoCsiNet. However, the input of the encoder is the complex CSI instead of the CSI magnitude. The BS concatenates the feedback codewords of two users and sends them into a joint feature decoder to reconstruct their CSI. The summation-based fusion branches in the distributed DeepCMC exploit the property of channel gains, which consist of the summation of multipath signal components.

### C. Bitstream Generation

In practical systems, the CSI is fed back in the form of bitstreams. If a 32-bit floating point codeword, that is, the encoder's output, is directly fed back, the overhead is very large. Therefore, the codeword needs to be discretized before feedback. Quantization error is considered a part of the errors introduced in the feedback process in [104]. The simulation results show that the feedback errors heavily affect the feedback accuracy.

Reference [105] introduces that a uniform module is added after the encoder in [105] and the number of quantization bits, $B$, is set as 4. Quantization can be written as[8]

$$\mathbf{s}_q = \frac{\text{round}(2^{B-1} \times \mathbf{s})}{2^{B-1}}. \tag{28}$$

This quantization method is easy to implement. However, the rounding operation is non-differentiable, making the quantization operation unable to be directly embedded into the end-to-end training based on gradient descent. Therefore, the gradient of the rounding operation is set as 1 during training in [105]; thus, the autoencoder can be end-to-end trained with the uniform quantization module. The rounding operation

[8]We assume that the activation function of the last layer at the encoder is Tanh, i.e., $\mathbf{s} \in (-1, 1)$
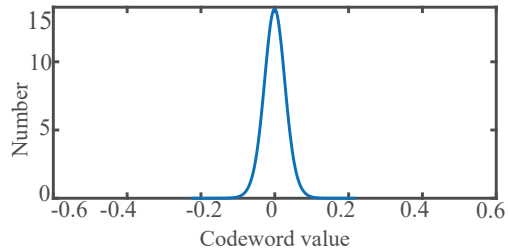


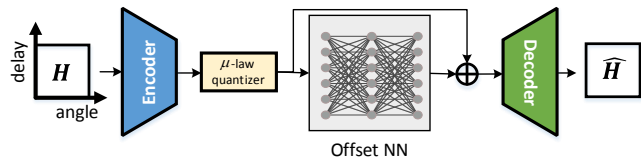Fig. 23. Distribution map of the compressed CSI values when CR is set as 1/4 for the indoor scenario [106].



Fig. 24. Bit-level autoencoder-based CSI feedback framework [41], [106], which adopts a $\mu$-law non-uniform quantizer. Upon receiving the quantized codeword, an offset NN at the BS refines the quantized codeword and then sends the refined codeword to the decoder.

is replaced by a differentiable Sigmoid-based approximate rounding function in [100].

Based on [41], [106], the values of most codeword elements are almost near-zero, as shown in Fig. 23. The uniform quantization provides unnecessary quantization performance for the high values that seldom appear in practical signals; thus, it is unsuitable for the quantization of CSI codeword. A quantizer, with smaller step sizes at lower amplitudes and larger step sizes at high amplitudes, is needed. A non-uniform quantization, that is, $\mu$-law quantizer [107], is introduced in [41], [106] to the CSI codeword quantization to meet the above requirement. Fig. 24 shows the bit-level CSI feedback framework proposed by [41], [106]. The downlink CSI is first compressed by an encoder and then quantized by a $\mu$-law non-uniform quantizer to generate a bitstream. Upon receiving the quantized codeword, the BS first refines the codeword by an offset NN with residual learning to reduce the effect of the quantization errors. Then, the refined codeword is sent to the decoder for CSI reconstruction. A two-stage training stage is used because the gradient cannot be passed during backpropagation. In the first stage, quantization is not considered, and the encoder and decoder are jointly trained with collected CSI samples. In the second stage, the codeword is discretized by the $\mu$-law quantizer. The offset NN is trained to minimize the errors introduced by quantization. Then, the decoder is finetuned with the refined codeword and the original high-quality CSI matrices. Simulation shows that the non-uniform quantization with an offset NN outperforms the uniform quantization by a large margin.

Entropy coding is introduced in [63] to CSI feedback to reduce feedback overhead further. In [63], the codeword is first quantized by a uniform quantizer. Then, context-based adaptive binary arithmetic coding [108] converts the quantized values into a bitstream. This entropy coding is dependent on the input probability model learned from the codeword. As mentioned before, the quantization is non-differentiable
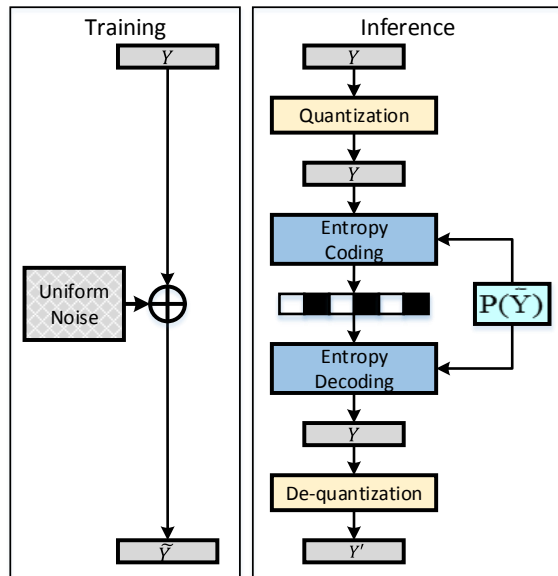
Fig. 25. Illustration of entropy bottlenecks during NN training and inference [109]. The uniform noise is added to the codeword during the NN training. During NN inference, the codeword is uniformly quantized and then entropy coded by the stored distribution $P(\tilde{Y})$. Finally, the generated bitstream is fed back. The codeword is reconstructed by entropy decoding and dequantization operations.
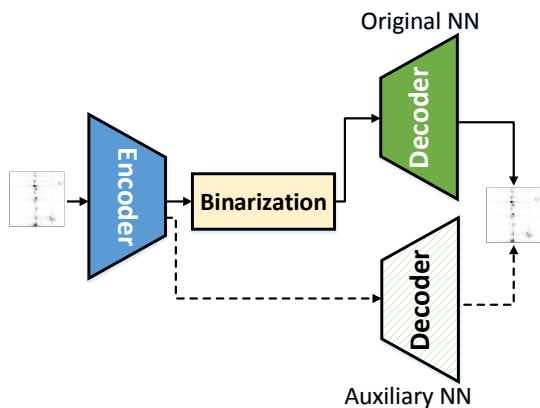


Fig. 26. Illustration of knowledge distillation-aided bit-level CSI feedback training framework, which consists of the original and highly complex NNs [111]

and cannot be directly embedded into end-to-end learning. Therefore, the quantization in this work is approximated as a random noise during the training instead of setting the gradient as one. The loss function in [63] consists of two parts: reconstruction error and entropy of feedback codeword.

Quantization and entropy coding are also used in the DL-based CSI feedback in [109] by an entropy bottleneck layer [110], which is composed of the quantizer, entropy model and coder, and dequantizer. The encoder and decoder in [109] are based on the CRNet [43]. During the NN training, a random uniform noise is added to the output of the encoder, which is similar to the operation in [63]. Fig. 25 depicts the entropy bottlenecks during NN training and inference. The numerical results show that the method in [109] performs better than those in [41], [106] for a wide range of bit rates.

The effect of binarization (1-bit quantization) on feedback performance is also considered in [111]. Knowledge distillation enables a teacher NN to distill and transfer dark knowledge to a simple student NN [112] for bit-level CSI feedback. In Fig. 26, an extra highly complex NN is introduced during the NN training. The highly complex NN is the same as the original NN except that it considers no binarization. Therefore, the gradient can be passed during the training of the highly complex NN. To utilize the highly complex NN during the training, the joint training method in [111] alternatively trains the original and the highly complex NNs. The highly complex NN can guide the original NN training and prevent the original NN to fall into a poor local minimum because the gradient of the highly complex NN is lossless. This process can be regarded as that the highly complex NN transfers its knowledge to the original NN.

### D. Joint Design with Other Modules

The above works assume that the user directly feeds back perfect downlink CSI to the BS. However, the CSI is estimated from the downlink pilot signals, and the channel estimation inevitably introduces errors to the downlink CSI. The quality of the recovered CSI at the BS becomes poor if the estimated CSI is fed back by the NNs trained with perfect CSI samples [121]. Therefore, the entire CSI acquisition needs to be considered during the design of CSI feedback. Fig. 27 shows the workflow of the downlink CSI acquisition and utilization in FDD massive MIMO systems. First, the BS designs the downlink pilots and transmits pilot signals to the users. Second, the user estimates the CSI from the received pilot signals. Third, the user compresses and quantizes the estimated CSI and transmits the bitstream to the BS. The BS reconstructs the downlink CSI from the feedback information. Finally, the BS designs the BF/precoding matrix based on the downlink CSI. In this part, some existing works on jointly designing the CSI feedback with other modules to maximize the performance gains, as shown in Table IV, are introduced.

#### 1) Joint Channel Acquisition:

*a) Joint channel estimation and feedback:* Two joint channel estimation and feedback frameworks are introduced in [113]. The first framework, namely, PFnet, regards the channel estimation and feedback as one module and directly compresses the received pilot signals with an NN-based encoder. The decoder reconstructs the downlink CSI from the feedback bitstream. This framework regards the CSI acquisition problem as an end-to-end black box. The second framework, namely CEFnet, combines the communication knowledge with NNs. The coarse downlink CSI is estimated by some simple
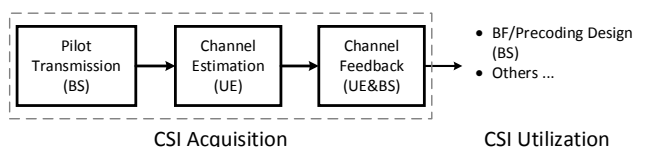


Fig. 27. Workflow of CSI acquisition and utilization in FDD massive MIMO systems

TABLE IV
JOINT DESIGN WITH OTHER MODULES

| Joint design types | Functions | Main contributions in joint design |
|---|---|---|
| **Joint Channel Acquisition** | Joint channel estimation and feedback | PFnet [113]: directly compressing and feedbacking the received pilot signals with an encoder; CEFnet [113]: refining the estimated coarse CSI and then feedbacking it via autoencoder; AnciNet [114]: introducing an two-stage training strategy when estimation and feedback are considered; |
| | Joint pilot design, channel estimation, and feedback | CAnet-J [115]: the BS first designs pilots based on uplink CSI magnitude. Upon receiving pilot signals, the user compresses the received signals using an encoder. Then, the decoder reconstructs downlink CSI by utilizing the information from feedback bitstreams and the uplink CSI magnitude; HyperRNN [90]: different from CAnet-J, the pilot is denoted by the weights of an FC layer; |
| **Joint Channel Feedback and Utilization** | Joint channel feedback and BF design | CsiFBnet, [116]: the decoder NNs directly produce a BF vector that maximizes the BF gain; [117]:jointly designing the feedback and hybrid precoding and pointing out that the gain achieved by joint design is large, especially when the feedback is extremely limited; CsiCPreNet [118]: the BSs exchange the feedback codewords with one another when multiple-cell is considered, and the precoding matrix is generated by a coordinated precoder design NN; |
| **Joint Channel Acquisition and Utilization** | Joint pilot design, channel estimation, feedback, and precoding design | [119]: proposing an end-to-end limited feedback framework, including channel estimation, feedback codebook design, and BF vector selection, for a single-user scenario; [111]: extending the work in [119] to a multiuser scenario; [120]: also including the pilot design module compared with [119]; |

algorithms, such as least-square estimation. Then, an extra estimation subnet is employed to refine the coarse CSI. Finally, the refined CSI is fed back to the BS by an autoencoder. A two-stage training strategy is used to train the CEFnet. During the first training stage, the estimation subnet is trained with coarse CSI as input data and perfect CS as target output. During the second stage, the feedback subnet is trained with the output of the estimation subnet as input data and the ideal CSI as ground truth. The CEFnet remarkably outperforms the black-box PFnet. The CEFnet framework can also be trained in an end-to-end manner. Based on the comparison of the one-stage end-to-end training with the two-stage one in [114], the two-stage training can bring considerable performance gains.

*b) Joint pilot design, channel estimation, and feedback:* Pilot length is limited due to the limited downlink training resource. A joint pilot design and channel estimation strategy is developed in [115] to reduce pilot overhead and channel estimation errors. Unlike the methods in [122], [123] that denote the downlink pilot by the weights of an FC layer, the pilot matrix in [115] is produced based on the uplink CSI magnitude because of the correlation between bidirectional channel magnitudes. Then, an uplink-aided entire CSI acquisition framework, CAnet-J, is shown in Fig. 28. The BS first designs pilots based on the uplink CSI magnitude in the angular domain. Upon receiving pilot signals, the user compresses and quantizes the received signals using an NN-based encoder without channel estimation. Then, the decoder reconstructs downlink CSI by utilizing the information from feedback bitstreams and the uplink CSI magnitude. The numerical results show that the joint design outperforms the method that separately estimates and feeds back CSI [113], [114]. If channel estimation and feedback are separately implemented, the pilot signals received by the user need to contain all information of the downlink CSI. By contrast, the signals do not need to contain the information shared with the uplink CSI if pilot signals are directly fed back and CSI is reconstructed with the uplink CSI at the BS. Therefore, CAnet-J can perform better specifically when the pilot length is limited.

HyperRNN in [90] also considers the entire CSI acquisition. The main difference between HyperRNN and CAnet-J is the method of producing pilots and introducing uplink CSI
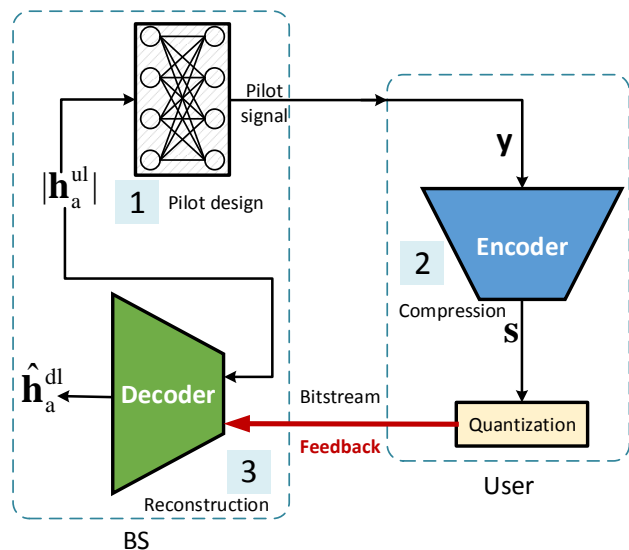


Fig. 28. Flowchart of DL-based uplink-aided joint CSI acquisition [115]

knowledge to the decoder at the BS. In HyperRNN, the pilot is denoted by the weights of an FC layer, which is similar to that in [122], [123]. As mentioned in Section IV-B2, HyperRNN introduces the uplink CSI to the downlink CSI reconstruction by a hypernetwork. The hypernetwork generates the weights of the NN used in downlink CSI reconstruction based on the input uplink CSI. Moreover, reference [90] also considers the effect of the imperfect uplink CSI on the downlink CSI acquisition and jointly designs the acquisition of uplink CSI.

*2) Joint Channel Feedback and Utilization:* As indicated in [116], the existing feedback methods only focus on obtaining as accurate downlink CSI as possible and ignore the physical meaning of downlink CSI. Therefore, the existing works reduce the CSI dimension by dropping the redundant information with a small effect on the reconstruction accuracy, that is, MSE/NMSE. However, signal fidelity cannot be exactly measured by MSE/NMSE [124]. Sometimes, the performance of communication systems may be poor with a good MSE/NMSE. Therefore, the effect of the CSI feedback on the next module, that is, the BF design, should be considered jointly. The feedback framework for BF design in [116], called

CsiFBnet, maximizes the BF gain instead of the feedback performance. In single-cell systems, the user compresses the CSI with an encoder and quantizes the compressed codeword. Upon receiving the codeword, the decoder NNs produce a BF vector. Considering the constant modulus constraint on the analog BF vector, the output of the NNs at the BS is the phase of the BF vector. The NNs are end-to-end trained by an unsupervised approach. In multicell systems, the soft hand-off model [125] is considered, and the user needs to feedback the desired and the interfering CSI matrices to maximize the sum rate, which is a complicated joint optimization problem and turns into a local one by the approximation proposed in [126]. Simulation shows that joint design can greatly increase the mean rate of massive MIMO systems, and the NNs show high generalization to signal-to-noise ratio (SNR) and path number of CSI. As mentioned in [117], the gain achieved by joint design is large, especially when the feedback is extremely limited. A more complicated scenario is considered in [118], where the user receives the signals from all cells instead of only the nearby cells in [116]. The framework in [118] consists of two modules: a CSI compression NN and a coordinated precoder NN. Upon receiving the feedback bitstreams, the BSs exchange the feedback codewords with one another. Then, the codewords are concatenated and sent to the coordinated precoder design NN to generate the precoding matrix.

*3) Joint Channel Acquisition and Utilization:* Reference [119] proposes an end-to-end limited feedback framework, including channel estimation, feedback codebook design, and BF vector selection. Specifically, the NNs at the user compress the pilot signals without channel estimation and discretize the compressed vector by a binarization operation. Upon yielding the binary feedback information, the NNs at the BS generate a BF vector, which can maximize the channel gain. The NNs at the user and the BS are jointly trained in an unsupervised manner. The work in [119] is extended to a multiuser scenario in [111]. The encoders at different users are the same. The feedback bitstreams of multiple users are concatenated at the BS and sent to the NNs to generate the precoding matrix. The optimization goal is to maximize the sum-rate instead of channel gain [119].

Parallel to [111], a joint CSI acquisition and precoding design framework is also proposed in [120]. This work includes the pilot design module. For the data-driven pilot design, the pilot is denoted by the weights of an FC layer as [90], [122], [123]. Moreover, a two-step training strategy is proposed to make the trained NNs generalizable to different feedback overheads. In the first step, the NNs at the UE and the BS are jointly trained by an end-to-end approach. Quantization is neglected in this step. In the second step, the weights of the user-side NNs are fixed. Different quantization steps are applied to the NN output at the user, that is, feedback vector. For each quantization step, a specific NN is trained at the BS.

### E. Practical Consideration

*1) Multirate Feedback:* Some practical communication systems need to adjust the number of feedback bits in accordance with the scenarios. For example, for the single-user scenario,
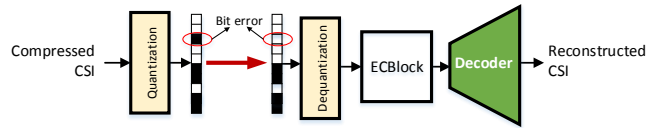


Fig. 29. Illustration of ECBlock-aided digital CSI feedback framework [131], in which quantization errors and feedback bit errors are considered. At the BS, the codeword is refined by the ECBlock before being sent to the decoder.

TYPE I feedback, with very low feedback overhead, is adopted in 5G NR systems. For a more complicated multiuser case, TYPE II feedback, with a much larger overhead, is preferred. Therefore, DL-based CSI feedback needs to generate codewords under different lengths/CRs and accuracy. A straightforward way is to train several an NN model for each CR, which occupies much space to store the NN parameters.

The authors of [41] focus on reducing the NN parameter number of the encoder and neglect that at the BS because of enough storage space at the BS. The FC layer in CsiNet+ contains nearly all NN parameters. For example, when CR is 1/4, the NN parameters of the FC layer occupy 99.91% of the entire encoder. Therefore, the FC layers are reused by all CRs. Two multirate frameworks, namely, serial multirate (SM-CsiNet+) and parallel multirate (PM-CsiNet+) frameworks, are proposed. The codeword under a low CR (such as 1/64) in SM-CsiNet+ is generated from that under a high CR (such as 1/32). By contrast, PM-CsiNet+ generates the codeword under a high CR from that under a low CR. The modular adaptive multirate (MAMR) framework in [127] also considers the complexity of the decoder at the BS. The input size of the decoder is fixed once trained. However, the codewords under different CRs have different sizes, which is solved in [127] by zero-padding. The NN parameter at the decoder can be reduced by approximately 42.5%. The framework in [128], called FOCU, combines the PM-CsiNet+ architecture in [41] with the padding operation in [127] to realize multirate reconstruction at the BS.

In practical systems, CR needs to be determined automatically. In [129], a CNN-based classification module is added before feedback. The classification model selects the suitable CR in accordance with the CSI. The key problem is how to generate the labels for NN training. IPredefining an accuracy threshold, such as $-10$ dB in [144], is suggested, and the lowest CR that meets the accuracy requirement is marked as the desired CR, that is, the label of the corresponding CSI. Then, a supervised end-to-end learning is employed for the classification NNs, in which the CSI is the input and the desired CR is the output.

*2) Imperfect Feedback Link:* In some practical environments, the feedback link suffers from various interference and non-linear effect, which disturb the feedback codeword. A plug-and-play denoise NN is added in [130] before the decoder, which is based on residual learning and similar to the offset network in [106]. The denoise NN consists of several FC layers and is trained to reduce the codeword noises introduced by imperfect uplink transmission. Compared with the NNs without consideration of imperfect feedback, the NNs with an extra denoise network show high robustness to the uplink

TABLE V
PRACTICAL CONSIDERATION

| Practical Consideration | NN Name or Method | Main contributions in practical consideration |
|---|---|---|
| **Multirate Feedback** | SM-CsiNet+ [41] | The codeword under a low CR (such as 1/64) is generated from that under a high CR (such as 1/32); |
| | PM-CsiNet+ [41] | The codeword under a high CR is generated from that under a low CR; |
| | MAMR [127] | The codewords under different CRs have different sizes and are padded with zeros at the BS; |
| | FOCU [128] | The PM-CsiNet+ architecture is combained with the padding operation; |
| | [129] | A classification module, which selects the suitable CR, is added before feedback; |
| **Imperfect Feedback Link** | DNNet [130] | A plug-and-play denoise NN is added before the decoder to reduce transmission errors; |
| | ECBlock [131] | The error correcting NN embedded before the decoder is trained with the autoencoder; |
| | AnalogDeepCMC [132] | The downlink CSI is directly mapped to the input of the uplink channel; |
| **NN Complexity** | NN weight pruning | [133]: the FC layer of the CsiNet+ encoder is pruned; |
| | NN weight quantization/binarization | [133]: the NN weights are quantized with 3-7 bits after pre-training; [61], [134]: the NN weights are quantized with 1 bit; |
| | Knowledge distillation | [135]: the knowledge of the complex CsiNet+ is transferred to the simple CsiNet; |
| **Data Collection and Online Training** | Data collection | [136]: the NMSE gap between the NNs trained with 3,200 and 800 CSI samples is 3.1 dB; [137]: the feedback NNs are trained using the uplink CSI samples due to the same characteristics; |
| | Online training | [59], [138]: transfer learning and meta learning is introduced to accelerate online training at the BS; [139]: the feedback NN is trained at the user side, and FL is adopted; [140]: a new encoder is trained at the user side for a specific area without changing the decoder, and gossip learning applied to multiuser scenario; |
| **Standardization** | ImCsiNet [141] and EVCsiNet [142] | The precoding matrix is fed back by an autoencoder instead of the whole CSI; |
| | AI4C$^2$F [143] | An NN module is added at the BS to refine the channel codeword obtained by codebook-based feedback; |

SNR. For example, the NMSE gap is up to 10 dB when uplink SNR is 5 dB.

Unlike [130], for the digital CSI feedback in [131], the codeword is first quantized by the method proposed by [145] and fed back in the form of the bitstream. Bit errors are inevitable due to imperfect transmission. Inspired by [130], an error correction block (ECBlock) is deployed before the decoder at the BS, as shown in Fig. 29. ECBlock consists of several FC layers, where residual learning is adopted. The NN training is divided into two stages, including pretraining and alternate training. In the first training stage, the encoder and decoder are first trained without feedback errors. Then, ECBlock is trained using the codewords generated by the encoder, which has been well trained. In the second stage, the entire model, including the encoder, quantization, adding bit errors, dequantization, ECBlock, and decoder, are connected together and trained by an end-to-end approach. Considering the operations of quantization and adding bit errors are non-differentiable, their gradient is set as 1 [145].

A CNN-based analog CSI feedback is adopted in [132]. It directly maps the downlink CSI to the input of the uplink channel and can be regarded as a joint source channel coding framework. This framework improves the robustness to the imperfect uplink transmission and simplifies the feedback process because of the joint source channel coding. Moreover, the end-to-end joint source channel coding framework for CSI feedback in [146] enhances NN robustness to the imperfect uplink transmission. Concretely, the uplink transmission SNR is input to the NNs, thereby making the trained NNs adaptive to the uplink channel condition.

*3) NN Complexity:* DL-based CSI feedback can improve feedback accuracy and reduce feedback overheads. Fig. 30 shows NMSE performance versus FLOP number of the entire NNs when CR is 1/16 for the indoor scenario. The FLOP number of TransNet [56] is approximately nine times that of CsiNet, and NMSE is reduced by 6.35 dB. Therefore, performance improvement is at the expense of NN complexity. However, the high requirement of DL-based algorithms in memory and computational complexity poses a major challenge to the deployment of DL-based feedback
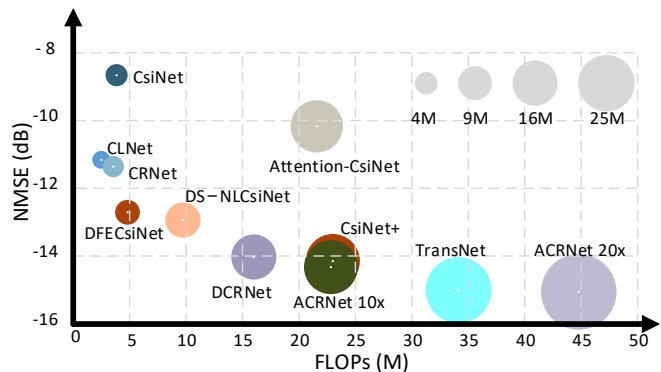


Fig. 30. NMSE (dB) versus FLOP number of entire NNs when CR is 1/16 for the indoor scenarios.
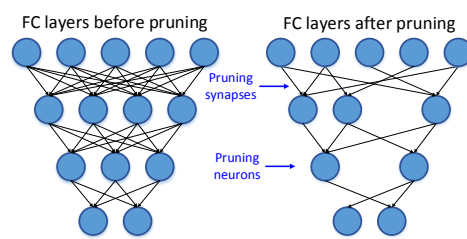


Fig. 31. Illustration of NN weight pruning in FC layers. NN complexity can be greatly reduced if redundant connections (synapses) and neurons are dropped.

to practical systems. Therefore, the NN weight pruning and quantization/binarization have been introduced to reduce the complexity of DL-based feedback.

*a) NN weight pruning:* As mentioned in Section IV-E1, the FC layer at the encoder occupies almost the entire weights of the encoder. In FC layers, most connections (synapses) and neurons are redundant. The weight number can be greatly decreased if redundant connections and neurons are dropped, as shown in Fig. 31. The basic idea of NN pruning is to remove the NN weights with small absolute values.

In [133], CsiNet+ in [41] is used as an example to prune the FC layer at the encoder. NN weight pruning has two kinds:

pruning during training and pruning after pre-training. The second pruning method is adopted by [133]. A binary mask, with the same shape as the FC weights, is added to the FC layer. The mask elements, with corresponding absolute values below a predefined threshold, are set as zero. Then, the entire NNs with the mask are finetuned with a small learning rate. The gradient flows through the fixed mask and the NN weights, with zero mask values, are not updated during the finetuning. Simulation shows that NN pruning can greatly reduce the weight number of FC layers with small effect on CSI feedback accuracy. When CR is 1/16 for the indoor scenario, accuracy drop is only 0.24 dB if 97.21% NN weights are pruned. This NN compression can be easily extended to other works. For example, the NNs for uplink-aided joint CSI acquisition in [115] are pruned similarly.

*b) NN weight quantization/binarization:* In most DL libraries, such as TensorFlow and PyTorch, the NN weights are set as 32-bit floating point, resulting in a waste in memory space and an increase in computational complexity. The computational power at the user is limited and cannot support high-precision computation. NN weight quantization is introduced in [133] to DL-based CSI feedback, where high-precision NN weights (such as 32-bit float point) are replaced with low-precision ones (such as 1-bit float point). In [133], the NN weights are quantized with 3-7 bits after pre-training. When the quantization bit of NN weights is set as 6 or 7 for the outdoor scenario, the performance gap between the original and the quantized NNs is small. Furthermore, BCsiNet in [61] and ACRNet in [134] quantize the NN weights with 1 bit, thereby offering over 30 times memory saving and approximate two times acceleration in inference speed with small effect on feedback performance.

*c) Efficient NN architecture design:* Early works, such as ConvCsiNet [40], improve NN performance by stacking the vanilla convolutional layers. A minute improvement sometimes is at the expense of a substantial increase in NN complexity. Therefore, the NN architecture should be carefully designed and the efficient NN architecture should be adopted instead of the redundant one. In [133] and [147], the vanilla convolutional layers in ConvCsiNet are replaced with more efficient convolution blocks, namely, the squeeze layer [148] and the shuffle layer [149], which can achieve a comparable feedback performance when FLOP numbers are 1/3 and 1/4 of ConvCsiNet.

As indicated in [51], the $3 \times 3$ convolution operation cannot offer enough receptive field, making the neurons in the deep layer unable to represent enough regions of the input CSI "images." Thus, dilated convolutions [150] are introduced in [51] to enhance the receptive field without a large increase in NN complexity. In Fig. 32, dilated convolutions inject holes into the vanilla convolution kernel. The dilated rate represents the interval number in the convolution kernel. If the dilated rate is set as 1, the dilated convolution is the same as the standard convolution. The dilated rate in Fig. 32 is set as 2, and the receptive field is $5 \times 5$. However, the NN complexity of this operation is the same as that of a standard convolution operation with $3 \times 3$ filters.
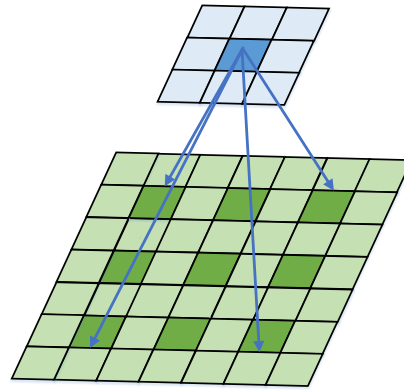


Fig. 32. Demonstration of dilated convolution in [51] when dilated rate is set as two. The receptive field is $5 \times 5$. However, the NN complexity of this operation is the same as that of a standard convolution operation with $3 \times 3$ filters.
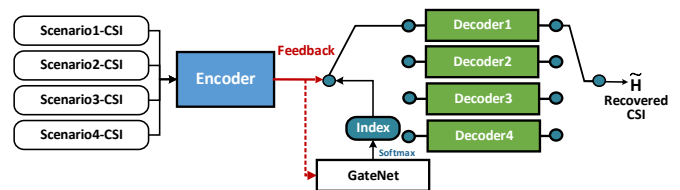


Fig. 33. Multitask learning-based CSI feedback framework for multiple scenarios in [160]

*d) Knowledge distillation:* In knowledge distillation, the knowledge learned by the complex teacher NN can be transferred to the simple student NN to improve the performance of the student NN. In [135], the knowledge of the complex CsiNet+ is transferred to the simple CsiNet. The performance of the CsiNet trained with knowledge distillation is greatly improved. For example, performance improvement is up to 7.5 dB when CR is 1/4 for the indoor scenario. Knowledge distillation can be regarded as an NN training trick, which is plug-and-play and can be easily extended to the existing works.

*4) Data Collection and Online Training:* Most existing works are conducted on simulation, in which the CSI samples can be easily obtained by channel generation software, such as COST 2100 [9] and QuaDRiGa[10]. When deploying DL-based feedback to practical systems, data collection and online training should be considered.

*a) Data collection:* NN performance depends on the number of CSI samples used during NN training. For example, the NMSE gap between the NNs trained with 3,200 and 800 CSI samples is 3.1 dB in [136]. A straightforward way for data collection is that the user sends the stored high-quality CSI samples to the BS as many as possible during idle time. However, this method contains two major problems. First, the user needs to store many CSI samples, occupying the storage space at the user. Second, the uplink transmission of CSI samples occupies precious uplink transmission resources. Therefore, it is difficult to deploy in practical systems.

[9]https://github.com/cost2100/cost2100
[10]https://quadriga-channel-model.de/

The NNs of downlink CSI feedback in [137] are trained using the uplink CSI samples with the same characteristics/statistics. Although bidirectional transmissions are operated over different frequency bands, they share the same propagation environment, which determines the CSI distributions. The numerical results show that the feedback accuracy of the NNs trained by uplink CSI samples is close to that trained by downlink CSI when the duplex distance is 200 MHz.

*b) Online training:* The propagation environment is usually stable for a long time. However, once the environment greatly changes, the CSI distribution also changes. NNs trained with the previous distribution cannot work well in a new environment. Therefore, online training is essential. Inspired by [151] that applies transfer learning to DL-based CSI prediction, several novel online training strategies are introduced to DL-based CSI feedback in [59] and [138] to accelerate the training convergence. Once the environment changes, pretrained NNs are fine-tuned using new CSI samples, which is regarded as transfer learning. Then, the model-agnostic meta learning algorithm, which trains NNs by alternating inner-task and across task updates and then adjusts the original NNs for a new environment with few CSI samples, is adopted to accelerate the NN training further. The training in [59], [138] is implemented at the BS. In [139], the feedback NN is trained at the user side. However, each user only stores some local CSI samples, which are not enough to train an NN for a whole cell. Thus, a distributed learning framework, that is, federated learning (FL) [152], is introduced to the NN training. In FL-aided online training, the user sends the NN gradient to the BS instead of CSI samples, thereby reducing the communication overhead. The BS, which can be regarded as an aggregation server, aggregates the received NN and then transmits a global NN model to each user.

As mentioned before, the user occasionally stays in an area (e.g., an office) for a long time. In this scenario, the propagation environment is relatively stable. An online training framework is proposed in [140] to utilize the above observation, where a new encoder is trained at the user side for a specific area without changing the decoder at the BS side. The NN training is employed at the user, and the CSI datasets do not need to be sent to the BS, thereby preventing the occupation of uplink transmission resources. Moreover, the training framework is further extended to the multiuser case. To utilize crowd intelligence, gossip learning [153] is applied to online learning, where the user exchanges the encoder weights with nearby users and then aggregates the local encoder with the received one.

*5) Standardization:* The current cellular systems, including 4G and 5G, are designed based on implicit feedback mechanisms. However, all mentioned works focus on explicit feedback, that is, full channel information feedback. The autoencoder architecture is introduced in [141], [142] to feedback CSI implicitly, in which the precoding matrix is fed back instead of the whole CSI in CsiNet-like works [23]. The simulation result in [141] shows that the DL-based implicit CSI feedback can reduce at least 25% and 30% of feedback overhead compared with TYPE I and TYPE II feedback codebooks that are adopted by practical systems.

The said autoencoder-based feedback framework needs to change the existing CSI feedback schemes completely, thereby making it difficult to be deployed in the next few years. Hence, developing a DL-based feedback framework, which does not change the existing codebook-based feedback strategy, is essential. The DL-aided codebook enhancement strategy in [143] meets the above requirement. An NN module is added at the BS to refine the channel codeword obtained by the codebook-based feedback. The performance of the original RVQ codebook is greatly improved with the aid of the NN-based enhancement module.

*F. Other Related Works*

In [154], DL is introduced to superimposed coding (SC)-based CSI feedback [155], where the user spreads and superimposes downlink CSI on the uplink user data sequence. Then, the BS recovers the downlink CSI and user data from the received signal with DL-based algorithms. Moreover, 1-bit CS and the partial bidirectional channel reciprocity are introduced by [156].

Feedback safety is considered in [157]. A bias layer is added after the encoder to simulate the attack noise on the air interface. The bias layer is trained in an end-to-end manner to maximize feedback errors and minimize attack noise power jointly. Simulation shows that the destructive effect of the adversarial attack is much higher than that of a jamming attack, which highlights the necessity to design an anti-attack method for DL-based algorithms.

LB-SciFi proposed by [158] is a novel feedback framework for multiuser MIMO in wireless local area networks. An autoencoder is used to compress the CSI in 802.11 protocols to lower airtime overhead and improve system spectral efficiency. Experiments in a wireless testbed show that the DL-based method can offer a 73% reduction in airtime overhead and a 69% increase in system throughput compared with the feedback protocol adopted by 802.11.

DL-based feedback is applied to RIS-assisted wireless systems by [159], where the user estimates the downlink CSI, including the channels of the BS-user, BS-RIS, and RIS-user. Given the substantial RIS element number, the dimension of the phase shift matrix is very high, making the overhead of feeding back the phase shift matrix unaffordable. Therefore, phase shift is fed back by an autoencoder. The main difference between reference [159] and other works is that the matrix fed back to the BS is the RIS phase shift instead of the downlink CSI.

Multitask learning is applied to the CSI feedback of multiple scenarios in [160]. The CSI feedback in different scenarios is regarded as different tasks. In Fig. 33, the user compresses the CSI using a fixed encoder in all scenarios. The GateNet at the BS is an NN-based classifier and determines which scenario the codeword comes from based on the distribution of the received codeword. Then, the corresponding decoder reconstructs the downlink CSI. The advantage of this method is the low complexity at the user because only an encoder is trained and stored for all scenarios.

A joint CSI compression and sensing framework is proposed by [161], in which CSI amplitude is compressed and recon-

structed by an autoencoder, and the sensing results are determined by the received codeword instead of the reconstructed CSI. The experiment results show that the classification accuracy of sensing tasks is comparable with that of the method that sends back CSI amplitude without compression.

## V. Future Research Directions

To accelerate the deployment of DL-based CSI feedback in future communication systems, many challenges must be tackled. Some of them are listed here.

### A. CSI Datasets from Realistic Systems

DL relies heavily on datasets. However, only simulated datasets are available. Most existing works use the datasets generated in [23], which adopts the COST 2100 channel model [162], to evaluate the performance of the proposed NNs. The remaining works, such as [74], [92], [143], [163], generate the CSI samples by the QuaDRiGa software. The NNs proposed for CSI feedback can obtain an excellent performance in some datasets. However, whether the NNs can perform well in other datasets is not clear. Therefore, how to check the robustness of the developed CSI feedback methods becomes critical.

No measured CSI samples except [145] are used to train and evaluate DL-based CSI feedback. The simulated CSI samples are generated by software, in which a certain channel distribution is adopted. The predefined channel distribution cannot exactly describe the characteristics of realistic systems. Although realistic datasets are introduced to DL-based CSI feedback in [145], [158], the channel environments are very simple, and the BS is equipped with several transmit antennas, which is far from the practical systems. Therefore, DL-based CSI feedback needs to be tested using realistic and complicated channel datasets.

Moreover, collecting CSI datasets from practical systems is difficult. The feedback NNs in [137], [163] are trained utilizing the uplink CSI samples based on the distribution reciprocity, which may not hold on in all systems. A dataset collection protocol should define how to select some appropriate users to transmit CSI samples, when to send back the CSI samples, and how to reduce the transmission overhead.

### B. Tradeoff between Performance and Complexity

Fig. 30 shows that the accuracy is improved usually at the expense of NN complexity. For example, the numbers of encoder FLOPs of CsiNet [23] and ConvCsiNet [40] are 0.56 M and 58.52 M when CR is 1/16, respectively. In this case, NMSE improvement is approximately 5 dB for the indoor scenario. This complexity increase is not affordable for the user with limited computational power. Although NN compression techniques can greatly reduce the NN complexity [133], NN complexity remains too high for users with extremely limited computational power, such as Internet of Things sensors. Therefore, NN complexity should be further reduced at the expense of performance, leading to a tradeoff between performance and complexity. The user with enough computational power can be equipped with a powerful NN and

the user with limited computational power needs a lightweight NN. For a certain user, the available computational power varies dynamically. Therefore, the feedback NNs need to be executable at different widths (that is, neuron/channel number in an FC/convolutional layer) to permit a performance complexity tradeoff during inference [164]. The NNs for the user with limited computational power are part of the entire NNs, and the BS transmits the partial NNs according to the user's computational power.

### C. Generalization

NNs are trained with the CSI samples following a certain distribution, which is determined by the propagation environment. However, the environment cannot always be stable [165]. The users do not always stay in a fixed cell and may move to different cells. The environment of a cell inevitably changes over time. Therefore, how to generate an NN with high generalization is one of the major challenges in DL-based CSI feedback. Two potential methods can be used to tackle this challenge. The first method is to build an NN with high generalization by carefully designing the training datasets to cover the most channel distributions. A deep generation model can be used to generate the CSI samples following a certain distribution, such as in [166]. The second potential solution is online training, but it needs to collect plenty of CSI samples, leading to an extra transmission overhead. Therefore, the CSI samples need to be sent to the BS selectively using methods such as the coreset selection algorithm in [167]. Moreover, domain adaptation techniques can be applied to reduce the dataset requirement further and accelerate training.

### D. Effect on Standardization

DL-based CSI feedback is incorporated into the 3GPP R18 study item [22]. The effect of DL-based CSI feedback on the existing standard needs to be evaluated. First, how many system gains (instead of feedback accuracy, such as NMSE) can be achieved should be provided through the link- and system-level simulation compared with the existing TYPE I and TYPE II codebook-based CSI feedback. Second, DL-based algorithms are different from the conventional algorithms and pose new requirements for the systems. Third, the evolution of the DL-based feedback framework needs to be discussed further. The existing standard cannot be totally changed and can only be revised. For example, explicit feedback is fully different from the existing feedback framework and is difficult to be deployed in 5G-Advanced.

### E. High-speed Scenario

Mobility of users becomes higher in the future. Channel aging is unavoidable and leads to a large drop in system performance. However, few DL-based feedback works consider the high-speed scenario. In this scenario, the decoder of the BS must not only be able to reconstruct the CSI accurately but also predict the future CSI to reduce the influence of channel aging [168]. The DL-based feedback method should be designed by considering the characteristics of the high-speed scenario. For

TABLE VI
THE WORKS WITH OPEN SOURCE CODE IN DL-BASED CSI FEEDBACK.

| Methods | Links |
|---|---|
| CsiNet [23] | https://github.com/sydney222/Python_CsiNet |
| CRNet [43] | https://github.com/Kylin9511/CRNet |
| DS-NLCsiNet [45] | https://github.com/yuxt1999/DS-NLCsiNet |
| CLNet [49] | https://github.com/SIJIEJI/CLNet |
| DCRNet [51] | https://github.com/recusant7/DCRNet |
| TransNet [56] | https://github.com/Treedy2020/TransNet |
| SALDR [65] | https://github.com/XS96/SALDR |
| P-SRNet [68] | https://github.com/MoliaChen/SRNet |
| CsiNet-LSTM [85] | https://www.ecsponline.com/goods.php?id=205629 |
| ConvlstmCsiNet [97] | https://github.com/Aries-LXY/ConvlstmCsiNet |
| DualNet [88] | https://github.com/DLinWL/Bi-Directional-Channel-Reciprocity |
| [120] | https://github.com/foadsohrabi/DL-DSC-FDD-Massive-MIMO |
| CHNet [128] | https://github.com/ch28/CHNet |
| ACRNet [134] | https://github.com/Kylin9511/ACRNet |
| PSCDN [159] | https://github.com/xian-hua/PSCDN/ |

example, the user in this scenario usually moves on a fixed path, such as in rails, because the environment around the fixed path is usually long-term stable.

*F. Other Emerging Techniques*

Many new techniques, such as RIS [169] and extra-large scale massive MIMO [170], are introduced to communications and regarded as potential key techniques in 6G. CSI feedback combined with these new techniques needs to be explored. For example, CSI acquisition (including feedback) is a major challenge of the RIS-assisted communication systems, in which the channel reciprocity may not hold on even in time-division duplexing systems [171]. The CSI dimension greatly increases because of the introduction of the RIS with a large element number. If the RIS has $100 \times 100$ elements and the user is equipped with a single antenna, the channel between the RIS and the user is $10,000 \times 1$, which is much larger than that in the current massive MIMO systems. Therefore, a more efficient DL framework needs to be explored to tackle the challenges introduced by these new techniques.

*G. Open Source Dataset and Code*

Table VI shows the most DL-based CSI feedback works with open source code. Reproducible research is essential in DL-based algorithms. Open source can make the works more convincing and help accelerate research. Therefore, more open source works are welcome. Wireless-Intelligence is a public channel dataset library, which has been built for DL-based wireless communications [172]. This library contains many channel datasets that satisfy the 3GPP standard. However, channel datasets measured from the practical massive MIMO systems are not publicly available. An open practical channel dataset is essential and urgent to accelerate the study of DL-based CSI feedback.

## VI. CONCLUSION

In this paper, an overview of DL-based CSI feedback has been provided. First, the basic DL concepts and representative NN architectures widely used in DL-based feedback have been briefly introduced to guide beginners. Then, the existing works have been divided into six different categories, and each has been comprehensively introduced and discussed. Finally, the new challenges and potential directions for future research in DL-based CSI feedback, especially focusing on practical deployment and standardization, have been elaborated.

## REFERENCES

[1] X. Lin, J. Li, R. Baldemair *et al.*, "5G new radio: Unveiling the essentials of the next generation wireless access technology," *IEEE Commun. Standards Mag.*, vol. 3, no. 3, pp. 30–37, Sep. 2019.

[2] M. Shafi, A. F. Molisch, P. J. Smith *et al.*, "5G: A tutorial overview of standards, trials, challenges, deployment, and practice," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 6, pp. 1201–1221, June 2017.

[3] A. Ghosh, A. Maeder, M. Baker, and D. Chandramouli, "5G evolution: A view on 5G cellular technology beyond 3GPP Release 15," *IEEE Access*, vol. 7, pp. 127 639–127 651, 2019.

[4] B. Bertenyi, "5G evolution: What's next?" *IEEE Wireless Commun.*, vol. 28, no. 1, pp. 4–8, Feb. 2021.

[5] X. Lin, "An overview of 5G Advanced evolution in 3GPP release 18," *arXiv preprint arXiv:2201.01358*, 2022. [Online]. Available: https://arxiv.org/abs/2201.01358

[6] J. Hoydis, F. A. Aoudia, A. Valcarce, and H. Viswanathan, "Toward a 6G AI-native air interface," *IEEE Commun. Mag.*, vol. 59, no. 5, pp. 76–81, May 2021.

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

[8] T. Wang, C.-K. Wen, H. Wang, F. Gao, T. Jiang, and S. Jin, "Deep learning for wireless physical layer: Opportunities and challenges," *China Commun.*, vol. 14, no. 11, pp. 92–111, Nov. 2017.

[9] S. Liu, T. Wang, and S. Wang, "Toward intelligent wireless communications: Deep learning-based physical layer technologies," *Digit. Commun. Netw.*, vol. 7, no. 4, pp. 589–597, Nov. 2021.

[10] A. Alkhateeb, S. Alex, P. Varkey, Y. Li, Q. Qu, and D. Tujkovic, "Deep learning coordinated beamforming for highly-mobile millimeter wave systems," *IEEE Access*, vol. 6, pp. 37 328–37 348, 2018.

[11] W. Xia, G. Zheng, Y. Zhu, J. Zhang, J. Wang, and A. P. Petropulu, "A deep learning framework for optimization of MISO downlink beamforming," *IEEE Trans. Commun.*, vol. 68, no. 3, pp. 1866–1880, March 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] X. Gao, S. Jin, C.-K. Wen, and G. Y. Li, "ComNet: Combination of deep learning and expert knowledge in OFDM receivers," *IEEE Commun. Lett.*, vol. 22, no. 12, pp. 2627–2630, Dec. 2018.

[14] P. Jiang, T. Wang, B. Han *et al.*, "AI-aided online adaptive OFDM receiver: Design and experimental results," *IEEE Trans. Wireless Commun.*, vol. 20, no. 11, pp. 7655–7668, Nov. 2021.

[15] T. O'shea and J. Hoydis, "An introduction to deep learning for the physical layer," *IEEE Trans. Cogn. Commun. Netw*, vol. 3, no. 4, pp. 563–575, Dec. 2017.

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

[17] S. Dörner, S. Cammerer, J. Hoydis, and S. Ten Brink, "Deep learning based communication over the air," *IEEE J. Sel. Topics Signal Process.*, vol. 12, no. 1, pp. 132–143, Feb. 2018.

[18] L. Lu, G. Y. Li, A. L. Swindlehurst, A. Ashikhmin, and R. Zhang, "An overview of massive MIMO: Benefits and challenges," *IEEE J. Sel. Topics Signal Process.*, vol. 8, no. 5, pp. 742–758, Oct. 2014.

[19] X. Rao and V. K. N. Lau, "Distributed compressive CSIT estimation and feedback for FDD multi-user massive MIMO systems," *IEEE Trans. Signal Process.*, vol. 62, no. 12, pp. 3261–3271, June 2014.

[20] D. J. Love, R. W. Heath, V. K. N. Lau, D. Gesbert, B. D. Rao, and M. Andrews, "An overview of limited feedback in wireless communication systems," *IEEE J. Sel. Areas Commun.*, vol. 26, no. 8, pp. 1341–1365, Oct. 2008.

[21] Z. Qin, J. Fan, Y. Liu, Y. Gao, and G. Y. Li, "Sparse representation for wireless communications: A compressive sensing approach," *IEEE Signal Process. Mag.*, vol. 35, no. 3, pp. 40–58, May 2018.

[22] 3GPP RP-213599, "New SI: Study on artificial intelligence (AI)/Machine Learning (ML) for NR air interface," Moderator (Qualcomm), Tech. Rep., Dec. 2021, accessed on May 1, 2022. [Online]. Available: https://www.3gpp.org/ftp/tsg_ran/TSG_RAN/TSGR_94e/Docs/RP-213599.zip

[23] 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.

[24] N. Jindal, "MIMO broadcast channels with finite-rate feedback," *IEEE Trans. Inf. Theory*, vol. 52, no. 11, pp. 5045–5060, Nov. 2006.

[25] A. Beck and M. Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems," *SIAM. J. Imaging Science*, vol. 2, no. 1, pp. 183–202, 2009.

[26] A. L. Maas, A. Y. Hannun, A. Y. Ng *et al.*, "Rectifier nonlinearities improve neural network acoustic models," in *Proc. ICML*, 2013, pp. 1–6.

[27] Z. Wu, C. Shen, and A. Van Den Hengel, "Wider or deeper: Revisiting the resnet model for visual recognition," *Pattern Recognit.*, vol. 90, pp. 119–133, June 2019.

[28] A. Araujo, W. Norris, and J. Sim, "Computing receptive fields of convolutional neural networks," *Distill*, vol. 4, no. 11, p. e21, 2019.

[29] X. Ding, X. Zhang, Y. Zhou, J. Han, G. Ding, and J. Sun, "Scaling up your kernels to 31x31: Revisiting large kernel design in CNNs," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2022.

[30] C. Olah, "Understanding LSTM networks," http://colah.github.io/posts/2015-08-Understanding-LSTMs/, 2015, accessed on May 1, 2022.

[31] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997.

[32] R. C. Staudemeyer and E. R. Morris, "Understanding LSTM–a tutorial into long short-term memory recurrent neural networks," *arXiv preprint arXiv:1909.09586*, 2019. [Online]. Available: https://arxiv.org/abs/1909.09586

[33] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013. [Online]. Available: https://arxiv.org/abs/1312.6114

[34] I. Goodfellow, J. Pouget-Abadie, M. Mirza *et al.*, "Generative adversarial nets," in *Proc. 28th Adv. Neural Inf. Process. Syst. (NIPS)*, vol. 27, 2014.

[35] A. Vaswani, N. Shazeer, N. Parmar *et al.*, "Attention is all you need," in *Proc. 31st Adv. Neural Inf. Process. Syst. (NIPS)*, vol. 30, 2017.

[36] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," in *Proc. ICLR*, 2015.

[37] K. Xu, J. Ba, R. Kiros *et al.*, "Show, attend and tell: Neural image caption generation with visual attention," in *Proc. ICML*, 2015, pp. 2048–2057.

[38] X. Yang, "An overview of the attention mechanisms in computer vision," *J. Phys. Conf. Ser.*, vol. 1693, no. 1, p. 012173, Dec. 2020.

[39] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon, "CBAM: Convolutional block attention module," in *Proc. ECCV*, Sept. 2018, pp. 3–19.

[40] W.-T. Shih, "Study on massive MIMO CSI feedback based on deep learning (in Traditional Chinese)," Master's thesis, National Sun Yat-sen University, 2018, Accessed on May 1, 2022. [Online]. Available: https://hdl.handle.net/11296/pvuea3

[41] J. Guo, C.-K. Wen, S. Jin, and G. Y. Li, "Convolutional neural network-based multiple-rate compressive sensing for massive MIMO CSI feedback: Design, simulation, and analysis," *IEEE Trans. Wireless Commun.*, vol. 19, no. 4, pp. 2827–2840, Apr. 2020.

[42] Q. Cai, C. Dong, and K. Niu, "Attention model for massive MIMO CSI compression feedback and recovery," in *Proc. IEEE Wireless Commun. Netw. Conf. (WCNC)*, 2019, pp. 1–5.

[43] Z. Lu, J. Wang, and J. Song, "Multi-resolution CSI feedback with deep learning in massive MIMO system," in *Proc. IEEE Int. Conf. Commun. (ICC)*, 2020, pp. 1–6.

[44] Q. Li, A. Zhang, P. Liu, J. Li, and C. Li, "A novel CSI feedback approach for massive mimo using LSTM-Attention CNN," *IEEE Access*, vol. 8, pp. 7295–7302, 2020.

[45] X. Yu, X. Li, H. Wu, and Y. Bai, "DS-NLCsiNet: Exploiting non-local neural networks for massive MIMO CSI feedback," *IEEE Commun. Lett.*, vol. 24, no. 12, pp. 2790–2794, Dec. 2020.

[46] B. Tolba, M. Elsabrouty, M. G. Abdu-Aguye, H. Gacanin, and H. M. Kasem, "Massive MIMO CSI feedback based on generative adversarial network," *IEEE Commun. Lett.*, vol. 24, no. 12, pp. 2805–2808, Dec. 2020.

[47] M. Hussien, K. K. Nguyen, and M. Cheriet, "PRVNet: Variational autoencoders for massive MIMO CSI feedback," *arXiv preprint arXiv:2011.04178*, 2020. [Online]. Available: https://arxiv.org/abs/2011.04178

[48] M. Gao, T. Liao, and Y. Lu, "Fully connected feedforward neural networks based CSI feedback algorithm," *China Commun.*, vol. 18, no. 1, pp. 43–48, Jan. 2021.

[49] S. Ji and M. Li, "CLNet: Complex input lightweight neural network designed for massive MIMO CSI feedback," *IEEE Wireless Commun. Lett.*, vol. 10, no. 10, pp. 2318–2322, Oct. 2021.

[50] Y. Sun, W. Xu, L. Liang, N. Wang, G. Y. Li, and X. You, "A lightweight deep network for efficient CSI feedback in massive MIMO systems," *IEEE Wireless Commun. Lett.*, vol. 10, no. 8, pp. 1840–1844, Aug. 2021.

[51] S. Tang, J. Xia, L. Fan, X. Lei, W. Xu, and A. Nallanathan, "Dilated convolution based CSI feedback compression for massive MIMO systems," *arXiv preprint arXiv:2106.04043*, 2021. [Online]. Available: https://arxiv.org/abs/2106.04043

[52] Y. Zhang, X. Zhang, and Y. Liu, "Deep learning based CSI compression and quantization with high compression ratios in FDD massive MIMO systems," *IEEE Wireless Commun. Lett.*, vol. 10, no. 10, pp. 2101–2105, Oct. 2021.

[53] Z. Hu, J. Guo, G. Liu, H. Zheng, and J. Xue, "MRFNet: A deep learning-based CSI feedback approach of massive MIMO systems," *IEEE Commun. Lett.*, vol. 25, no. 10, pp. 3310–3314, Oct. 2021.

[54] B. Cao, Y. Yang, P. Ran, D. He, and G. He, "ACCsiNet: Asymmetric convolution-based autoencoder framework for massive MIMO CSI feedback," *IEEE Commun. Lett.*, vol. 25, no. 12, pp. 3873–3877, Dec. 2021.

[55] Y. Xu, M. Zhao, S. Zhang, and H. Jin, "DFECsiNet: Exploiting diverse channel features for massive MIMO CSI feedback," in *Proc. 13th WCSP*, 2021, pp. 1–5.

[56] Y. Cui, A. Guo, and C. Song, "TransNet: Full attention network for CSI feedback in FDD massive MIMO system," *IEEE Wireless Commun. Lett.*, pp. 1–1, 2022.

[57] X. Bi, S. Li, C. Yu, and Y. Zhang, "A novel approach using convolutional transformer for massive MIMO CSI feedback," *IEEE Wireless Commun. Lett.*, 2022, Early access.

[58] H. Li, B. Zhang, H. Chang, X. Liang, and X. Gu, "CVLNet: A complex-valued lightweight network for CSI feedback," *IEEE Wireless Commun. Lett.*, 2022, Early access.

[59] Y. Wang, J. Sun, J. Wang *et al.*, "Multi-rate compression for downlink CSI based on transfer learning in FDD massive MIMO systems," in *Proc. IEEE 94th VTC-Fall*, 2021, pp. 1–5.

[60] P. Liang, J. Fan, W. Shen, Z. Qin, and G. Y. Li, "Deep learning and compressive sensing-based CSI feedback in FDD massive MIMO systems," *IEEE Trans. Veh. Technol.*, vol. 69, no. 8, pp. 9217–9222, Aug. 2020.

[61] Z. Lu, J. Wang, and J. Song, "Binary neural network aided CSI feedback in massive MIMO system," *IEEE Wireless Commun. Lett.*, vol. 10, no. 6, pp. 1305–1308, June 2021.

[62] B. Cheng, J. Zhao, and Y. Hu, "Multi-scale and multi-channel networks for CSI feedback in massive MIMO system," *J. Comput. Commun.*, vol. 9, no. 10, pp. 132–141, Oct. 2021.

[63] Q. Yang, M. B. Mashhadi, and D. Gündüz, "Deep convolutional compression for massive MIMO CSI feedback," in *Proc. IEEE 29th Int. Workshop Mach. Learn. Signal Process. (MLSP)*, 2019, pp. 1–6.

[64] G. Fan, J. Sun, G. Gui, H. Gacanin, B. Adebisi, and T. Ohtsuki, "Fully convolutional neural network based CSI limited feedback for FDD massive MIMO systems," *IEEE Trans. on Cogn. Commun. Netw.*, 2021, Early access.

[65] X. Song, J. Wang, J. Wang *et al.*, "SALDR: Joint self-attention learning and dense refine for massive MIMO CSI feedback with multiple compression ratio," *IEEE Wireless Commun. Lett.*, vol. 10, no. 9, pp. 1899–1903, Sept. 2021.

[66] Y. Xu, M. Yuan, and M.-O. Pun, "Transformer empowered CSI feedback for massive MIMO systems," in *26th Wireless Opt. Commun. Conf. (WOCC)*, 2021, pp. 157–161.

[67] S. Jo, J. Lee, and J. So, "Deep learning-based massive multiple-input multiple-output channel state information feedback with data normalisation using clipping," *Electron. Lett.*, vol. 57, no. 3, pp. 151–154, Feb. 2021.

[68] X. Chen, C. Deng, B. Zhou, H. Zhang, G. Yang, and S. Ma, "High-accuracy CSI feedback with super-resolution network for massive MIMO systems," *IEEE Wireless Commun. Lett.*, vol. 11, no. 1, pp. 141–145, Jan. 2022.

[69] Y. Wang, X. Chen, H. Yin, and W. Wang, "Learnable sparse transformation-based massive MIMO CSI recovery network," *IEEE Commun. Lett.*, vol. 24, no. 7, pp. 1468–1471, July 2020.

[70] J. Guo, L. Wang, F. Li, and J. Xue, "CSI feedback with model-driven deep learning of massive MIMO systems," *IEEE Commun. Lett.*, pp. 1–1, 2021, Early access.

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

[72] X. Ding, Y. Guo, G. Ding, and J. Han, "ACNet: Strengthening the kernel skeletons for powerful cnn via asymmetric convolution blocks," in *Proc. IEEE/CVF Interface Conf. Comput. Vis. (ICCV)*, 2019, pp. 1911–1920.

[73] Z. Zhang, Y. Zheng, C. Gan, and Q. Zhu, "Massive MIMO CSI reconstruction using CNN-LSTM and attention mechanism," *IET Commun.*, vol. 14, no. 18, pp. 3089–3094, 2020.

[74] D. J. Ji and D.-H. Cho, "ChannelAttention: Utilizing attention layers for accurate massive MIMO channel feedback," *IEEE Wireless Commun. Lett.*, vol. 10, no. 5, pp. 1079–1082, May 2021.

[75] H. Zhao, J. Jia, and V. Koltun, "Exploring self-attention for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2020, pp. 10073–10082.

[76] G. Huang, Z. Liu, L. V. D. Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2017, pp. 2261–2269.

[77] M. Abadi, P. Barham, J. Chen *et al.*, "TensorFlow: A system for large-scale machine learning," in *Proc. 12th OSDI*, 2016, pp. 265–283.

[78] A. Paszke, S. Gross, F. Massa *et al.*, "PyTorch: An imperative style, high-performance deep learning library," *NeurIPS*, vol. 32, 2019.

[79] J. Sola and J. Sevilla, "Importance of input data normalization for the application of neural networks to complex industrial problems," *IEEE Trans.Nucl. Sci.*, vol. 44, no. 3, pp. 1464–1468, June 1997.

[80] 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.

[81] N. Shlezinger, J. Whang, Y. C. Eldar, and A. G. Dimakis, "Model-based deep learning," *arXiv preprint arXiv:2012.08405*, 2020. [Online]. Available: https://arxiv.org/abs/2012.08405

[82] K. Gregor and Y. LeCun, "Learning fast approximations of sparse coding," in *Proc. 27th ICLR*, 2010, pp. 399–406.

[83] J. Liu, X. Chen, Z. Wang, and W. Yin, "ALISTA: Analytic weights are as good as learned weights in LISTA," in *Proc. ICLR*, 2019. [Online]. Available: https://openreview.net/forum?id=B1lnzn0ctQ

[84] D. Donoho, "De-noising by soft-thresholding," *IEEE Trans. Inf. Theory*, vol. 41, no. 3, pp. 613–627, May 1995.

[85] T. Wang, C.-K. Wen, S. Jin, and G. Y. Li, "Deep learning-based CSI feedback approach for time-varying massive MIMO channels," *IEEE Wireless Commun. Lett.*, vol. 8, no. 2, pp. 416–419, Apr. 2019.

[86] C. Lu, W. Xu, H. Shen, J. Zhu, and K. Wang, "MIMO channel information feedback using deep recurrent network," *IEEE Commun. Lett.*, vol. 23, no. 1, pp. 188–191, Jan. 2019.

[87] S. Hong, S. Jo, and J. So, "Machine learning-based adaptive CSI feedback interval," *ICT Express*, 2021, Early access.

[88] Z. Liu, L. Zhang, and Z. Ding, "Exploiting bi-directional channel reciprocity in deep learning for low rate massive MIMO CSI feedback," *IEEE Wireless Commun. Lett.*, vol. 8, no. 3, pp. 889–892, June 2019.

[89] T. Wang, "Research on key technology of massive MIMO channel feedback for intelligent communications (in Chinese)," Master's thesis, Southeast University, 2019, Accessed on May 1, 2022. [Online]. Available: https://kns.cnki.net/kcms/detail/detail.aspx?dbcode=CMFD&dbname=CMFD202001&filename=1020612377.nh&uniplatform=NZKPT&v=bxKIu4EzKygyTCR6jG9Js0YheR7qc2TAD5tx_MBkZkNnopR6AiW0Fe4CRDB-NTth

[90] Y. Liu and O. Simeone, "HyperRNN: Deep learning-aided downlink CSI acquisition via partial channel reciprocity for FDD massive MIMO," in *Proc. IEEE 22nd Int. Workshop Signal Process. Adv. Wireless Commun. (SPAWC)*, 2021, pp. 31–35.

[91] J. Wang, G. Gui, T. Ohtsuki, B. Adebisi, H. Gacanin, and H. Sari, "Compressive sampled CSI feedback method based on deep learning for FDD massive MIMO systems," *IEEE Trans. Commun.*, vol. 69, no. 9, pp. 5873–5885, Sept. 2021.

[92] J. Guo, X. Yang, C.-K. Wen, S. Jin, and G. Y. Li, "DL-based CSI feedback and cooperative recovery in massive MIMO," *arXiv preprint arXiv:2003.03303*, 2020. [Online]. Available: https://arxiv.org/abs/2003.03303

[93] M. B. Mashhadi, Q. Yang, and D. Gündüz, "Distributed deep convolutional compression for massive MIMO CSI feedback," *IEEE Trans. Wireless Commun.*, vol. 20, no. 4, pp. 2621–2633, Apr. 2021.

[94] H. Yin and D. Gesbert, "A partial channel reciprocity-based codebook for wideband FDD massive MIMO," *IEEE Trans. Wireless Commun.*, 2022, Early access.

[95] M. Stojanovic, J. Proakis, and J. Catipovic, "Analysis of the impact of channel estimation errors on the performance of a decision-feedback equalizer in fading multipath channels," *IEEE Trans. Commun.*, vol. 43, no. 2/3/4, pp. 877–886, Feb./March/April 1995.

[96] Z. Liu, M. del Rosario, and Z. Ding, "A Markovian model-driven deep learning framework for massive MIMO CSI feedback," *IEEE Trans. Wireless Commun.*, vol. 21, no. 2, pp. 1214–1228, Feb. 2022.

[97] X. Li and H. Wu, "Spatio-temporal representation with deep neural recurrent network in MIMO CSI feedback," *IEEE Wireless Commun. Lett.*, vol. 9, no. 5, pp. 653–657, May 2020.

[98] Z. Huang and M. Li, "Research on channel feedback algorithm in UAV inspection communication subsystem of smart grid," in *Proc. 2nd IEEE ISCEIC*, 2021, pp. 236–240.

[99] Y.-C. Lin, Z. Liu, T.-S. Lee, and Z. Ding, "Deep learning phase compression for MIMO CSI feedback by exploiting FDD channel reciprocity," *IEEE Wireless Commun. Lett.*, vol. 10, no. 10, pp. 2200–2204, Oct. 2021.

[100] Z. Liu, L. Zhang, and Z. Ding, "An efficient deep learning framework for low rate massive MIMO CSI reporting," *IEEE Trans. Commun.*, vol. 68, no. 8, pp. 4761–4772, Aug. 2020.

[101] D. Ha, A. Dai, and Q. V. Le, "Hypernetworks," *arXiv preprint arXiv:1609.09106*, 2016. [Online]. Available: https://arxiv.org/abs/1609.09106

[102] M. Latva-aho, K. Leppänen, F. Clazzer, and A. Munari, "Key drivers and research challenges for 6G ubiquitous wireless intelligence," *White Paper*, 2020. [Online]. Available: http://jultika.oulu.fi/Record/isbn978-952-62-2354-4

[103] X. Du and A. Sabharwal, "Massive MIMO channels with inter-user angle correlation: Open-access dataset, analysis and measurement-based validation," *IEEE Trans. Veh. Technol.*, vol. 71, no. 2, pp. 1602–1616, Feb. 2022.

[104] Y. Jang, G. Kong, M. Jung, S. Choi, and I.-M. Kim, "Deep autoencoder based CSI feedback with feedback errors and feedback delay in FDD massive MIMO systems," *IEEE Wireless Commun. Lett.*, vol. 8, no. 3, pp. 833–836, June 2019.

[105] C. Lu, W. Xu, S. Jin, and K. Wang, "Bit-level optimized neural network for multi-antenna channel quantization," *IEEE Wireless Commun. Lett.*, vol. 9, no. 1, pp. 87–90, Jan. 2020.

[106] T. Chen, J. Guo, S. Jin, C.-K. Wen, and G. Y. Li, "A novel quantization method for deep learning-based massive MIMO CSI feedback," in *Proc. IEEE Global Conf. Signal Inf. Process. (GlobalSIP)*, 2019, pp. 1–5.

[107] C. Recommendation, "Pulse code modulation (PCM) of voice frequencies," in *ITU*, 1988.

[108] D. Marpe, H. Schwarz, and T. Wiegand, "Context-based adaptive binary arithmetic coding in the H.264/AVC video compression standard," *IEEE Trans. Circuits and Syst. for Video Tech.*, vol. 13, no. 7, pp. 620–636, July 2003.

[109] S. Ravula and S. Jain, "Deep autoencoder-based massive MIMO CSI feedback with quantization and entropy coding," in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, 2021, pp. 1–6.

[110] J. Ballé, D. Minnen, S. Singh, S. J. Hwang, and N. Johnston, "Variational image compression with a scale hyperprior," in *Proc. ICLR*, 2018.

[111] K. Kong, W.-J. Song, and M. Min, "Knowledge distillation-aided end-to-end learning for linear precoding in multiuser MIMO downlink

systems with finite-rate feedback," *IEEE Trans. Veh. Technol.*, vol. 70, no. 10, pp. 11 095–11 100, Oct. 2021.

[112] G. Hinton, O. Vinyals, J. Dean *et al.*, "Distilling the knowledge in a neural network," in *Proc. Adv. Neural Inf. Process. Syst. Workshops (NIPSW)*, 2015.

[113] J. Guo, T. Chen, C.-K. Wen, S. Jin, G. Y. Li, X. Wang, and X. Hou, "Deep learning for joint channel estimation and feedback in massive MIMO systems," *Digit. Commun. Netw.*, 2022, Language/Minor revision.

[114] Y. Sun, W. Xu, L. Fan, G. Y. Li, and G. K. Karagiannidis, "AnciNet: An efficient deep learning approach for feedback compression of estimated CSI in massive MIMO systems," *IEEE Wireless Commun. Lett.*, vol. 9, no. 12, pp. 2192–2196, Dec. 2020.

[115] J. Guo, C.-K. Wen, and S. Jin, "CAnet: Uplink-aided downlink channel acquisition in FDD massive MIMO using deep learning," *IEEE Trans. Commun.*, vol. 70, no. 1, pp. 199–214, Jan. 2022.

[116] ——, "Deep learning-based CSI feedback for beamforming in single- and multi-cell massive MIMO systems," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 7, pp. 1872–1884, July 2021.

[117] Q. Sun, H. Zhao, J. Wang, and W. Chen, "Deep learning-based joint CSI feedback and hybrid precoding in FDD mmWave massive MIMO systems," *Entropy*, vol. 24, no. 4, p. 441, 2022.

[118] A.-A. Lee, Y.-S. Wang, and Y.-W. P. Hong, "Deep CSI compression and coordinated precoding for multicell downlink systems," in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, 2020, pp. 1–6.

[119] J. Jang, H. Lee, S. Hwang, H. Ren, and I. Lee, "Deep learning-based limited feedback designs for MIMO systems," *IEEE Wireless Commun. Lett.*, vol. 9, no. 4, pp. 558–561, Apr. 2020.

[120] F. Sohrabi, K. M. Attiah, and W. Yu, "Deep learning for distributed channel feedback and multiuser precoding in FDD massive MIMO," *IEEE Trans. Wireless Commun.*, vol. 20, no. 7, pp. 4044–4057, July 2021.

[121] M. Boloursaz Mashhadi and D. Gündüz, "Deep learning for massive mimo channel state acquisition and feedback," *J. Indian Inst. Sci.*, vol. 100, no. 2, pp. 369–382, 2020.

[122] X. Ma and Z. Gao, "Data-driven deep learning to design pilot and channel estimator for massive MIMO," *IEEE Trans. Veh. Technol.*, vol. 69, no. 5, pp. 5677–5682, March 2020.

[123] M. B. Mashhadi and D. Gündüz, "Pruning the pilots: Deep learning-based pilot design and channel estimation for MIMO-OFDM systems," *IEEE Trans. Wireless Commun.*, vol. 20, no. 10, pp. 6315–6328, Oct. 2021.

[124] Z. Wang and A. C. Bovik, "Mean squared error: Love it or leave it? a new look at signal fidelity measures," *IEEE Signal Process. Mag.*, vol. 26, no. 1, pp. 98–117, Jan. 2009.

[125] O. Somekh, B. M. Zaidel, and S. Shamai, "Sum rate characterization of joint multiple cell-site processing," *IEEE Trans. Inf. Theory*, vol. 53, no. 12, pp. 4473–4497, Dec. 2007.

[126] R. Bhagavatula and R. W. Heath, "Adaptive limited feedback for sum-rate maximizing beamforming in cooperative multicell systems," *IEEE Trans. Signal Process.*, vol. 59, no. 2, pp. 800–811, Feb. 2011.

[127] Y. Wang, Y. Zhang, J. Sun, G. Gui, T. Ohtsuki, and F. Adachi, "A novel compression CSI feedback based on deep learning for FDD massive MIMO systems," in *Proc. IEEE Wireless Commun. Netw. Conf. (WCNC)*, 2021, pp. 1–5.

[128] X. Liang, H. Chang, H. Li, X. Gu, and L. Zhang, "Changeable rate and novel quantization for CSI feedback based on deep learning," *arXiv preprint arXiv:2202.13627*, 2022. [Online]. Available: https://arxiv.org/abs/2202.13627

[129] S. Jo and J. So, "Adaptive lightweight CNN-based CSI feedback for massive MIMO systems," *IEEE Wireless Commun. Lett.*, vol. 10, no. 12, pp. 2776–2780, Dec. 2021.

[130] H. Ye, F. Gao, J. Qian, H. Wang, and G. Y. Li, "Deep learning-based denoise network for CSI feedback in FDD massive MIMO systems," *IEEE Commun. Lett.*, vol. 24, no. 8, pp. 1742–1746, Aug. 2020.

[131] H. Chang, X. Liang, H. Li, J. Shen, X. Gu, and L. Zhang, "Deep learning-based bitstream error correction for CSI feedback," *IEEE Wireless Commun. Lett.*, vol. 10, no. 12, pp. 2828–2832, Dec. 2021.

[132] M. B. Mashhadi, Q. Yang, and D. Gündüz, "CNN-based analog CSI feedback in FDD MIMO-OFDM systems," in *Proc. IEEE 45th Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, 2020, pp. 8579–8583.

[133] J. Guo, J. Wang, C.-K. Wen, S. Jin, and G. Y. Li, "Compression and acceleration of neural networks for communications," *IEEE Wireless Commun.*, vol. 27, no. 4, pp. 110–117, Aug. 2020.

[134] Z. Lu, X. Zhang, H. He, J. Wang, and J. Song, "Binarized aggregated network with quantization: Flexible deep learning deployment for CSI

feedback in massive MIMO system," *IEEE Trans. Wireless Commun.*, 2022, Early access.

[135] H. Tang, J. Guo, M. Matthaiou, C.-K. Wen, and S. Jin, "Knowledge-distillation-aided lightweight neural network for massive MIMO CSI feedback," in *Proc. IEEE 94th VTC-Fall*, 2021, pp. 1–5.

[136] H. Sun, Z. Zhao, X. Fu, and M. Hong, "Limited feedback double directional massive MIMO channel estimation: From low-rank modeling to deep learning," in *Proc. IEEE 19th Int. Workshop Signal Process. Adv. Wireless Commun. (SPAWC)*, 2018, pp. 1–5.

[137] N. Song and T. Yang, "Machine learning enhanced CSI acquisition and training strategy for FDD massive MIMO," in *Proc. IEEE Wireless Commun. Netw. Conf. (WCNC)*, 2021, pp. 1–6.

[138] J. Zeng, J. Sun, G. Gui *et al.*, "Downlink CSI feedback algorithm with deep transfer learning for FDD massive MIMO systems," *IEEE Trans. Cogn. Commun. Netw.*, vol. 7, no. 4, pp. 1253–1265, Dec. 2021.

[139] J. Jiang, R. Han, B. Liu, and D. Feng, "Federated learning-based codebook design for massive MIMO communication system," in *Proc. International Conference on Natural Computation, Fuzzy Systems and Knowledge Discovery*. Springer, 2021, pp. 1198–1205.

[140] J. Guo, Y. Zuo, C.-K. Wen, and S. Jin, "User-centric online gossip training for autoencoder-based CSI feedback," *IEEE J. Sel. Topics Signal Process.*, 2022, Early access.

[141] M. Chen, J. Guo, C.-K. Wen, S. Jin, G. Y. Li, and A. Yang, "Deep learning-based implicit CSI feedback in massive MIMO," *IEEE Trans. Commun.*, vol. 70, no. 2, pp. 935–950, Feb. 2022.

[142] W. Liu, W. Tian, H. Xiao, S. Jin, X. Liu, and J. Shen, "EVCsiNet: Eigenvector-based CSI feedback under 3GPP link-level channels," *IEEE Wireless Commun. Lett.*, vol. 10, no. 12, pp. 2688–2692, Dec. 2021.

[143] J. Guo, C.-K. Wen, M. Chen, and S. Jin, "Environment knowledge-aided massive MIMO feedback codebook enhancement using artificial intelligence," 2022, Early access.

[144] H. Xiao, Z. Wang, W. Tian *et al.*, "AI enlightens wireless communication: Analyses, solutions and opportunities on CSI feedback," *China Commun.*, vol. 18, no. 11, pp. 104–116, Nov. 2021.

[145] J. Guo, X. Li, M. Chen *et al.*, "AI enabled wireless communications with real channel measurements: Channel feedback," *J. Commun. Inf. Netw.*, vol. 5, no. 3, pp. 310–317, Sept. 2020.

[146] J. Xu, B. Ai, N. Wang, and W. Chen, "Deep joint source-channel coding for CSI feedback: An end-to-end approach," *arXiv preprint arXiv:2203.16005*, 2022. [Online]. Available: https://arxiv.org/abs/2203.16005

[147] Z. Cao, W.-T. Shih, J. Guo, C.-K. Wen, and S. Jin, "Lightweight convolutional neural networks for CSI feedback in massive MIMO," *IEEE Commun. Lett.*, vol. 25, no. 8, pp. 2624–2628, Aug. 2021.

[148] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer, "SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and <0.5 mb model size," *arXiv preprint arXiv:1602.07360*, 2016. [Online]. Available: https://arxiv.org/abs/1602.07360

[149] X. Zhang, X. Zhou, M. Lin, and J. Sun, "ShuffleNet: An extremely efficient convolutional neural network for mobile devices," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2018, pp. 6848–6856.

[150] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 40, no. 4, pp. 834–848, Apr. 2018.

[151] Y. Yang, F. Gao, Z. Zhong, B. Ai, and A. Alkhateeb, "Deep transfer learning-based downlink channel prediction for FDD massive MIMO systems," *IEEE Trans. Commun.*, vol. 68, no. 12, pp. 7485–7497, Dec. 2020.

[152] M. B. Mashhadi, M. Jankowski, T.-Y. Tung, S. Kobus, and D. Gündüz, "Federated mmwave beam selection utilizing LIDAR data," *IEEE Wireless Commun. Lett.*, vol. 10, no. 10, pp. 2269–2273, Oct. 2021.

[153] L. Giaretta and Š. Girdzijauskas, "Gossip learning: Off the beaten path," in *Proc. IEEE Int. Conf. Big Data*, 2019, pp. 1117–1124.

[154] C. Qing, B. Cai, Q. Yang, J. Wang, and C. Huang, "Deep learning for CSI feedback based on superimposed coding," *IEEE Access*, vol. 7, pp. 93 723–93 733, 2019.

[155] D. Xu, Y. Huang, and L. Yang, "Feedback of downlink channel state information based on superimposed coding," *IEEE Commun. Lett.*, vol. 11, no. 3, pp. 240–242, March 2007.

[156] C. Qing, Q. Ye, W. Liu, and J. Wang, "Fusion learning for 1-bit CS-based superimposed CSI feedback with bi-directional channel reciprocity," *IEEE Commun. Lett.*, vol. 26, no. 4, pp. 813–817, Apr. 2022.

[157] Q. Liu, J. Guo, C.-K. Wen, and S. Jin, "Adversarial attack on DL-based massive MIMO CSI feedback," *J. Commun. Netw.*, vol. 22, no. 3, pp. 230–235, June 2020.

[158] P. K. Sangdeh, H. Pirayesh, A. Mobiny, and H. Zeng, "LB-SciFi: Online learning-based channel feedback for MU-MIMO in wireless LANs," in *Proc. IEEE 28th Int. Conf. Netw. Protocols (ICNP)*, 2020, pp. 1–11.

[159] X. Yu, D. Li, Y. Xu, and Y.-C. Liang, "Convolutional autoencoder-based phase shift feedback compression for intelligent reflecting surface-assisted wireless systems," *IEEE Commun. Lett.*, vol. 26, no. 1, pp. 89–93, Jan. 2022.

[160] X. Li, J. Guo, C.-K. Wen, S. Jin, and S. Han, "Multi-task learning-based CSI feedback design in multiple scenarios," *arXiv preprint arXiv:2204.12698*, 2022. [Online]. Available: https://arxiv.org/abs/2204.12698

[161] J. Yang, X. Chen, H. Zou, D. Wang, Q. Xu, and L. Xie, "EfficientFi: Towards large-scale lightweight WiFi sensing via CSI compression," *IEEE Internet Things J.*, 2021, Early access.

[162] L. Liu, C. Oestges, J. Poutanen *et al.*, "The COST 2100 MIMO channel model," *IEEE Wireless Commun.*, vol. 19, no. 6, pp. 92–99, Dec. 2012.

[163] W. Utschick, V. Rizzello, M. Joham, Z. Ma, and L. Piazzi, "Learning the CSI recovery in FDD systems," *IEEE Trans. Wireless Commun.*, 2022, Early access.

[164] J. Yu, L. Yang, N. Xu, J. Yang, and T. Huang, "Slimmable neural networks," in *Proc. ICLR*, 2019. [Online]. Available: https://openreview.net/forum?id=H1gMCsAqY7

[165] W. Tong and G. Y. Li, "Nine challenges in artificial intelligence and wireless communications for 6G," *IEEE Wireless Commun.*, pp. 1–10, 2022, Early access.

[166] H. Xiao, W. Tian, W. Liu, and J. Shen, "ChannelGAN: Deep learning-based channel modeling and generating," *IEEE Wireless Commun. Lett.*, vol. 11, no. 3, pp. 650–654, March 2022.

[167] T. Campbell and T. Broderick, "Bayesian coreset construction via greedy iterative geodesic ascent," in *Proc. ICML*, 2018, pp. 698–706.

[168] Y. Yang, F. Gao, X. Ma, and S. Zhang, "Deep learning-based channel estimation for doubly selective fading channels," *IEEE Access*, vol. 7, pp. 36 579–36 589, 2019.

[169] S. Basharat, S. A. Hassan, H. Pervaiz, A. Mahmood, Z. Ding, and M. Gidlund, "Reconfigurable intelligent surfaces: Potentials, applications, and challenges for 6G wireless networks," *IEEE Wireless Commun.*, vol. 28, no. 6, pp. 184–191, Dec. 2021.

[170] E. D. Carvalho, A. Ali, A. Amiri, M. Angjelichinoski, and R. W. Heath, "Non-stationarities in extra-large-scale massive MIMO," *IEEE Wireless Commun.*, vol. 27, no. 4, pp. 74–80, Aug. 2020.

[171] W. Tang, X. Chen, M. Z. Chen *et al.*, "On channel reciprocity in reconfigurable intelligent surface assisted wireless networks," *IEEE Wireless Commun.*, vol. 28, no. 6, pp. 94–101, Dec. 2021.

[172] OPPO, "Wireless-Intelligence," https://wireless-intelligence.com/, accessed on May 1, 2022.