

Deep Learning for Massive MIMO Channel State Acquisition and Feedback

Mahdi Boloursaz Mashhadi, and Deniz Gündüz

Dept. of Electrical and Electronic Eng., Imperial College London, UK

Email: {m.boloursaz-mashhadi, d.gunduz}@imperial.ac.uk

Abstract

Massive multiple-input multiple-output (MIMO) systems are a main enabler of the excessive throughput requirements in 5G and future generation wireless networks as they can serve many users simultaneously with high spectral and energy efficiency. To achieve this, massive MIMO systems require accurate and timely channel state information (CSI), which is acquired by a training process that involves pilot transmission, CSI estimation and feedback. This training process incurs a training overhead, which scales with the number of antennas, users and subcarriers. Reducing this training overhead in massive MIMO systems has been a major topic of research since the emergence of the concept. Recently, deep learning (DL)-based approaches for massive MIMO training have been proposed and showed significant improvements compared to traditional techniques. This paper provides an overview of how neural networks (NNs) can be used in the training process of massive MIMO systems to improve the performance by reducing the CSI acquisition overhead and to reduce complexity.

I. INTRODUCTION

Massive multiple-input multiple-output (MIMO) systems are an important component of 5G and future generation wireless networks due to their ability to serve many users simultaneously with high spectral and energy efficiency. The main idea in massive MIMO is to equip base stations (BSs) in wireless networks with large arrays of cooperating antennas to facilitate spatial multiplexing of many user equipments (UEs) within the same time-frequency resources. Since the number of antennas at the BS is typically assumed to be significantly larger than the number of users, a large number of degrees of freedom are available in the downlink, which can be used to shape the transmitted signals in a specific direction or to null interference. This yields a beamforming gain that translates into increased energy efficiency, reduced interference and

improved coverage. In the uplink, each single-antenna user in a massive MIMO system can scale down its transmit power proportional to the number of antennas at the BS while maintaining the same performance as the corresponding single-input single-output (SISO) system. This leads to higher energy efficiency which is a major benefit in next generation wireless networks where excessive energy consumption is a growing concern. On the other hand, if adequate transmit power is available, then a massive MIMO system can significantly increase the range of operation compared with a single antenna system.

In communication systems, channel state information (CSI) is required at the receiver to be able to decode the information transmitted over a time-varying channel. CSI is acquired by a training process which involves pilot transmission and CSI estimation at the receiver. This imposes a training overhead on the communication system which scales up with the number of antennas, receivers and subcarriers. In massive MIMO systems, to achieve the aforementioned performance gains, accurate and timely CSI is required both at the BS and the UEs. Availability of downlink CSI at the massive MIMO BS is crucial to enable beamforming and achieve spatial multiplexing gains. Reducing the training overhead in massive MIMO has been a major topic of research since the emergence of the concept.

Massive MIMO was originally introduced in a time division duplex (TDD) setting where the uplink and downlink channels are separated in time [1], [2]. In the TDD mode of operation, due to uplink/downlink channel reciprocity, which holds under certain conditions [3], downlink CSI does not induce extra training overhead. However, motivated by spectrum regulation issues, FDD operation gained significant interest [4], [5], and there has been a long-standing debate on the relative performance of TDD and FDD schemes [6]–[8]. Although the FDD scheme is favourable due to its improved coverage and reduced interference, these benefits come at the price of increased complexity of the training process for FDD massive MIMO. Unlike in TDD, in FDD the uplink and downlink channels are separated in frequency, and hence are not reciprocal. Consequently, in FDD massive MIMO, downlink CSI need to be first estimated at each UE, and then fed back to the BS through the uplink channel, which significantly increases the CSI overhead. Fig. 1 depicts the downlink training process in FDD mode.

For smaller number of BS antennas, simple vector quantization (VQ) along with exhaustive search may work sufficiently well for MIMO CSI. In the fourth generation long term evolution (4G-LTE) advanced standard, a 4-bit channel quality index (CQI) and the pre-coding matrix indicator (PMI) are fed back to the BS to reveal the CSI [9]. However, with the increased

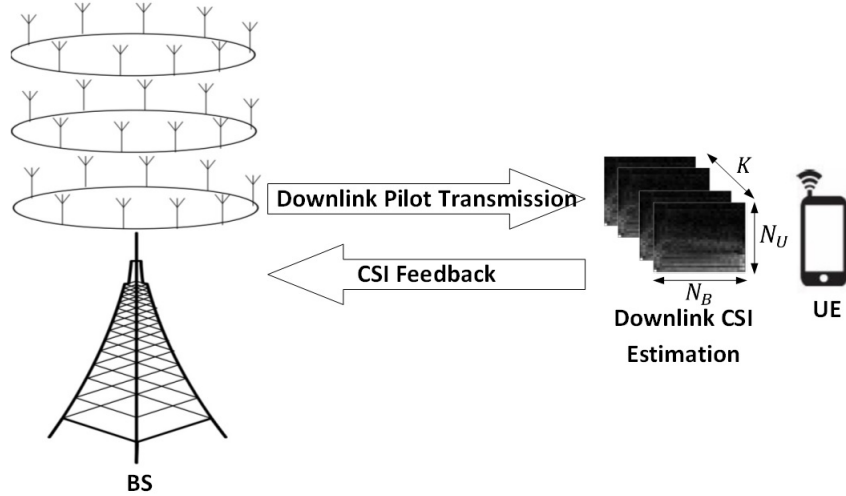


Fig. 1: Massive MIMO downlink training process in FDD mode.

number of massive MIMO antennas, CSI dimensions increase drastically and the traditional VQ-based approaches are no longer practical. This has encouraged great interest in more efficient training and compression techniques for massive MIMO CSI acquisition and feedback. Initial efforts in this direction, took a model-based approach assuming sparse or low-rank models on the CSI matrix. However, a sparse model on the channel may not be very accurate when MIMO dimensions are not sufficiently large, which degrades the performance of sparsity-based techniques. The same discussion holds for low-rank based techniques [10], [11] where there is a model mismatch. These approaches do not take into account the inherent statistical correlations and structures beyond sparse or low-rank patterns. Moreover, sparse and low-rank reconstruction techniques are computationally demanding iterative algorithms, which may further limit their practical implementation.

More recently, data-driven approaches have been proposed as an alternative paradigm based on the common structures observed in typical massive MIMO CSI matrices. Data-driven approaches train neural network (NN) structures over large datasets of CSI matrices to capture these structures and use them to reduce the CSI acquisition overhead. NN-based approaches have recently shown significant improvements over their model-based counterparts in wireless communications [12]–[14].

Consider the following massive MIMO channel matrix $\mathbf{H}(\tau) \in \mathbb{C}^{N_U \times N_B}$ in the delay domain:

$$\mathbf{H}(\tau) = \sqrt{\frac{N_U N_B}{L}} \sum_{l=1}^L \alpha_l \delta(\tau - \tau_l) \mathbf{a}_U(\theta_l) \mathbf{a}_B^H(\phi_l), \quad (1)$$

where N_U and N_B denote the number of antennas at the BS and UE, respectively, L is the number of multi-path components with α_l denoting the propagation gain of the l th path. The propagation gain α_l is assumed to follow a zero-mean complex Gaussian distribution $\mathcal{CN}(0, \sigma_\alpha^2)$, where σ_α^2 denotes the average power gain. Also, \mathbf{a}_B and \mathbf{a}_U are the array response vectors for the BS and user with θ_l and ϕ_l denoting the azimuth angles of arrival and/or departure (AoA/AoD) respectively, and $(\cdot)^H$ denotes the conjugate transpose operation. For uniform linear arrays, we have

$$\begin{aligned} \mathbf{a}_U(\theta_l) &= [1, e^{-j\frac{2\pi d}{\lambda} \sin \theta_l}, \dots, e^{-j\frac{2\pi d}{\lambda} (N_U-1) \sin \theta_l}]^T, \\ \mathbf{a}_B(\phi_l) &= [1, e^{-j\frac{2\pi d}{\lambda} \sin \phi_l}, \dots, e^{-j\frac{2\pi d}{\lambda} (N_B-1) \sin \phi_l}]^T, \end{aligned} \quad (2)$$

where d and λ denote the distance between adjacent antennas and the carrier wavelength, respectively. Equivalently, the MIMO channel matrix at the k th subcarrier in OFDM $\mathbf{H}_k \in \mathbb{C}^{N_U \times N_B}$ is given by (3)

$$\mathbf{H}_k = \sqrt{\frac{N_U N_B}{L}} \sum_{l=1}^L \alpha_l e^{-j2\pi\tau_l f_s \frac{k}{K}} \mathbf{a}_U(\theta_l) \mathbf{a}_B^H(\phi_l), \quad (3)$$

where f_s denotes the sample rate and K is the total number of subcarriers.

According to (3), the CSI values for nearby users, sub-carriers and antennas are correlated due to similar propagation paths, gains, delays and AoDs/AoAs. Apart from the correlations governed by (3), there exists inherent characteristics in MIMO environments due to specific user distributions, scattering parameters, geometry, materials, etc., that cause common structures among MIMO CSI matrices. We note that the joint statistics of the channel gains across antennas, subcarriers and users is extremely complicated. Even if accurate models are known on the statistics in (3), identifying a lossy compression scheme to optimally exploit structures and correlations in (3) is challenging. On the other hand, NNs are extremely powerful in learning complex distributions and exploiting them for various classification/regression (supervised learning) or compression (unsupervised learning) tasks. NNs can be used to learn the common structures and inherent correlations to leverage them for efficient CSI estimation, compression and feedback, reducing the overall MIMO training overhead.

Success of data-driven approaches depends critically on the datasets used to train the NN models. Unlike some more popular applications of NNs, rich and standardized datasets of CSI measurements in actual massive MIMO scenarios do not yet exist. However, there exists MIMO channel models that have proved to be very accurate in statistically modeling actual CSI measurements in practical MIMO scenarios. Among these are the third generation partnership project (3GPP) spatial channel model (SCM) [15], WINNER II [16] and COST 2100 [17]. Unfortunately, different researches have used different channel models to generate CSI datasets ranging from the simple formula in (3) to more sophisticated channel models like COST 2100, 3GPP TR 38.901 release 15 [18] or the DeepMIMO ray-tracing propagation model [19], which makes comparison between results difficult. The most widely used channel model so far has been COST2100 which will also be used in this paper. We would like to emphasize that different datasets hamper the comparison of different results and there is a pressing need for standard datasets.

This paper provides an overview of how NNs can be used in massive MIMO systems to improve the performance of CSI acquisition and feedback while reducing both the complexity and overhead. In the following sections, we shall review recently proposed data-driven approaches for CSI estimation, compression and feedback and suggest future research directions.

II. MIMO CHANNEL ESTIMATION BY DL

Consider uplink MIMO training where the user transmits a block of P pilot signals, denoted by $\mathbf{X} \in \mathbb{C}^{N_U \times P}$, which is known at both the UE and the BS. The BS needs to estimate the channel matrix $\mathbf{H} \in \mathbb{C}^{N_B \times N_U}$ from received measurements $\mathbf{Y} \in \mathbb{C}^{N_B \times P}$, given by

$$\mathbf{Y} = \mathbf{H}\mathbf{X} + \mathbf{Z}, \quad (4)$$

where $\mathbf{Z} \in \mathbb{C}^{N_B \times P}$ is the complex additive white Gaussian noise (AWGN).

Standard channel estimation techniques are typically based on linear minimum mean square error (LMMSE) estimation method. A common assumption in LMMSE-based channel estimation techniques is that the pilot length is larger than the number of transmit antennas, which may be prohibitive in downlink training of massive MIMO systems ($P \geq N_B$). For downlink massive MIMO channel estimation, where N_B is large, it is challenging to ensure $P \geq N_B$ not only because it shall increase the training overhead and computational complexity for channel estimation, but also because a large P may even exceed the channel coherence interval. If this assumption does not hold, LMMSE-based channel estimation performance degrades significantly.

Many previous works take a model-based estimation approach assuming sparse [20]–[22] or low-rank [10], [11] models on the channel matrix. Sparsity of the channel in the angular-delay domain has been assumed in [20]–[22] where compressive sensing based reconstruction techniques are used to reduce the pilot length and training overhead. Sparsity based techniques can decrease the pilot length required to sense and estimate the channel by an order of magnitude [23] compared to a simple exhaustive search approach. However as mentioned earlier, these techniques rely on the sparse or low-rank properties of the channels, which may not be very accurate and do not take into account the inherent statistical correlations and structures beyond sparse or low-rank patterns.

This motivates the use of data-driven approaches based on NNs to learn these complex structures and correlations [24]. The authors in [24], [25] use convolutional NNs to improve the quality of a coarse initial estimate of the channel matrix exploiting temporal and inter-frequency correlations. Denoting the MIMO channel matrix for the k th subcarrier at temporal slot n by $\mathbf{H}_k(n) \in \mathbb{C}^{N_B \times N_U}$ where the channel is assumed constant during each slot. A coarse initial estimate of $\mathbf{H}_k(n)$ is given by $\mathbf{R}_k(n) = \mathbf{X}_k(n)^\dagger \mathbf{H}_k(n)$, where $\mathbf{X}_k(n)^\dagger$ denotes the pseudo-inverse of the pilot signals transmitted over the k th subcarrier at time n . The authors form large tensors by concatenating $\mathbf{R}_k(n)$ s along time and frequency dimensions, and then apply multi-dimensional convolution kernels on it. During training, these kernels capture temporal and inter-frequency correlations, and can be exploited to provide accurate estimates of the channel matrix. This idea outperforms non-ideal minimum mean square error (MMSE) (with estimated covariance matrix) estimation and achieves performance very close to the ideal MMSE (with true covariance matrix) that is very difficult to be implemented in practical situations.

The NN architecture used in [24] consists of 12 convolutional layers. There is still much work to be done to design NN architectures with reduced complexity and improved performance to guarantee that the channel estimation task can be carried out rapidly within the channel coherence time.

On the other hand, many massive MIMO structures use low-resolution analog-to-digital converters (ADCs) to reduce the power consumption and hardware complexity at the BS; and hence, only a coarsely quantized version of \mathbf{Y} shall be available for channel estimation at the BS. For the quantized case, we have (5)

$$\mathbf{Y}_q = \mathcal{Q}(\mathbf{H}\mathbf{X} + \mathbf{Z}), \quad (5)$$

where $\mathcal{Q}(\cdot)$ denotes quantization performed element-wise on the real and imaginary parts of the received signals independently. Low-resolution ADCs incur nonlinear distortion, which poses significant challenges to channel estimation from highly quantized measurements. Hence, efficient estimation techniques from quantized received signals \mathbf{Y} are needed.

With coarsely quantized measurements, the pilot length required for reliable estimation of the channel, and hence, the training overhead increases significantly. Previously proposed model-based estimation techniques generally minimize a cost function (e.g. maximum likelihood, square error, etc.) iteratively subject to sparsity [26]–[28] or low-rank [29] constraints on the channel matrix \mathbf{H} . Due to the additional non-linearity introduced by quantization, NN-based techniques can be even more beneficial in channel estimation from low resolution received signals.

For the extreme case of 1-bit ADCs, reconstruction is possible only up to a scale factor. The initial results reported in [30], [31] show that a simple fully-connected network trained in a supervised setting to estimate the channel directly from sign measurements can reduce the required pilot length roughly by an order of magnitude, while achieving similar reconstruction performance in comparison with previous sparse or low-rank based techniques. In [32], the authors consider a mixed-ADC scenario, where several BS antennas are equipped with high resolution ADCs and others with few-bit ADCs to achieve a trade off between the performance and power consumption. They input an initial least square (LS) channel estimate to a 5-layer fully-connected NN and show that the NN can learn to utilize the correlation between antennas to improve the estimation performance for the low-resolution branches. The above works utilize inter-antenna correlations for channel estimation; however, there is still room for improvement utilizing temporal and inter-frequency correlations by convolutional kernels in a setting similar to [24].

While fully-connected NNs have been commonly used in previous works and they have the potential to learn and exploit complex joint distributions across all antennas and subcarriers, they do not easily scale with MIMO dimensions and need separate training for different number of antennas, subcarriers, etc. However, as discussed earlier, correlations in typical MIMO channels exhibit locality among antennas and subcarriers, which encourages utilizing convolutional architectures, which can significantly reduce the complexity in both training and inference. This is especially critical in wireless applications, as it is important to acquire an accurate channel estimate within the channel coherence time. Moreover, convolutional kernels, once trained, work for different input dimensions; that is, we do not need to train and use a different NN when

the number of antennas in either side of the channel, or the number of subcarriers allocated for communication changes.

III. DL-BASED MIMO CSI REDUCTION AND FEEDBACK

Once the channel matrix \mathbf{H} is estimated at the UE, it needs to be transmitted back to the other side through a feedback channel, which incurs further overhead. With massive number of antennas and increased number of subcarriers, the CSI matrix dimensions, and the resulting overhead, increase significantly, which motivate CSI reduction techniques. Many works in the literature focused on more traditional compression techniques, such as VQ, sparsifying transforms (e.g., discrete cosine transform (DCT), Karhunen-Loeve transform (KLT)), principal component analysis (PCA)-based dimensionality reduction and compressed sensing (CS) to compress the CSI using spatio-temporal MIMO channel correlations.

However, as we have mentioned earlier, lossy compression is a challenging task even when the underlying source distributions is known perfectly. While we have a relatively good understanding of the fundamental rate-distortion performance for independent and identically distributed source in the asymptotic limit, lossy compression for practical sources, such as image, audio, or video, has been a research channel for many decades. This is where dimensionality reducing autoencoders can be borrowed from ML theory and put into practice to efficiently reduce the massive MIMO CSI overhead. These autoencoder architectures can be trained to learn a lower dimensional representation of the original CSI matrix to be transmitted over the feedback channel with a reduced overhead. An initial study using this autoencoder approach showed significant improvement in comparison with the best performing sparsity-based techniques [33]. The authors in [33] proposed CSINet, which has since been adopted as a benchmark architecture for performance comparisons by subsequent researches. CSINet includes convolutional layers as well as dense layers and Refine-Net architectures. In [34], [35], the authors combine CSINet and long short-term memory (LSTM) cells to improve upon the basic CSINet architecture by exploiting the temporal correlations in CSI matrices for consecutive time instances. The authors in [36] use the uplink CSI (which is already available to the BS by uplink training) as a side information to further improve CSI reconstruction performance utilizing the correlations between downlink and uplink channels.

There are two main approaches to cope with the limitations in the CSI feedback channel, i.e., the *digital* and *analog* CSI schemes. Digital schemes, which have traditionally received

more attention, are based on the separation approach: CSI is first compressed into as few bits as possible and these bits are reliably fed back to the transmitter using a low-rate channel code, which adds redundancy in a way to cope with the channel noise and error in the feedback link. On the other hand, analog CSI follows a joint source-channel coding approach, and directly maps the downlink CSI to the uplink channel input in an unquantized and uncoded manner. The analog scheme simplifies the feedback operation as it does not require explicit quantization, coding, and modulation. If the uplink feedback channel is an additive white Gaussian noise channel, and the downlink CSI is Gaussian and perfectly known at the UE, the analog CSI scheme (that incurs zero delay) is optimal in that it achieves the same minimum mean-squared error distortion for the reconstructed CSI at the BS as a scheme that optimally quantizes and encodes the CSI, while incurring infinite delay. The low-latency of the analog CSI scheme makes it a favourable alternative in rapidly changing MIMO channels where the CSI needs to be estimated and fed back to the BS periodically.

Many DL-based CSI reduction techniques [33], [34], [36]–[38] train an end-to-end auto-encoder assuming ideal feedback of the reduced CSI. However, the estimated CSI (forward channel state) is fed back to the transmitter through the uplink channel which suffers from noise, interference and fading. It becomes crucial to design CSI compression and feedback schemes that not only reduce the CSI overhead efficiently, but are also robust against the feedback channel impairments.

The authors in [39] train the autoencoder over an AWGN feedback channel. However, a simple AWGN model is far from accurate for the feedback channel, which is another MIMO fading channel. In [40], the COST2100 model is used to generate simultaneous uplink and downlink channel matrices, which enables taking feedback link impairments into account explicitly. By training over datasets of both uplink and downlink channels, the end-to-end CSI reconstruction quality improves as the NN shall not only capture the structures in the DL channel, but also in the uplink feedback channel and the uplink/downlink inter-channel correlations.

We shall consider both analog and digital CSI schemes in presence of feedback channel impairments in the following subsections.

A. Digital CSI feedback

The earlier autoencoder-based CSI reduction techniques [33], [34], [36] overlooked subsequent feedback of the reduced CSI and mainly focused on the dimensionality reduction by a direct

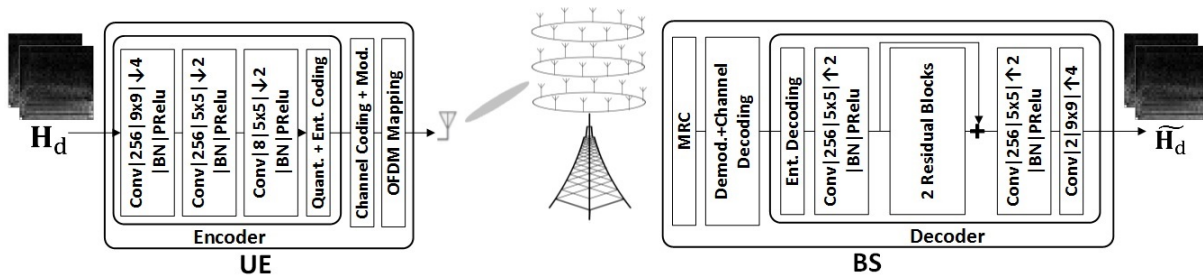


Fig. 2: DeepCMC for MIMO CSI compression.

application of the autoencoder architecture. These works seem to be based on the assumption that reducing the dimension of CSI matrix would result in reduced feedback overhead. This is not necessarily correct since the reduced representation consists of real numbers, which may still need to be compressed further, and the impact of such compression on the final CSI accuracy is not taken into account. Several subsequent works assume that the reduced CSI is quantized before being digitally fed back to the BS. The authors consider simple uniform quantization in [41] and non-uniform μ -law quantization in [42]. Since quantization is a non-differentiable function, the gradient cannot pass through it in the backpropagation step of the learning algorithm. This makes it challenging to train digital CSI feedback schemes in an end-to-end manner and requires further considerations to overcome the gradient backpropagation issue. A widely used solution is to set the quantization gradient to a constant and train end-to-end for a specific number of quantization bits. The authors in [42] add an offset module to the decoder to compensate for the quantization distortion where the network is trained in multiple stages: end to end training without quantization with a larger learning rate, adding quantization and optimizing the offset module, and finally fine tuning the offset and decoder by further training with a small learning rate.

Although the authors in [41], [42] consider quantization of the reduced CSI to convert it into bits to be transmitted over the feedback link, the simple scalar quantization approach cannot fully exploit the potential correlations remaining among the components of the reduced CSI. Indeed the quantizer output does not produce equally probable bits; and hence, additional lossless compression of the bits would further reduce the CSI and reduce the feedback overhead.

In [43], we employ entropy coding to further compress the quantizer outputs. Fig. 2 provides the end-to-end block diagram for a downlink digital CSI feedback scheme employing our

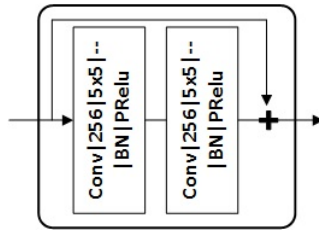


Fig. 3: The residual block model.

proposed deep learning based CSI compression technique, called DeepCMC [43]. In this figure, \mathbf{H}_d and $\widetilde{\mathbf{H}}_d$ denote the downlink CSI matrix at the UE and its estimate at the BS, respectively, and the two model input matrices represent $\Re(\mathbf{H}_d)$ and $\Im(\mathbf{H}_d)$. According to this figure, the UE applies a CNN-based feature encoder on \mathbf{H}_d to obtain its low dimensional representation, which is subsequently quantized and compressed using context-adaptive binary arithmetic coding (CABAC) [44]. The resulting bit stream passes through channel coding and digital modulation. The modulation output is then mapped over OFDM subcarriers and transmitted back to the BS over the uplink channel. The BS performs maximum ratio combining (MRC) on the received signals to maximize the SNR and benefit diversity of the feedback channel. The resulting signal then passes through the demodulation, channel decoding, entropy decoding and CNN-based feature decoder to reconstruct $\widetilde{\mathbf{H}}_d$.

In the CNN architecture in Fig. 2, “Conv|256| 9×9 | $\downarrow 4$ |BN|PReLU” represents a convolutional layer with 256 features and kernel size of 9×9 followed by downsampling by a factor of 4, batch normalization and parametric rectified linear activation unit (PReLU). As depicted in Fig. 2, the feature decoder consists of three convolutional layers and two residual blocks with shortcut connections where “+” denotes simple element-wise addition. Fig. 3 illustrates the architecture for each residual block where “|–” means the corresponding convolution output is not downsampled. The residual and shortcut structures ease training of the network by preventing vanishing gradients along the stacked non-linear layers and improve the performance according to our simulation results.

The training cost for DeepCMC is a weighted sum of the mean square error (MSE) of the CSI reconstruction and the quantizer’s output entropy. A weight parameter λ controls the tradeoff between the reconstruction quality and the feedback bit rate, with a larger value resulting in improved MSE at an increased bit rate. For a good quality feedback channel with larger capacity,

utilizing a network trained with a larger λ results in improved CSI quality at the BS. However, if the feedback channel capacity is smaller than the resulting bit rate, the feedback channel will fail to deliver the CSI. To avoid this, a network trained to work at a lower bit rate (trained with smaller λ) should be used. Different λ values will provide networks that work on different points on a rate-distortion curve. The UE will store different networks, and use the proper one depending on the uplink channel state and the capacity achievable for CSI feedback.

We note that in contrast to the literature on CSI feedback, which has mainly focused on minimizing the reconstruction error, DeepCMC is trained with a rate-distortion cost that takes into account both the compression rate (in terms of bits per channel dimension) and the reconstruction MSE. As we will see below, this additional compression step leads to a significant improvement in the achieved performance. It also allows adapting the CSI quality to the available feedback channel quality.

Another important benefit of the DeepCMC architecture is that it is fully convolutional, and has no densely connected layers, which makes it flexible for a wide range of MIMO scenarios with different number of sub-channels and antennas. As shown by the simulation results, although DeepCMC is trained for a specific number of sub-channels and antennas it generalizes well to different number of sub-channels and antennas [43]. This is very important for practical implementation of NN-based CSI compression techniques, as otherwise the nodes would have to store a large number of NN parameters for every possible combination of antenna and subcarrier numbers.

In Fig. 4 we compare the output rate-distortion curves for DeepCMC [43], CSINet [33] and CRNet [38]. In this comparison, we use the normalized mean square error defined as $\text{NMSE} \triangleq \mathbb{E} \left[\frac{\|\mathbf{H}_d - \widetilde{\mathbf{H}}_d\|_2^2}{\|\mathbf{H}_d\|_2^2} \right]$. We plot the achieved NMSE, in dBs, as a function of the average number of bits used to encode each CSI entry. Note that the outputs for CSINet and CRNet are feature vectors of type “float32”, and hence 32-bit quantization is considered to calculate the resulting bit rate for them.

For the comparison in Fig. 4 we consider downlink training for a single-antenna user in an FDD MIMO setting. We set $K = 256$, $N_B = 32$, $N_U = 1$, and use the COST2100 channel model [17] to generate sample channel matrices for training and testing. We consider the indoor picocellular scenario at 5.3 GHz, where the BS is equipped with a ULA of dipole antennas positioned at the center of a $20\text{m} \times 20\text{m}$ square. The user is placed within this square uniformly at random. All other parameters follow the default settings in [17]. The number of training and

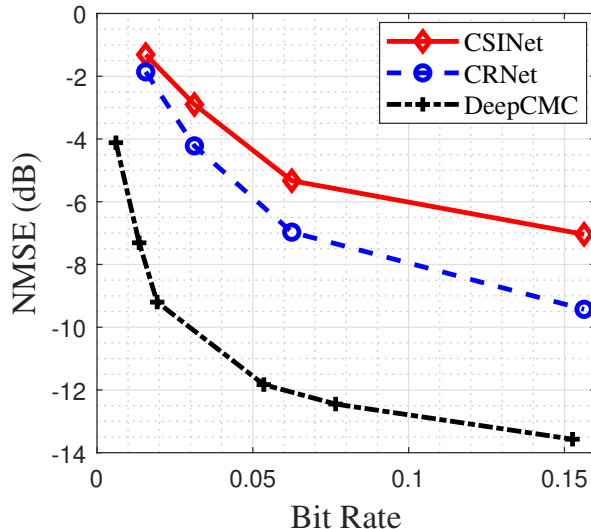


Fig. 4: Bit rate-NMSE trade-off comparison, $K = 256$, $N_B = 32$, $N_U = 1$.

testing samples are 80000 and 20000, respectively, and the batch size is 100.

As it can be observed from Fig. 4, DeepCMC provides significant improvement in the quality of the reconstructed CSI at the BS with respect to CSINet and CRNet at all bit rate values. We remark here that, CSINet itself provides 3 – 6dB improvement in NMSE compared to model-based CSI compression techniques in the literature exploiting sparsity of the channel gain matrix [33]. However, the gains from DeepCMC are even more drastic, achieving remarkably good reconstruction of the channel gain matrix with NMSE of -13 dB at a bit rate lower than 0.16 bits per CSI entry. These results show that DeepCMC outperforms CSINet 4 to 6 dB in NMSE for the range of compression rates considered here. For example, for a target value of $\text{NMSE} = -5$ dB, DeepCMC can provide more than 5 times reduction in the number of bits that must be fed back from the UE to the BS.

B. Analog CSI feedback

The analog CSI feedback scheme follows a joint source-channel coding approach, and directly maps the downlink CSI to the uplink channel input in an unquantized and uncoded manner. A CNN-based analog CSI feedback scheme, namely AnalogDeepCMC is proposed in [40], to do the CSI compression and feedback tasks simultaneously taking into account the feedback channel impairments. It uses a fully convolutional autoencoder model to efficiently map the downlink

CSI at the UE to the uplink channel inputs, and to reconstruct them at the BS. The model is trained treating the uplink feedback channel as a non-trainable layer in the autoencoder. In this section, we provide performance comparisons between AnalogDeepCMC and the digital approach using DeepCMC for CSI compression, based on the quality of the reconstructed CSI at the BS when the same amount of uplink channel resources is devoted to CSI feedback. We will observe that the analog scheme improves the CSI reconstruction quality and consequently the achievable downlink rate without requiring the UL CSI at the UE for feedback transmission.

Consider CSI feedback from a single-antenna user to a BS with N_B antennas utilizing OFDM over K subcarriers. Denote the uplink and downlink channel matrices by $\mathbf{H}_u \in \mathbb{C}^{K \times N_B}$ and $\mathbf{H}_d \in \mathbb{C}^{K \times N_B}$, respectively. Assume that the downlink CSI \mathbf{H}_d available at the UE is fed back to the BS over N_F uplink OFDM subcarriers devoted to CSI feedback picked uniformly at random, with $\rho \triangleq \frac{N_F}{K}$ denoted as the *feedback overhead*. The feedback channel over the j -th uplink subcarrier denoted by $\mathbf{h}_F^j \in \mathbb{C}^{N_B \times 1}$, $j = 1, \dots, N_F$, is obtained from the corresponding row of \mathbf{H}_u , which specifies a SIMO channel with its output given by

$$\mathbf{y}_j = \mathbf{h}_F^j x_j + \mathbf{z}_j, \quad (6)$$

in which $\mathbf{y}_j \in \mathbb{C}^{N_B \times 1}$ is the received signal at the BS antennas, x_j is the symbol fed back over the j -th subcarrier and $\mathbf{z}_j \in \mathbb{C}^{N_B \times 1}$ is the independent AWGN component. With N_F uplink subcarriers dedicated for CSI feedback, a maximum rate of $C_{FB} = \sum_{j=1}^{N_F} \log_2(1 + SNR_{FB} \|\mathbf{h}_F^j\|^2)$ is available for CSI feedback, where SNR_{FB} is the signal to noise ratio (SNR) in the uplink channel. However, note that C_{FB} depends on the uplink channel state which is not known by the UE. In a digital CSI feedback scheme, the UE will typically take a conservative approach and transmit at a rate that can be decoded with high probability. Here, we use C_{FB} as the feedback rate to obtain an upper bound on the performance of any digital CSI feedback scheme.

Fig. 5 depicts the architecture of our proposed analog CSI feedback model, Analog-DeepCMC. The UE applies a CNN-based feature encoder composed of three convolutional layers, which outputs real-valued features. Each pair of these real numbers are then grouped to form a complex-valued symbol, which are subsequently normalized to ensure the input power constraint over the feedback channel is met. These normalized symbols are then directly mapped into the corresponding subcarriers, and transmitted over the CSI feedback channel. The BS then performs maximum ratio combining and feature decoding to reconstruct the original CSI matrix \mathbf{H}_d .

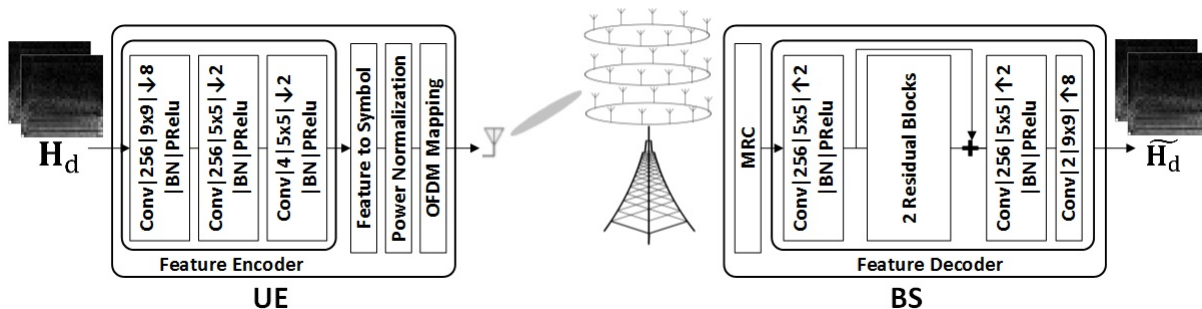


Fig. 5: The AnalogDeepCMC architecture for MIMO analog CSI feedback.

AnalogDeepCMC is trained including the feedback channel (noise and fading) and MRC blocks as non-trainable layers in between the autoencoder structure.

As we have discussed above, in the case of digital CSI feedback, the UE has to know the uplink channel capacity to choose the highest compression rate that can be reliably fed back to the BS. However, during downlink channel training in an FDD MIMO scenario, the UE does not yet know the uplink CSI, and hence, will typically take a cautious channel coding and modulation approach which works at a rate significantly below C_{FB} . Moreover, we assume error-free transmission at the capacity despite a codelength of only N_F symbols. In our simulations we will assume that UE has knowledge of the instantaneous uplink channel capacity and can transmit at this rate reliably. Hence, the corresponding NMSE results achieved in simulations serve as a rather generous lower bound on the actual NMSE performance of any practical digital CSI feedback scheme.

Fig. 6 depicts the average NMSE (dB) as a function of the CSI overhead ρ for different values of the uplink SNR using the digital scheme based on separate CSI compression using the DeepCMC algorithm followed by capacity-achieving channel coding. In this figure, DeepCMC is trained for two different λ values resulting in NN1 and NN2. NN1 corresponds to a point with better reconstruction quality at a higher rate on the rate-distortion curve given in Fig. 4. The simulation scenario is the same as in Fig. 4 with $K = 256, N_B = 32, N_U = 1$. Note that DeepCMC is a variable-length lossy compression scheme; that is, for each CSI matrix, the UE obtains different number of bits at the output of the entropy coder. On the other hand, the capacity of the feedback channel is also random, depending on the states of the N_F subcarriers dedicated to CSI feedback. Therefore, if the number of bits at the encoder's output exceeds C_{FB} , the feedback channel fails to deliver the CSI, called an *outage* event, and the NMSE

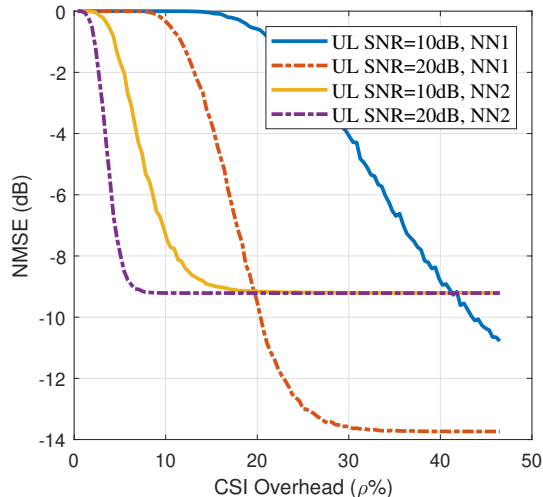


Fig. 6: Average NMSE vs. ρ for digital CSI feedback using DeepCMC followed by capacity-achieving channel code.

will equal 0dB. This is why the NMSE curves all saturate at 0dB for low ρ values. If the CSI overhead decreases below a threshold, outages will occur with increasing frequency resulting in an increased NMSE. As ρ increases beyond this threshold value, outage probability decreases with ρ . Beyond another higher threshold value, outage probability approaches zero, and the autoencoder reconstructs the CSI at the NMSE that it has been trained for (depending on the λ value which controls the rate-distortion trade-off). This is the reason why the NMSE curves also saturate at high ρ values. According to the figure, as the uplink SNR decreases, thresholds for both saturation regions increase. We would like to highlight that, for the setting considered here ($K = 256$), $\rho = 20\%$ would correspond to a channel code of length 51 symbols, in which case the code rates with reasonable reliability are significantly below the capacity [45]; that is, the NMSE values in this figure are quite generous for the digital scheme.

As observed in Fig. 6, for efficient digital CSI feedback, the UE requires the uplink CSI not only to decide on the appropriate channel coding rate, but also to use a NN trained with the proper λ value to achieve the minimum possible NMSE. Networks trained for different reconstruction qualities result in different threshold behaviours. A network trained for better reconstruction quality results in an increased performance threshold but achieves a smaller NMSE for overhead values above the threshold. If uplink CSI is not available, which is the case during downlink training of FDD massive MIMO, the UE will typically need to take conservative

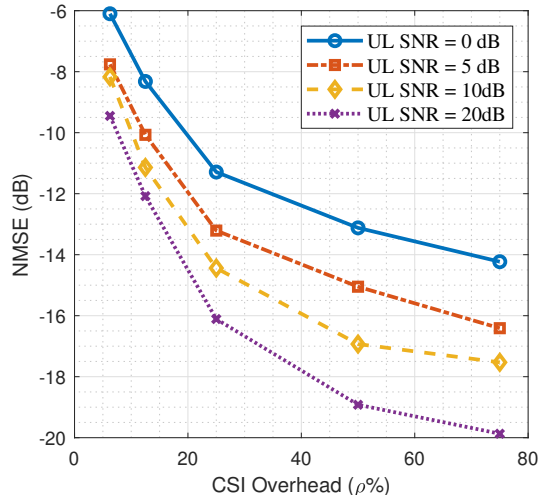


Fig. 7: Average NMSE vs. ρ for the AnalogDeepCMC architecture.

source and channel coding approaches, which will result in considerable degradation of the CSI reconstruction quality.

Fig. 7 depicts the average CSI reconstruction NMSE (dB) as a function of the CSI overhead ρ for different values of the uplink SNR using AnalogDeepCMC. The curves for different SNR values in Fig. 7 correspond to NN models trained for the corresponding uplink SNR. As observed in Fig. 7, there is no threshold behaviour in the analog CSI scheme and the NMSE curves exhibit graceful performance degradation with decreasing SNR in the uplink channel. This is unlike the digital case which may result in severely degraded CSI quality (NMSE= 0dB) due to outages if the uplink SNR decreases below a threshold. Hence, unlike the digital CSI scheme, AnalogDeepCMC does not require uplink CSI to send the downlink CSI back to the BS. The analog CSI scheme is much more favourable not only due to avoiding the performance thresholds and eliminating the need for explicit uplink CSI, but also for avoiding the channel coding and modulation delays.

IV. CSI TRAINING WITH SIDE INFORMATION

In the previous sections, we have focused on exploiting the joint distribution of CSI matrices to reduce the overhead for CSI estimation and feedback using NNs for lossy CSI compression. In this section, we will explore how we can exploit joint distribution across antennas and subcarriers, or across time and space to further reduce, or even completely remove the amount of required

CSI feedback. This is based on the idea of using the available CSI information at the BS at a certain point in time, space or frequency or a subset of antennas as correlated side information to improve the compression efficiency with NNs. If the side information proves to be sufficient for predicting the required CSI with an acceptable distortion using a NN, then the NN can characterize a mapping function to predict the required CSI from the available side information with zero overhead ($\rho = 0$). This can be considered as the Wyner-Ziv lossy compression [46], in the presence of correlated side information at the receiver.

As an example, consider a FDD massive MIMO scenario, where channel reciprocity does not hold and separate downlink and uplink training would normally be necessary. Although the uplink and downlink channels are not fully reciprocal, the uplink and downlink signals traverse the same geometrical paths with different frequencies, which imposes a correlation between the uplink and downlink CSI. The authors in [36] use the uplink CSI (which is already available at the BS by uplink training) as a side information to further improve CSI reconstruction performance utilizing the correlations between downlink and uplink channels. Exploiting the uplink/downlink correlation, the authors in [47] train a fully connected NN to predict downlink from the uplink CSI, and hence totally eliminate the downlink training and feedback overhead. The authors in [41] use CSI which has been delayed due to the limited communication rate in the feedback channel as the correlated side information.

As another example, in TDD massive MIMO, the BS can exploit joint distribution of the CSI for nearby UEs to estimate, compress and feedback the CSI jointly at a reduced overhead. If a NN can be trained to characterize a mapping to predict CSI at a certain set of UEs totally from another set of UEs, it will significantly reduce pilot transmission and CSI feedback overhead [48]. Similarly, in the case of FDD massive MIMO, the BS can use a NN to learn and exploit joint distribution of the CSI for nearby UEs and the correlation among their channels to reduce the CSI feedback overhead similar to a distributed lossy compression scheme.

V. CONCLUSION

Massive MIMO systems are considered as the key technology to enable the excessive throughput requirements in 5G and future generation wireless networks due to their ability to serve many users simultaneously with high spectral and energy efficiency. However, due to the drastic increase in the number of antennas, CSI acquisition and feedback become challenging requiring excessive time, frequency and computational resources potentially crippling benefits of massive

MIMO systems. Many previous works have taken model-driven approaches assuming sparse or low-rank models on the CSI matrix to reduce the overhead. However, these techniques cannot exploit statistical structures that go beyond sparsity. This encouraged data-driven approaches based on training NN architectures over large datasets of CSI matrices, generated using accurate channel models or even from channel measurements, to capture these structures and use them to reduce CSI acquisition and feedback overhead. Deep learning based approaches have shown significant improvements in comparison with traditional CSI acquisition methods. In this paper we have provided an overview of recent results using NNs in the training process of massive MIMO systems that show significant performance improvements while reducing the complexity and the overhead of CSI acquisition. Several interesting future directions are also indicated in this highly promising research area.

REFERENCES

- [1] T. L. Marzetta, "Noncooperative cellular wireless with unlimited numbers of base station antennas," *IEEE Transactions on Wireless Communications*, vol. 9, no. 11, pp. 3590–3600, November 2010.
- [2] F. Rusek, D. Persson, B. K. Lau, E. G. Larsson, T. L. Marzetta, O. Edfors, and F. Tufvesson, "Scaling up MIMO: Opportunities and challenges with very large arrays," *IEEE Signal Processing Magazine*, vol. 30, no. 1, pp. 40–60, Jan 2013.
- [3] T. L. Marzetta, "How much training is required for multiuser MIMO?" in *2006 Fortieth Asilomar Conference on Signals, Systems and Computers*, Oct 2006, pp. 359–363.
- [4] A. Adhikary, J. Nam, J. Ahn, and G. Caire, "Joint spatial division and multiplexing—the large-scale array regime," *IEEE Transactions on Information Theory*, vol. 59, no. 10, pp. 6441–6463, Oct 2013.
- [5] Z. Jiang, A. F. Molisch, G. Caire, and Z. Niu, "Achievable rates of FDD massive MIMO systems with spatial channel correlation," *IEEE Transactions on Wireless Communications*, vol. 14, no. 5, pp. 2868–2882, May 2015.
- [6] J. Flordelis, F. Rusek, F. Tufvesson, E. G. Larsson, and O. Edfors, "Massive MIMO performance—TDD versus FDD: What do measurements say?" *IEEE Transactions on Wireless Communications*, vol. 17, no. 4, pp. 2247–2261, April 2018.
- [7] E. Björnson, E. G. Larsson, and T. L. Marzetta, "Massive MIMO: Ten myths and one critical question," *IEEE Communications Magazine*, vol. 54, no. 2, pp. 114–123, February 2016.
- [8] L. Lu, G. Y. Li, A. L. Swindlehurst, A. Ashikhmin, and R. Zhang, "An overview of massive MIMO: Benefits and challenges," *IEEE Journal of Selected Topics in Signal Processing*, vol. 8, no. 5, pp. 742–758, Oct 2014.
- [9] S. Ahmadi, "LTE-Advanced: A practical systems approach to understanding 3GPP LTE releases 10 and 11 radio access technologies," *Academic Press*, 2013.
- [10] P. A. Elias, S. Rangan, and T. S. Rappaport, "Low-rank spatial channel estimation for millimeter wave cellular systems," *IEEE Transactions on Wireless Communications*, vol. 16, no. 5, pp. 2748–2759, May 2017.
- [11] X. Li, J. Fang, H. Li, and P. Wang, "Millimeter wave channel estimation via exploiting joint sparse and low-rank structures," *IEEE Transactions on Wireless Communications*, vol. 17, no. 2, pp. 1123–1133, Feb 2018.
- [12] T. O’Shea and J. Hoydis, "An introduction to deep learning for the physical layer," *IEEE Transactions on Cognitive Communications and Networking*, vol. 3, no. 4, pp. 563–575, Dec 2017.

- [13] S. Dörner, S. Cammerer, J. Hoydis, and S. t. Brink, “Deep learning based communication over the air,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 12, no. 1, pp. 132–143, Feb 2018.
- [14] D. Gündüz, P. Kerret, N. D. Sidiropoulos, D. Gesbert, C. Murthy, and M. Schaar, “Machine learning in the air,” *CoRR*, vol. abs/1904.12385, 2019. [Online]. Available: <http://arxiv.org/abs/1904.12385>
- [15] “3rd generation partnership project (3gpp), tr 25.996 v6.1.0 (2003-09) technical report,” available: <http://www.3gpp.org>., 2003.
- [16] P. Kyösti, “WINNER II channel models, DR1.1.2,” available: <https://www.ist-winner.org/WINNER2-Deliverables/D1.1.2v1.1.pdf>, 2007.
- [17] L. Liu, C. Oestges, J. Poutanen, K. Haneda, P. Vainikainen, F. Quitin, F. Tufvesson, and P. D. Doncker, “The COST 2100 MIMO channel model,” *IEEE Wireless Commun.*, vol. 19, no. 6, pp. 92–99, December 2012.
- [18] 3GPP, “Study on channel model for frequencies from 0.5 to 100 ghz,” *3rd Generation Partnership Project (3GPP), TR 38.901 V15.0.0*, June 2018.
- [19] A. Alkhateeb, “DeepMIMO: A generic deep learning dataset for millimeter wave and massive mimo applications,” in *Proc. of Information Theory and Applications Workshop (ITA)*, San Diego, CA, Feb 2019, pp. 1–8.
- [20] A. Alkhateeb, O. El Ayach, G. Leus, and R. W. Heath, “Channel estimation and hybrid precoding for millimeter wave cellular systems,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 8, no. 5, pp. 831–846, Oct 2014.
- [21] P. Schniter and A. Sayeed, “Channel estimation and precoder design for millimeter-wave communications: The sparse way,” in *2014 48th Asilomar Conference on Signals, Systems and Computers*, Nov 2014, pp. 273–277.
- [22] J. Lee, G. Gil, and Y. H. Lee, “Exploiting spatial sparsity for estimating channels of hybrid MIMO systems in millimeter wave communications,” in *2014 IEEE Global Communications Conference*, Dec 2014, pp. 3326–3331.
- [23] A. Alkhateeb, G. Leus, and R. W. Heath, “Compressed sensing based multi-user millimeter wave systems: How many measurements are needed?” in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, April 2015, pp. 2909–2913.
- [24] P. Dong, H. Zhang, G. Y. Li, I. S. Gaspar, and N. NaderiAlizadeh, “Deep CNN-based channel estimation for mmwave massive MIMO systems,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 13, no. 5, pp. 989–1000, Sep. 2019.
- [25] N. Shlezinger and Y. C. Eldar, “Deep task-based quantization,” *arXiv: 1908.06845[eess.SP]*, 2019.
- [26] A. Kaushik, E. Vlachos, J. Thompson, and A. Perelli, “Efficient channel estimation in millimeter wave hybrid MIMO systems with low resolution ADCs,” in *2018 26th European Signal Processing Conference (EUSIPCO)*, Sep. 2018, pp. 1825–1829.
- [27] Y. Ding, S. Chiu, and B. D. Rao, “Bayesian channel estimation algorithms for massive MIMO systems with hybrid analog-digital processing and low-resolution ADCs,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 12, no. 3, pp. 499–513, June 2018.
- [28] J. Mo, P. Schniter, and R. W. Heath, “Channel estimation in broadband millimeter wave MIMO systems with few-bit ADCs,” *IEEE TRANSACTIONS ON SIGNAL PROCESSING*, vol. 66, no. 5, pp. 1141–1154, MARCH 2018.
- [29] N. J. Myers, K. N. Tran, and R. W. H. Jr, “Low-rank mmwave MIMO channel estimation in one-bit receivers,” 2019.
- [30] Y. Zhang, M. Alrabeiah, and A. Alkhateeb, “Deep learning for massive MIMO with 1-bit ADCs: When more antennas need fewer pilots,” 2019.
- [31] E. Balevi and J. G. Andrews, “Two-stage learning for uplink channel estimation in one-bit massive MIMO,” 2019.
- [32] S. Gao, P. Dong, Z. Pan, and G. Y. Li, “Deep learning based channel estimation for massive MIMO with mixed-resolution ADCs,” *IEEE Communications Letters*, vol. 23, no. 11, pp. 1989–1993, Nov 2019.
- [33] 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, 2018.

- [34] T. Wang, C. 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, April 2019.
- [35] C. Lu, W. Xu, H. Shen, J. Zhu, and K. Wang, "MIMO channel information feedback using deep recurrent network," *IEEE Wireless Commun. Lett.*, vol. 23, no. 1, pp. 188–191, Jan 2019.
- [36] 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.*, 2019.
- [37] Y. Liao, H. Yao, Y. Hua, and C. Li, "CSI feedback based on deep learning for massive MIMO systems," *IEEE Access*, vol. 7, pp. 86 810–86 820, 2019.
- [38] Z. Lu, J. Wang, and J. Song, "Multi-resolution CSI feedback with deep learning in massive MIMO system," 2019.
- [39] Q. Sun, Y. Wu, J. Wang, C. Xu, and K. Wong, "CNN-based CSI acquisition for FDD massive MIMO with noisy feedback," *Electronics Lett.*, vol. 55, no. 17, pp. 963–965, 2019.
- [40] M. B. Mashhadi, Q. Yang, and D. Gunduz, "CNN-based analog CSI feedback in FDD MIMO-OFDM systems," 2019.
- [41] Y. Jang, G. Kong, M. Jung, S. Choi, and I. 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.
- [42] 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," *arXiv: 1906.06007[eess.SP]*, 2019.
- [43] Q. Yang, M. B. Mashhadi, and D. Gündüz, "Deep convolutional compression for massive MIMO CSI feedback," in *2019 IEEE 29th International Workshop on Machine Learning for Signal Processing (MLSP)*, Oct 2019, pp. 1–6.
- [44] 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, 2003.
- [45] Y. Polyanskiy, H. V. Poor, and S. Verdú, "Channel coding rate in the finite blocklength regime," *IEEE Trans. Inform. Theory*, vol. 56, no. 5, pp. 2307–2359, may 2010.
- [46] A. Wyner and J. Ziv, "The rate-distortion function for source coding with side information at the decoder," *IEEE Transactions on Information Theory*, vol. 22, no. 1, pp. 1–10, January 1976.
- [47] Y. Yang, F. Gao, G. Y. Li, and M. Jian, "Deep learning-based downlink channel prediction for FDD massive MIMO system," *IEEE Communications Letters*, vol. 23, no. 11, pp. 1994–1998, Nov 2019.
- [48] M. Alrabeiah and A. Alkhateeb, "Deep learning for TDD and FDD massive MIMO: Mapping channels in space and frequency," 2019.