# Joint Antenna Selection and Hybrid Beamformer Design using Unquantized and Quantized Deep Learning Networks

Ahmet M. Elbir and Kumar Vijay Mishra

*Abstract*—In millimeter wave communications, multiple-input-multiple-output (MIMO) systems use large antenna arrays to achieve high gain and spectral efficiency. These massive MIMO systems employ hybrid beamformers to reduce power consumption associated with fully digital beamforming in large arrays. Further savings in cost and power is possible through use of subarrays. Unlike prior works which resort to large latency methods such as optimization and greedy search for subarray selection, we propose a deep-learning-based approach in order to overcome the complexity issue without causing significant performance loss. We formulate antenna selection and hybrid beamformer design as a classification/prediction problem for convolutional neural networks (CNNs). For antenna selection, the CNN accepts the channel matrix as input and outputs a subarray with an optimal spectral efficiency. The resultant subarray channel matrix is then again fed to a CNN to obtain analog and baseband beamformers. We train the CNNs with several noisy channel matrices that have different channel statistics in order to achieve a robust performance at the network output. Numerical experiments show that our CNN framework provides an order better spectral efficiency and is 10 times faster than the conventional techniques. Further investigations with quantized-CNNs show that the proposed network, saved in no more than 5 bits, is also suited for digital mobile devices.

*Index Terms*—Antenna selection, CNN, deep learning, hybrid beamforming, massive MIMO.

## I. INTRODUCTION

The conventional cellular communications systems suffer from spectrum shortage while the demand for wider bandwidth and higher data rates is continuously increasing [1]. In this context, millimeter wave (mm-Wave) band, formally defined with the frequency range 30-300 GHz, is a preferred candidate for fifth-generation (5G) communications technology [2]–[5]. Compared to sub-6 GHz transmissions envisaged in 5G, the mm-Wave signals encounter a more complex propagation environment that is characterized by higher scattering, severe penetration losses, lower diffraction, and higher path loss for fixed transmitter and receiver gains [1]. These losses are compensated by providing beamforming power gain through massive number of antennas at both transmitter and receiver. Such a massive multiple-input-multiple-output (MIMO) structure [6], [7] enhances the signal-to-noise ratio (SNR) at the reception.

The wide mm-Wave bandwidth enables higher data rates for communications. Since the Nyquist sampling rate is twice the baseband bandwidth, mm-Wave receivers require expensive, high-rate analog-to-digital converters (ADCs) [8], [9]. The power consumption of an ADC increases with the sampling rate, more so at high frequencies [10], for a given architecture. At baseband, each full-resolution ADC consumes 15-795 mW at 36 MHz-1.8 GHz bandwidths. In addition, power consumed by other RF elements such as power amplifiers and data interface circuits in conjunction with large arrays renders it infeasible to utilize a separate radio- and intermediate-frequency (RF-IF) chain for each element. To reduce these cost-power-hardware overheads and yet provide reasonable performance, hybrid beamforming architectures have been proposed for massive MIMO. Here, the signal is processed by both analog and digital beamformers [11]–[15].

In the analog processing section of hybrid systems, it is common to employ phase shifters with constant modulus. Using analog switches, which are much simpler and cheaper than the phase shifters, it is possible to further make the overall system more energy-efficient by using subarrays of the larger full antenna array [16], [17]. Optimal selection of subarray elements reduces the power consumption of the analog phase shifters and low-noise amplifiers (LNAs) [18]–[20]. Very recent works consider the problem of antenna selection jointly with hybrid beamformer design to optimally trade-off cost and power efficiency [19]–[22]. In particular, [21] proposed antenna selection and analog precoder design with low-resolution phase shifters for multiple-input-single-output (MISO) systems. The hybrid beamformer designs suggested in [19] and [20] involve a sub-optimum antenna selection strategy through quadratic approximation with smooth optimization. Similarly, the massive MIMO architecture in [22] employs greedy search for choosing sub-optimal subarrays.

Nearly all of these works provide sub-optimum solutions despite attempting various antenna selection criteria and optimization strategies. Even while using branch-and-bound (BAB) algorithms - which provide good estimation of the lower and upper bounds of regions/branches of the search space in polynomial time - obtaining an optimum solution for massive MIMO subarray selection requires high computational burden [23]. In this paper, to reduce the complexity (in cases where the optimum solution can still be obtained), we introduce an approach based on deep learning (DL) to find an optimum subarray jointly with the design of hybrid beamformers; the optimality is in the sense of achieving maximum spectral efficiency.

As a class of machine learning techniques, DL methods have gained much interest recently for solving many challenging problems such as visual object recognition [24], rainfall estimation [25], and language processing [26]. These techniques offer advantages such as low computational complexity while solving optimization-based or combinatorial search problems as well as the ability to extrapolate new features from a limited set of features contained in a training set [24]. Very recently, DL has received significant attention in addressing problems in communications signal processing such as channel estimation [27], direction-of-arrival (DoA) estimation [28] analog beam selection [29]–[31] and beam management in dense mm-Wave networks [32]. At the physical layer of wireless communications, DL has been applied for signal detection [33] and channel estimation [34]. An end-to-end single-input-single-output (SISO) communications scenario is modeled in [34] and [35] by using auto-encoders. In [36], auto-encoders are employed for channel state information (CSI) feedback. A sub-optimum method based on support vector machines (SVMs) is proposed in [28] for selecting analog beamforming vector. In a recent work [37], multilayer perceptrons (MLP) are used for precoder design in a single-user mm-Wave scenario.

In this paper, we exploit DL to simultaneously select antenna elements and design hybrid beamformer. This joint problem as well as the stand-alone hybrid beamforming remain unexamined in the previous DL works. Specifically, we design a convolutional neural

A. M. E. is with the Department of Electrical and Electronics Engineering, Duzce University, Duzce, Turkey. E-mail: ahmetelbir@duzce.edu.tr.

K. V. M. is with The University of Iowa, Iowa City, IA 52242 USA. E-mail: kumarvijay-mishra@uiowa.edu.

network (CNN) to achieve both tasks sequentially. The element selection problem is cast as a classification problem [38]. A similar DL approach was adopted for radar antenna arrays recently in [39]. We further incorporate the hybrid beamformer design in this DL framework by exploiting the structure of analog beamformers which are obtained by minimizing the cost between hybrid and unconstrained beamformers. The optimization problem is cast jointly with the antenna selection problem and it is solved by MATLAB-based Manopt algorithm [40] via manifold optimization (MO) [15].

In our formulation, a CNN accepts channel matrix as input and provides the subarray that maximizes the spectral efficiency. Once the antenna selection is finalized, the corresponding partial channel matrix is fed to a second CNN which then chooses the best RF beamformer and constructs the corresponding baseband beamformer. To train both CNN models, different realizations of the channel matrix are used and the input data are labeled by the selected subarray/RF chains with the highest spectral efficiency. Even though our proposed network structures require channel matrix as an input, precise knowledge of this matrix is only required in the training stage to obtain the labels of the network. In the prediction stage, where the RF beamformers are estimated, precise channel knowledge is not necessary. Both CNNs are trained with channel matrices generated for different user location, channel gains and number of user clusters. Furthermore, each realization of channel matrix in the training data is corrupted by synthetic noise so that the performance of the learning network does not deteriorate with noisy test inputs.

We evaluate the performance of the proposed framework over several experiments and show that proposed CNN approach provides significantly better performance as compared to the optimization and greedy based techniques [9], [14], [37], [41]. In order to account for time-varying channel and user parameters, we use several channel realizations with added noise. We train the networks with huge training data (~240000 input samples) with noisy channel matrices. As a result, the classification accuracy quickly reaches 100% wherein optimum antenna selection and RF beamformer design are accomplished. The CNN is trained offline and, hence, all the computational overhead is taken into account for data generation and training. The classification and prediction time for our proposed approach is at least 10 times faster than the conventional antenna selection techniques as well as hybrid beamformer design algorithms.

Finally, our approach is helpful in reducing the computational burden involved in hybrid beamformers by simply feeding the channel matrix to the network. This requires using CNNs in mobile devices where the data are collected in digital form. Since existing deep neural network models are computationally and memory intensive, they cannot be deployed in devices with low memory resources and low overhead requirements. These constraints have driven investigation into compression of deep neural networks. One of the common approaches is to quantize the CNN weights [42], [43]. In this paper, we investigate the performance of the proposed framework when the weights of the CNNs are quantized. While quantized-CNN structures are recently studied for image classification purposes, ours is the first work that examines quantized-CNNs for communications. Preliminary results of our work appeared in [38] and [42]; while a basic formulation suggested in [38] solved the joint problem for a specific hybrid beamforming scheme, [42] proposed a quantized-CNN approach for the same scheme but did not include antenna selection.

Throughout the paper, we denote the identity matrix of size $N \times N$ as $\mathbf{I}_N$. $(\cdot)^T$ and $(\cdot)^H$ denote transpose and the conjugate transpose operations, respectively. For a matrix $\mathbf{A}$ and a vector $\mathbf{a}$, $[\mathbf{A}]_{:,i}$ and $[\mathbf{A}]_{i,j}$ denote the $i$th column and $(i,j)$th element of matrix $\mathbf{A}$, $[\mathbf{a}]_i$ means the $i$th element of vector $\mathbf{a}$, respectively. The notation $|\mathbf{A}|$

denotes the determinant of matrix $\mathbf{A}$ whereas $|a|$ is the absolute value of the scalar $a$. The function $\mathbb{E}\{\cdot\}$ provides the statistical expectation of its argument and $\angle\{\cdot\}$ measures the angle of complex quantity.

## II. System Model For mm-Wave MIMO Systems

Consider a single user mm-Wave MIMO system with $N_T$ and $N_R$ transmit and receive antennas, respectively (Fig. 1). Assume that $N_S$ data streams are desired to be transmitted to the receiver where the antenna selection is performed to select a subarray with $N_{RS}$ antennas out of $N_R$. There are $N_T^{RF}$ and $N_R^{RF}$ RF beamformers at transmit and receive sides such that $N_S \leq N_T^{RF} \leq N_T$ and $N_S \leq N_R^{RF} \leq N_{RS} \leq N_R$. The hybrid precoder structure applies the baseband precoder $\mathbf{F}_{BB} \in \mathbb{C}^{N_T^{RF} \times N_S}$ to the transmit signal vector $\mathbf{s} \in \mathbb{C}^{N_S}$, where $\mathbb{E}\{\mathbf{ss}^H\} = \mathbf{I}_{N_S}/N_S$. Then, the signal is passed through RF precoders $\mathbf{F}_{RF} \in \mathbb{C}^{N_T \times N_T^{RF}}$ (constructed using phase shifters) to $N_T$ transmit antennas. The RF precoder has equal-norm elements so that $[[\mathbf{F}_{RF}]_{:,i}[\mathbf{F}_{RF}]_{:,i}^H]_{i,i} = 1/N_T$. The power of the transmitter is constrained to $||\mathbf{F}_{RF}\mathbf{F}_{BB}||_{\mathcal{F}} = N_S$. The transmitted signal at RF stage is $\mathbf{x} = \mathbf{F}_{RF}\mathbf{F}_{BB}\mathbf{s} \in \mathbb{C}^{N_T}$. Assuming a narrowband block-fading channel, the received signal at $N_R$ antennas is [14], [44]

$$\mathbf{y}^{\text{Full}} = \sqrt{\rho}\mathbf{H}\mathbf{F}_{RF}\mathbf{F}_{BB}\mathbf{s} + \mathbf{n}, \qquad (1)$$

where $\mathbf{y}^{\text{Full}} \in \mathbb{C}^{N_R}$ is the output of $N_R$ antennas at the receiver, $\rho$ is average received power, $\mathbf{n} \in \mathbb{C}^{N_R}$ is the additive white Gaussian noise (AWGN) with $\mathbf{n} \sim \mathcal{CN}(\mathbf{0}, \sigma_n^2 \mathbf{I}_{N_R})$, and $\mathbf{H} \in \mathbb{C}^{N_R \times N_T}$ is the channel matrix with $||\mathbf{H}||_{\mathcal{F}} = N_R N_T$.

A standard analog-digital hybrid beamformer operates on a full array. When a subarray is employed, then beamformers must be derived for reduced dimensions. We model the antenna selection problem to select the *best* antenna subset from all possible antenna subarrays. Throughout this paper, we use the terms "subarray" and "subset" interchangeably; a subarray being a subset of indices (or antenna positions) of a full array configuration. A popular scheme is to employ a predetermined subarray with $N_{RS}$ antennas is selected from a full array of $N_R$ elements (Fig. 1a). Each subarray feeds into a fully-connected phase shifter network of size $N_R^{RF}$ with a single RF chain. This has the complexity of phase shifters but the antenna selection process is not optimized. Another common receiver architecture feeds the antennas directly to the RF chains thereby eliminating the phase shifters completely. Here, each RF chain is connected to the $N_R$ antennas of which $N_{RS}$ elements are selected using switches (Fig. 1b). In this case, the entries of the combiner matrix are either 1 or 0 to indicate the selected or unselected antennas, respectively. This is the simplest structure with no phase shifters. However, the antenna selection is not optimzed and the elements are determined by simply choosing the largest absolute values in each column of the channel matrix. Finally, Fig. 1c shows a receiver that employs a switching network with phase shifters. In this system, a subarray with $N_{RS}$ antennas is selected from a full array comprising $N_R$ antennas. The subarray is connected to a phase shifter network of size $N_R^{RF}$ which may apply an optimization procedure for antenna selection to achieve greater efficiency.

Our mm-Wave channel representation is based on the Saleh-Valenzuela (SV) model that utilizes the clustered channel model [45], [46]. Here, the channel matrix $\mathbf{H}$ includes the contributions of $N_c$ scattering clusters, each of which has $N_{\text{ray}}$ paths. We have

$$\mathbf{H} = \gamma \sum_{i,j} \alpha_{ij} g_R(\Theta_R^{(ij)}) g_T(\Theta_T^{(ij)}) \mathbf{a}_R(\Theta_R^{(ij)}) \mathbf{a}_T^H(\Theta_T^{(ij)}), \quad (2)$$

where $\Theta_R^{(ij)} = (\phi_R^{(ij)}, \theta_R^{(ij)})$ and $\Theta_T^{(ij)} = (\phi_T^{(ij)}, \theta_T^{(ij)})$, respectively, denote the angle of arrivals and angle of departures
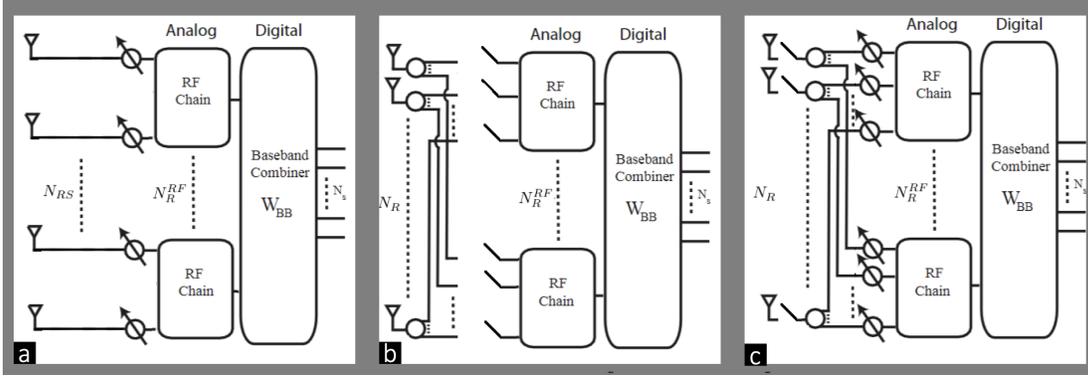
Fig. 1. Receiver architectures with antenna selection for single user mmWave MIMO systems. (a) Scheme 1: Fixed subarray with fully-connected phase shifters. (b) Scheme 2: Switching network without optimized antenna selection and no phase shifters. (c) Scheme 3: Switching network with optimized antenna selection and phase shifters.

wherein the azimuth (elevation) angle is denoted by $\phi$ ($\theta$), $\gamma = \sqrt{N_T N_{RS}/(N_c N_{\text{ray}})}$ is the normalization factor, and $\alpha_{ij}$ is the complex channel gain associated with the $i$th scattering cluster and $j$th path for $i = 1, \ldots, N_c$ and $j = 1, \ldots, N_{\text{ray}}$. The antenna element gains for receive and transmit antennas are $g_R(\Theta_R^{(ij)})$ and $g_T(\Theta_T^{(ij)})$, respectively. The steering vector representing the array response at the transmitter (receiver) is $\mathbf{a}_T(\Theta_T^{(ij)}) \in N_T \times 1$ ($\mathbf{a}_R(\Theta_R^{(ij)}) \in N_R \times 1$). The $n$th element of $\mathbf{a}_R(\Theta_R^{(ij)})$ is

$$[\mathbf{a}_R(\Theta_R^{(ij)})]_n = \exp\left\{-\frac{2\pi}{\lambda}\mathbf{p}_n^T \mathbf{r}(\Theta_R^{(ij)})\right\}, \qquad (3)$$

where $\mathbf{p}_n = [x_n, y_n, z_n]^T$ is the position of the $n$th antenna in Cartesian coordinate system and $\mathbf{r}(\Theta_R^{(ij)}) = [\sin(\phi_R^{(ij)})\cos(\theta_R^{(ij)}), \sin(\phi_R^{(ij)})\sin(\theta_R^{(ij)}), \cos(\theta_R^{(ij)})]^T$. The transmit steering vector $\mathbf{a}_T(\Theta_T^{(ij)})$ is defined similarly.

In practice, the estimation process of the channel matrix is a challenging task, especially in case of a large number of antennas in massive MIMO systems [47], [48]. Additionally, the mm-Wave channel has short coherence times [1], [49]. In practice, the estimated channel matrix could be obtained via one of the several channel estimation techniques [12], [46], [48], [50], [51]. For robust performance against imperfect channel estimates, our proposed DL framework feeds the deep network with several channel realizations which are corrupted by synthetic noise in the training stage. This is an offline process. Note that the perfect knowledge of channel matrix $\mathbf{H}$ is only required[1] at the training stage to obtain the labels (e.g., hybrid beamformer matrices). During the testing stage when the network predicts the beamformer weights, the network does not necessarily require the perfect CSI. Our numerical experiments demonstrate that the proposed approach can handle the corrupted channel matrix case and exhibits satisfactory performance regarding the achievable spectral efficiency.

In the hybrid beamformer, analog and digital beamformers are obtained to maximize the spectral efficiency. Often, this is achieved by exploiting the structure of the mm-Wave channel matrix [51]. Using the full antenna array at the receiver, the received signal in (1)

is processed by analog and baseband combiners to yield

$$\bar{\mathbf{y}} = \mathbf{W}_{BB}^{\text{Full}\,H} \mathbf{W}_{RF}^{\text{Full}\,H} \mathbf{y}^{\text{Full}}$$
$$= \sqrt{\rho}\mathbf{W}_{BB}^{\text{Full}\,H} \mathbf{W}_{RF}^{\text{Full}\,H} \mathbf{H}\mathbf{F}_{RF}\mathbf{F}_{BB}\mathbf{s} + \mathbf{W}_{BB}^{\text{Full}\,H} \mathbf{W}_{RF}^{\text{Full}\,H}\mathbf{n}, \qquad (4)$$

where $\mathbf{W}_{RF}^{\text{Full}} \in \mathbb{C}^{N_R \times N_R^{RF}}$ is the analog combiner with the constrained $[[\mathbf{W}_{RF}^{\text{Full}}]_{:,i}[\mathbf{W}_{RF}^{\text{Full}}]_{:,i}^H]_{i,i} = 1/N_R$ and $\mathbf{W}_{BB}^{\text{Full}} \in \mathbb{C}^{N_R^{RF} \times N_S}$ denotes the baseband combiner matrix. Assuming that the Gaussian symbols are transmitted through the mm-Wave channel, we define the spectral efficiency [11]–[14] achieved from the full array as

$$R^{\text{Full}} = \log_2 \left| \mathbf{I}_{N_S} + \frac{\rho}{N_S}\mathbf{\Lambda}_n^{\text{Full}^{-1}} \mathbf{W}_{BB}^{\text{Full}\,H} \mathbf{W}_{RF}^{\text{Full}\,H} \mathbf{H} \right.$$
$$\left. \times \mathbf{F}_{RF}\mathbf{F}_{BB}\mathbf{F}_{BB}^H \mathbf{F}_{RF}^H \mathbf{H}^H \mathbf{W}_{RF}^{\text{Full}} \mathbf{W}_{BB}^{\text{Full}} \right|, \qquad (5)$$

where $\mathbf{\Lambda}_n^{\text{Full}} = \sigma_n^2 \mathbf{W}_{BB}^{\text{Full}\,H} \mathbf{W}_{RF}^{\text{Full}\,H} \mathbf{W}_{RF}^{\text{Full}} \mathbf{W}_{BB}^{\text{Full}} \in \mathbb{C}^{N_S \times N_S}$ is the covariance matrix of the noise term in (4) after analog combining. We now formulate the problem for subarray selection and obtaining the corresponding analog-digital beamformer in the following section.

## III. Joint Antenna and RF Chain Selection

Among the antenna selection schemes presented in the previous section, we focus on the Scheme 3 because this architecture requires optimization (the remaining configurations consider selecting a fixed subarray with/without phase shifters). In particular, our goal is to select the outputs of $N_{RS}$ antennas from the full array output $\mathbf{y}^{\text{Full}}$. Consequently, this also requires designing transmit and receive analog and baseband beamformers $\mathbf{F}_{RF} \in \mathbb{C}^{N_T \times N_T^{RF}}$, $\mathbf{W}_{RF} \in \mathbb{C}^{N_{RS} \times N_R^{RF}}$ and $\mathbf{F}_{BB} \in \mathbb{C}^{N_T^{RF} \times N_S}$, $\mathbf{W}_{BB} \in \mathbb{C}^{N_R^{RF} \times N_S}$. In other words, the solution of joint antenna selection and hybrid beamformer design satisfies

$$\underset{\mathbf{Q},\mathbf{F}_{RF},\mathbf{F}_{BB},\mathbf{W}_{RF},\mathbf{W}_{BB}}{\text{maximize}} \log_2 \left| \mathbf{I}_{N_S} + \frac{\rho}{N_S \sigma_n^2} \left(\mathbf{W}_{BB}^H \mathbf{W}_{RF}^H \mathbf{W}_{RF} \mathbf{W}_{BB}\right)^{-1} \right.$$
$$\left. \mathbf{W}_{BB}^H \mathbf{W}_{RF}^H \mathbf{H}_{\text{sub}}\mathbf{F}_{RF}\mathbf{F}_{BB}\mathbf{F}_{BB}^H \mathbf{F}_{RF}^H \mathbf{H}_{\text{sub}}^H \mathbf{W}_{RF}\mathbf{W}_{BB} \right|$$

$$\text{subject to:} \quad \mathbf{F}_{RF} \in \mathcal{F}_{RF}, \ \|\mathbf{F}_{RF}\mathbf{F}_{BB}\|_{\mathcal{F}}^2 = N_S,$$
$$\mathbf{W}_{RF} \in \mathcal{W}_{RF}, \ \|\text{diag}\{\mathbf{Q}\}\|_0 = N_{RS}, \qquad (6)$$

where $\mathcal{F}_{RF}$ and $\mathcal{W}_{RF}$ denote the feasible sets of analog beamformers, $\mathbf{H}_{\text{sub}} = \mathbf{Q}\mathbf{H}$ is $N_{RS} \times N_T$ channel matrix of the selected antennas, and $\mathbf{Q}$ is the $N_{RS} \times N_R$ selection matrix whose $(i, j)$th entry is either 1 or 0. Even without antenna selection, the problem in (6) is difficult to solve because of several matrix variables $\mathbf{Q}$, $\mathbf{F}_{RF}$, $\mathbf{W}_{RF}$

---

[1]In order to achieve antenna selection and train the deep network, full channel information of size $N_R \times N_T$ is required. This is also a common requirement in most antenna selection algorithms [19], [20], [23], [52], [53]. After antenna selection, the hybrid beamforming network requires only the partial channel matrix of size $N_{RS} \times N_T$ corresponding to the selected antennas.

and $\mathbf{F}_{BB}$, $\mathbf{W}_{BB}$ [23], [54]. Since obtaining a solution to (6) in real-time is deemed infeasible, we propose a deep learning approach here to achieve an optimum solution with less computational complexity. We first cast the antenna selection stage as a classification problem as follows.

### A. Antenna Selection

In subarray selection, we are interested in picking $N_{RS}$ out of $N_R$ elements. This yields $Q_A = \begin{pmatrix} N_R \\ N_{RS} \end{pmatrix}$ possible solutions. Therefore, choosing subarrays can be viewed as a classification problem with $Q_A$ classes. We define $\mathbb{S}$ as the set of all possible antenna subarray configurations, i.e., $\mathbb{S} = \{\mathbb{S}_1, \mathbb{S}_2, \ldots, \mathbb{S}_{Q_A}\}$, where $\mathbb{S}_{q_A} = \{\mathbf{p}_1^{q_A}, \mathbf{p}_2^{q_A}, \ldots, \mathbf{p}_{N_{RS}}^{q_A}\}$ includes the antenna positions of the $q_A$th subarray configuration with $q_A \in \mathbb{Q}_A = \{1, \ldots, Q_A\}$. Let $\mathbf{y}_{q_A}$ be an $N_{RS} \times 1$ vector containing the output signal of the selected antennas for the $q_A$th subarray configuration of the full array output $\mathbf{y}^{\text{Full}}$ with positions $\mathbb{S}_{q_A}$. Then,

$$\mathbf{y}_{q_A} = \sqrt{\rho}\mathbf{H}_{q_A}\mathbf{F}_{RF}\mathbf{F}_{BB}\mathbf{s} + \mathbf{n}_{q_A}, \tag{7}$$

where $\mathbf{H}_{q_A}$ is the $N_{RS} \times N_T$ channel matrix with selected antennas and $\mathbf{n}_{q_A}$ is similarly defined. At the receiver, analog and the baseband combiners - $\mathbf{W}_{RF} \in \mathbb{C}^{N_{RS} \times N_R^{RF}}$ and $\mathbf{W}_{BB} \in \mathbb{C}^{N_R^{RF} \times N_S}$, respectively - are applied to the received signal to produce the $N_S \times 1$ discrete-time signal $\bar{\mathbf{y}}_{q_A} = \mathbf{W}_{BB}^H \mathbf{W}_{RF}^H \mathbf{y}_{q_A}$ as

$$\bar{\mathbf{y}}_{q_A} = \sqrt{\rho}\mathbf{W}_{BB}^H \mathbf{W}_{RF}^H \mathbf{H}_{q_A}\mathbf{F}_{RF}\mathbf{F}_{BB}\mathbf{s} + \mathbf{W}_{BB}^H \mathbf{W}_{RF}^H \mathbf{n}_{q_A}. \tag{8}$$

When the $q_A$th subarray is selected, the spectral efficiency [55] of the mm-Wave channel is

$$R(q_A) = \log_2\left|\mathbf{I}_{N_S} + \frac{\rho}{N_S}\mathbf{\Lambda}_n^{-1}\mathbf{W}_{BB}^H \mathbf{W}_{RF}^H \mathbf{H}_{q_A}\right.$$
$$\left. \times \mathbf{F}_{RF}\mathbf{F}_{BB}\mathbf{F}_{BB}^H \mathbf{F}_{RF}^H \mathbf{H}_{q_A}^H \mathbf{W}_{RF}\mathbf{W}_{BB}\right|, \tag{9}$$

where $\mathbf{\Lambda}_n = \sigma_n^2 \mathbf{W}_{BB}^H \mathbf{W}_{RF}^H \mathbf{W}_{RF}\mathbf{W}_{BB} \in \mathbb{C}^{N_S \times N_S}$ corresponds to the noise term of the subarray output in (7). Note that $R(q_A)$ depends on $q_A$ through $\mathbf{H}_{q_A}$. By maximizing the spectral efficiency in (9) over all subarray configurations, the best antenna subarray is obtained as

$$\bar{q}_A = \underset{q_A \in \mathbb{Q}_A}{\operatorname{argmax}} R(q_A), \tag{10}$$

where $\bar{q}_A$ denotes the subarray index with antenna positions $\mathbb{S}_{\bar{q}_A}$ which provide the maximum spectral efficiency.

While the optimization problem in (10) yields the best subarray configuration, it does not impose any constraint on the hybrid beamformers. Here the problem in (10) is cast by using unconstrained beamformers $\mathbf{F}_{q_A}^{\text{opt}} \in \mathbb{C}^{N_T \times N_S}$ and $\mathbf{W}_{q_A}^{\text{opt}} \in \mathbb{C}^{N_{RS} \times N_S}$. These can be obtained from the singular value decomposition (SVD) of the $N_{RS} \times N_T$ complex-valued channel matrix: $\mathbf{H}_{q_A} = \mathbf{U}_{q_A}\mathbf{\Sigma}_{q_A}\mathbf{V}_{q_A}^H$, where $\mathbf{U}_{q_A} \in \mathbb{C}^{N_{RS} \times \operatorname{rank}(\mathbf{H}_{q_A})}$ and $\mathbf{V}_{q_A} \in \mathbb{C}^{N_T \times \operatorname{rank}(\mathbf{H}_{q_A})}$ are the left and the right singular value matrices of the $q_A$th channel matrix, respectively, and $\mathbf{\Sigma}_{q_A}$ is $\operatorname{rank}(\mathbf{H}_{q_A}) \times \operatorname{rank}(\mathbf{H}_{q_A})$ matrix composed of the singular values of $\mathbf{H}_{q_A}$ in descending order. By decomposing $\mathbf{\Sigma}_{q_A}$ and $\mathbf{V}_{q_A}$ as $\mathbf{\Sigma}_{q_A} = \operatorname{diag}\{\mathbf{\Sigma}_{q_A}^{(1)}, \mathbf{\Sigma}_{q_A}^{(2)}\}$, $\mathbf{V}_{q_A} = [\mathbf{V}_{q_A}^{(1)}, \mathbf{V}_{q_A}^{(2)}]$, where $\mathbf{V}_{q_A}^{(1)} \in \mathbb{C}^{N_T \times N_S}$ and $\mathbf{V}_{q_A}^{(2)} \in \mathbb{C}^{N_T \times N_{RS} - N_S}$, one can readily select the unconstrained precoder as $\mathbf{F}_{q_A}^{\text{opt}} = \mathbf{V}_{q_A}^{(1)}$ [14]. Using the unconstrained beamformer $\mathbf{F}_{q_A}^{\text{opt}}$, $\mathbf{W}_{q_A}^{\text{opt}}$ is computed as [56]

$$\mathbf{W}_{q_A}^{\text{opt}} = \left(\frac{1}{\rho}\left(\mathbf{F}_{q_A}^{\text{opt}H}\mathbf{H}_{q_A}^H\mathbf{H}_{q_A}\mathbf{F}_{q_A}^{\text{opt}} + \frac{N_S\sigma_n^2}{\rho}\mathbf{I}_{N_S}\right)^{-1}\mathbf{F}_{q_A}^{\text{opt}H}\mathbf{H}_{q_A}^H\right)^H.$$

Using $\mathbf{F}_{q_A}^{\text{opt}}$ and $\mathbf{W}_{q_A}^{\text{opt}}$, the following problem can be written, i.e.,

$$\underset{q_A \in \mathbb{Q}_A}{\operatorname{maximize}} \log_2\left|\mathbf{I}_{N_S} + \frac{\rho}{N_S\sigma_n^2}(\mathbf{W}_{q_A}^{\text{opt}H}\mathbf{W}_{q_A}^{\text{opt}})^{-1}\mathbf{W}_{q_A}^{\text{opt}H}\mathbf{H}_{q_A}\right.$$
$$\left. \times \mathbf{F}_{q_A}^{\text{opt}}\mathbf{F}_{q_A}^{\text{opt}H}\mathbf{H}_{q_A}^H\mathbf{W}_{q_A}^{\text{opt}}\right|. \tag{11}$$

The optimization problem in (11) uses unconstrained beamformers $\mathbf{F}^{\text{opt}}$, $\mathbf{W}^{\text{opt}}$ for antenna selection and hybrid beamformer design. However, the "best" subarray obtained from (11) would be different if hybrid beamformers are used in the problem. Hence, we consider the joint problem with hybrid beamformers and write the joint antenna selection and hybrid beamformer design problem as

$$\underset{q_A \in \mathbb{Q}_A}{\operatorname{maximize}} \log_2\left|\mathbf{I}_{N_S} + \frac{\rho}{N_S\sigma_n^2}(\mathbf{W}_{BB_{q_A}}^H\mathbf{W}_{RF_{q_A}}^H\mathbf{W}_{RF_{q_A}}\mathbf{W}_{BB_{q_A}})^{-1}\mathbf{W}_{BB_{q_A}}^H\right.$$
$$\left. \times \mathbf{W}_{RF_{q_A}}^H\mathbf{H}_{q_A}\mathbf{F}_{RF_{q_A}}\mathbf{F}_{BB_{q_A}}\mathbf{F}_{BB_{q_A}}^H\mathbf{F}_{RF_{q_A}}^H\mathbf{H}_{q_A}^H\mathbf{W}_{RF_{q_A}}\mathbf{W}_{BB_{q_A}}\right|$$
$$\text{subject to:} \quad \mathbf{F}_{RF_{q_A}} \in \mathcal{F}_{RF}, \quad ||\mathbf{F}_{RF_{q_A}}\mathbf{F}_{BB_{q_A}}||_{\mathcal{F}}^2 = N_S,$$
$$\mathbf{W}_{RF_{q_A}} \in \mathcal{W}_{RF}, \tag{12}$$

which requires to solve hybrid beamformer design problem for each $q_A \in \mathbb{Q}_A$. Let us, for now, assume that the hybrid beamformers are estimated as described in the next subsection. Then, the antenna subarray that provides the maximum spectral efficiency is obtained and this subarray is represented by the subarray index $\bar{q}_A$. When the problem in (12) is solved for different channel matrices, some of the $\bar{q}_A$ values turn out to be the same for different channel matrices which are similar to each other and the same antenna subarray provides the maximum spectral efficiency for these channel matrices [39]. As a result, the number of subarrays providing the maximum spectral efficiency, say $\bar{Q}_A$, is much less than the number of all subarray configurations, i.e., $\bar{Q}_A \ll Q_A$. In [39], a similar observation is made for cognitive radar scenario. Therefore, we define another subset of subarray configurations as $\mathbb{A} = \{\mathbb{A}_1, \ldots, \mathbb{A}_{\bar{Q}_A}\}$, which is composed of the subarrays providing maximum spectral efficiency for different channel matrices. In other words, $\mathbb{A}_{\bar{q}_A}$ includes the antenna positions of the subarray obtained by solving (11), hence we have $\mathbb{A} \subset \mathbb{S}$.

For very large number of antennas, say, $N_R > 64$, the computation of all possible antenna subsets is computationally prohibitive and requires very large amount of memory. To tackle this problem, we use close-to-optimum branch-and-bound (BAB) techniques [23]. To further reduce the complexity, we also partition $\mathbb{S}$ into $B$ non-overlapping blocks as $\mathbb{S}^{(b)} = \mathbb{S}_{N_A(b-1)+1}, \ldots, \mathbb{S}_{N_Ab}$ where $N_A$ is the block size and $b = 1, \ldots, B$. Then the antenna selection problem is solved for $N_A$ nodes at a times hence less memory is used. In particular, the following strategy is used:

1) For $b = 1$, construct $\mathbb{S}^{(b)}$ and solve (12) as $\bar{q}_A^{(b)} = \operatorname{argmax}_{q_A \in \mathbb{S}^{(b)}} R(q_A)$.
2) For $b > 1$, clear $\mathbb{S}^{(b-1)}$ from the memory and construct $\mathbb{S}^{(b)}$, then solve (12) as $\bar{q}_A^{(b)} = \operatorname{argmax}_{q_A \in \mathbb{S}^{(b)}} R(q_A)$ and obtain the new best subarray index as

$$\bar{q}_A^{(b)} := \begin{cases} \bar{q}_A^{(b)}, & \text{if } R(\bar{q}_A^{(b-1)}) < R(\bar{q}_A^{(b)}) \\ \bar{q}_A^{(b-1)}, & \text{otherwise} \end{cases}. \tag{13}$$

### B. Hybrid Beamformer Design

In (12), antenna selection is performed by estimating the beamformer weights for each $q_A$. Let $\mathbf{H}_{q_A} \in \mathbb{C}^{N_{RS} \times N_T}$ be the selected channel matrix, then the hybrid design problem can be written as

$$\underset{\mathbf{F}_{RF}, \mathbf{F}_{BB}, \mathbf{W}_{RF}, \mathbf{W}_{BB}}{\operatorname{maximize}} \log_2\left|\mathbf{I}_{N_S} + \frac{\rho}{N_S\sigma_n^2}(\mathbf{W}_{BB}^H\mathbf{W}_{RF}^H\mathbf{W}_{RF}\mathbf{W}_{BB})^{-1}\right.$$
$$\left. \times \mathbf{W}_{BB}^H\mathbf{W}_{RF}^H\mathbf{H}_{q_A}\mathbf{F}_{RF}\mathbf{F}_{BB}\mathbf{F}_{BB}^H\mathbf{F}_{RF}^H\mathbf{H}_{q_A}^H\mathbf{W}_{RF}\mathbf{W}_{BB}\right|$$
$$\text{subject to: } \mathbf{F}_{RF} \in \mathcal{F}_{RF}, ||\mathbf{F}_{RF}\mathbf{F}_{BB}||_{\mathcal{F}}^2 = N_S, \mathbf{W}_{RF} \in \mathcal{W}_{RF}. \tag{14}$$

Above problem can be written as two decoupled optimization problems to find precoders and combiners separately [14], [15]. In this case, the Euclidean distance between the unconstrained beamformers and the hybrid beamformers is minimized. In other words, the hybrid precoder design problem can be written as follows

$$\underset{\mathbf{F}_{RF}, \mathbf{F}_{BB}}{\text{minimize}} \quad ||\mathbf{F}_{q_A}^{\text{opt}} - \mathbf{F}_{RF}\mathbf{F}_{BB}||_{\mathcal{F}}^2$$
$$\text{subject to: } \mathbf{F}_{RF} \in \mathcal{F}_{RF}. \quad (15)$$

In a similar way, the combiner design problem can be written as

$$\underset{\mathbf{W}_{RF}, \mathbf{W}_{BB}}{\text{minimize}} \quad ||\mathbf{W}_{q_A}^{\text{opt}} - \mathbf{W}_{RF}\mathbf{W}_{BB}||_{\mathcal{F}}^2$$
$$\text{subject to: } \mathbf{W}_{RF} \in \mathcal{W}_{RF},$$
$$\mathbf{W}_{BB} = (\mathbf{W}_{RF}^H \mathbf{\Lambda}_{q_A} \mathbf{W}_{RF})^{-1}(\mathbf{W}_{RF}^H \mathbf{\Lambda}_{q_A} \mathbf{W}_{q_A}^{\text{opt}}), \quad (16)$$

where $\mathbf{\Lambda}_{q_A} = \frac{\rho}{N_S}\mathbf{H}_{q_A}\mathbf{F}_{RF}\mathbf{F}_{BB}\mathbf{F}_{BB}^H\mathbf{F}_{RF}^H\mathbf{H}_{q_A}^H + \sigma_n^2\mathbf{I}_{N_{RS}}$, denotes the covariance of the array output in (7) corresponding to the selected subarray index $q_A$.

Above optimization problems in (15) and (16) can be effectively solved by MATLAB-based Manopt algorithm [40] via manifold optimization [15][2]. Once we obtain $\mathbf{F}_{RF}, \mathbf{F}_{BB}$ and $\mathbf{W}_{RF}, \mathbf{W}_{BB}$, the labels of the output layer of the network can be formed as

$$\mathbf{z} = [\text{vec}^T\{\angle\mathbf{F}_{RF}\}, \text{Re}\{\text{vec}^T\{\mathbf{F}_{BB}\}\}, \text{Im}\{\text{vec}^T\{\mathbf{F}_{BB}\}\},$$
$$\text{vec}^T\{\angle\mathbf{W}_{RF}\}, \text{Re}\{\text{vec}^T\{\mathbf{W}_{BB}\}\}, \text{Im}\{\text{vec}^T\{\mathbf{W}_{BB}\}\}]^T, \quad (17)$$

which is a $G \times 1$ real-valued vector and $G = N_T N_T^{RF} + N_{RS} N_R^{RF} + 2N_S(N_T^{RF} + N_R^{RF})$.

Even with the above memory-friendly approach, it is still computationally complex to search all possible antenna subsets in real-time. In order to circumvent this problem, we design a deep learning approach where the network is trained offline with the overhead containing the computation of all possible subarray and RF chain combinations. Then, the trained network can simply be employed as a classification network to select the best antenna subarray and the corresponding RF beamformers for the given channel matrix. We introduce the proposed deep learning technique in the following section.

## IV. TRAINING THE DL NETWORK

The proposed deep network comprises two CNNs (Fig. 2). The first ($\text{CNN}_{\text{AS}}$) accepts the input of channel matrix with the goal to select best antenna subarray $\bar{q}_A$. The second CNN ($\text{CNN}_{\text{RF}}$) takes the input of the subsequent channel matrix with selected rows to choose RF beamformers. For both CNNs, training data are selected from different channel matrix realizations each of which is assigned with the corresponding output classes.

Let $\mathbf{X}$ be $N_R \times N_T \times 3$ input data of the network with $c = 3$ channels. We define the first channel of the input as the absolute value of the imperfect channel matrix $\tilde{\mathbf{H}}$ as $[[\mathbf{X}]_{:,:,1}]_{i,j} = |[\tilde{\mathbf{H}}]_{i,j}|$ where $\tilde{\mathbf{H}} \sim \mathcal{CN}(\mathbf{H}, \mathbf{\Gamma})$. $\mathbf{\Gamma} \in \mathbb{R}^{N_R \times N_T}$ denotes the variance of AWGN with its $(i, j)$th entry as $\mathbf{\Gamma}_{i,j} = \frac{10^{\text{SNR}_{\text{TRAIN}}/20}}{[\mathbf{H}]_{i,j}}$ where $\text{SNR}_{\text{TRAIN}}$ denotes the SNR for the AWGN in the training process. Similarly, the second and the third channels are defined as the real and imaginary parts of $\tilde{\mathbf{H}}$, i.e., $[[\mathbf{X}]_{:,:,2}]_{i,j} = \text{Re}\{[\tilde{\mathbf{H}}]_{i,j}\}$ and $[[\mathbf{X}]_{:,:,3}]_{i,j} = \text{Im}\{[\tilde{\mathbf{H}}]_{i,j}\}$. We generate $NL$ realizations of the channel matrix where $N$ different channel matrices are generated with different user locations and channel gains. Each channel matrix is generated for $L_c$ different

---

**Algorithm 1** Training data generation for $\text{CNN}_{\text{AS}}$ and $\text{CNN}_{\text{RF}}$.

**Input:** $L_c, L_n, N, N_T, N_R, N_{RS}, \text{SNR}_{\text{TRAIN}}$.
**Output:** Training data $\mathcal{D}_{\text{AS}}$ and $\mathcal{D}_{\text{RF}}$.
1: Initialize with $t = 1$ for $t = 1, \ldots, T = NL$ where $L = L_c L_n$.
2: **for** $1 \leq n \leq N$ **do**
3:     **for** $1 \leq l_c \leq L_c$ **do**
4:         Generate $\mathbf{H}^{(l_c,n)}$ with $N_c = \mathbb{N}_{\text{cluster}}^{(l_c)}$.
5:         **for** $1 \leq l \leq L_n$ **do**
6:             $[\mathbf{H}^{(t)}]_{i,j} \sim \mathcal{CN}([\mathbf{H}^{(l_c,n)}]_{i,j}, \sigma_{\text{TRAIN}}^2)$.
7:             **for** $1 \leq q_A \leq Q_A$ **do**
8:                 $\mathbf{H}_{q_A}^{(t)} = \mathbf{U}_{q_A}^{(t)}\mathbf{\Sigma}_{q_A}^{(t)}\mathbf{V}_{q_A}^{(t)^H}$.
9:                 $\mathbf{F}_{q_A}^{\text{opt}(t)} = \mathbf{V}_{q_A}^{(1)^{(t)}}$.
10:               $\mathbf{W}_{q_A}^{\text{opt}(t)} = \left(\frac{1}{\rho}\left(\mathbf{F}_{q_A}^{\text{opt}(t)^H}\mathbf{H}_{q_A}^{(t)^H}\mathbf{H}_{q_A}^{(t)}\mathbf{F}_{q_A}^{\text{opt}(t)}\right.\right.$
                      $\left.\left. + \frac{N_S\sigma_n^2}{\rho}\mathbf{I}_{N_S}\right)^{-1}\mathbf{F}_{q_A}^{\text{opt}(t)^H}\mathbf{H}_{q_A}^{(t)^H}\right)^H$.
11:               Use $\mathbf{F}_{q_A}^{\text{opt}(t)}$, and find $\mathbf{F}_{RF_{q_A}}^{(t)}, \mathbf{F}_{BB_{q_A}}^{(t)}$ by solving (15).
12:               Use $\mathbf{W}_{q_A}^{\text{opt}(t)}$, and find $\mathbf{W}_{RF_{q_A}}^{(t)}, \mathbf{W}_{BB_{q_A}}^{(t)}$ in (16).
13:               Compute $R_A^{(t)}(q_A)$ in (9).
14:             **end for**
15:             $\mathbf{H}_{\bar{q}_A}^{(t)} \leftarrow \bar{q}_A^{(t)} = \arg\max_{q_A} R_A^{(t)}(q_A)$.
16:             $\mathbf{F}_{RF}^{(t)} \leftarrow \mathbf{F}_{RF_{\bar{q}_A}}^{(t)}, \mathbf{F}_{BB}^{(t)} \leftarrow \mathbf{F}_{BB_{\bar{q}_A}}^{(t)}$.
17:             $\mathbf{W}_{RF}^{(t)} \leftarrow \mathbf{W}_{RF_{\bar{q}_A}}^{(t)}, \mathbf{W}_{BB}^{(t)} \leftarrow \mathbf{W}_{BB_{\bar{q}_A}}^{(t)}$.
18:             $[[\mathbf{X}^{(t)}]_{:,:,1}]_{i,j} = |[\mathbf{H}^{(t)}]_{i,j}|$.
19:             $[[\mathbf{X}^{(t)}]_{:,:,2}]_{i,j} = \text{Re}\{[\mathbf{H}^{(t)}]_{i,j}\}$.
20:             $[[\mathbf{X}^{(t)}]_{:,:,3}]_{i,j} = \text{Im}\{[\mathbf{H}^{(t)}]_{i,j}\} \; \forall ij$.
21:             $\mathbf{z}^{(t)} = [\text{vec}^T\{\angle\mathbf{F}_{RF}^{(t)}\}, \text{Re}\{\text{vec}^T\{\mathbf{F}_{BB}^{(t)}\}\}, \text{Im}\{\text{vec}^T\{\mathbf{F}_{BB}^{(t)}\}\},$
              $\text{vec}^T\{\angle\mathbf{W}_{RF}^{(t)}\}, \text{Re}\{\text{vec}^T\{\mathbf{W}_{BB}^{(t)}\}\}, \text{Im}\{\text{vec}^T\{\mathbf{W}_{BB}^{(t)}\}\}]^T$.
22:             Construct the input-output pair $(\mathbf{X}^{(t)}, \bar{q}_A^{(t)})$ for $\text{CNN}_{\text{AS}}$ and $(\mathbf{X}_{\bar{q}_A}^{(t)}, \mathbf{z}^{(t)})$ for $\text{CNN}_{\text{RF}}$.
23:         $t \leftarrow t + 1$.
24:         **end for**
25:     **end for**
26: **end for**
27: Training data for $\text{CNN}_{\text{AS}}$ and $\text{CNN}_{\text{RF}}$ is obtained from the collection of the input-output pairs as
$$\mathcal{D}_{\text{AS}} = ((\mathbf{X}^{(1,1)}, \bar{q}_A^{(1,1)}), \ldots, (\mathbf{X}^{(T)}, \bar{q}_A^{(T)})),$$
$$\mathcal{D}_{\text{RF}} = ((\mathbf{X}_{\bar{q}_A}^{(1,1)}, \mathbf{z}^{(1,1)}), \ldots, (\mathbf{X}_{\bar{q}_A}^{(T)}, \mathbf{z}^{(T)})).$$

---

number of cluster to enrich the training data. In addition, the channel matrix is corrupted by synthetic noise for $L_n$ realizations where the element-wise noise is defined by $\text{SNR}_{\text{TRAIN}}$. Hence, the total size of the training input data is $N_R \times N_T \times 3 \times NL$ where $L = L_c L_n$. For each generated channel matrix, say $\mathbf{H}^{(n)}$, the best antenna subarrays with positions $\mathbb{A}_{\bar{q}_A}$ and the best RF beamformers $\mathbf{F}_{RF}, \mathbf{W}_{RF}$ are obtained by solving (11), (15) and (16) offline. This gives input-output pairs of the training data. The training process of both CNNs is identical except that they have different input dimensions. Algorithm 1 summarizes the steps of training data generation.

The $\text{CNN}_{\text{AS}}$ accepts the input of size $N_R \times N_T \times 3$ with labels $q_A$ whereas the input of $\text{CNN}_{\text{RF}}$ is of size $N_{RS} \times N_T \times 3$ with label $\mathbf{z}$. For each CNN, the network is composed of 14 layers. The first layer is the input layer with appropriate size. The second, fourth and the sixth layers are convolutional layers with 64 filters of size $2 \times 2$. The eight and eleventh layers are fully connected layers with 512 units. The tenth and thirteenth layers are dropout layers with 50% probability placed after each fully connected layers. There are ReLU (Rectified Linear Unit) layers after each convolutional and fully connected layers where $\text{ReLU}(x) = \max(x, 0)$. The final layer of $\text{CNN}_{\text{AS}}$ is the classification layer with size $\bar{Q}_A$ which is the number of subarrays that yield maximum spectral efficiency. In the classification layer, a softmax function is used to obtain the probability distribution of the classes. The output layer of $\text{CNN}_{\text{RF}}$ is a regression layer of

---

[2]The proposed antenna and RF chain selection framework can also be applied to the uplink case where the received signal model is $\bar{\mathbf{y}}^{\text{UL}} = \sqrt{\rho}\mathbf{F}_{BB}^{\text{Full}^H}\mathbf{F}_{RF}^{\text{Full}^H}\mathbf{H}^{\text{UL}}\mathbf{W}_{RF}\mathbf{W}_{BB}\mathbf{s} + \mathbf{F}_{BB}^{\text{Full}^H}\mathbf{F}_{RF}^{\text{Full}^H}\mathbf{n}$, which is obtained by switching the precoders $\mathbf{F}_{RF}, \mathbf{F}_{BB}$ and combiners $\mathbf{W}_{RF}, \mathbf{W}_{BB}$ in (4) and the uplink channel matrix is represented by the $N_T \times N_R$ matrix $\mathbf{H}^{\text{UL}} = \mathbf{H}^T$ [11].
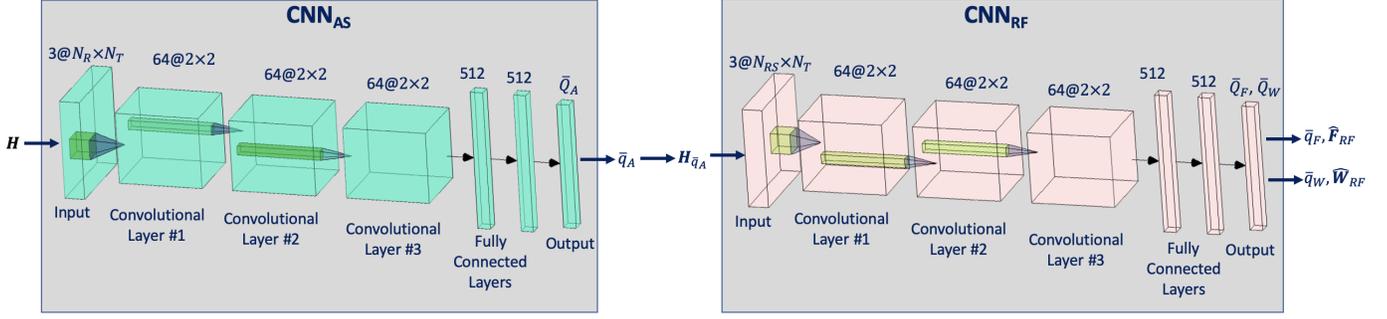
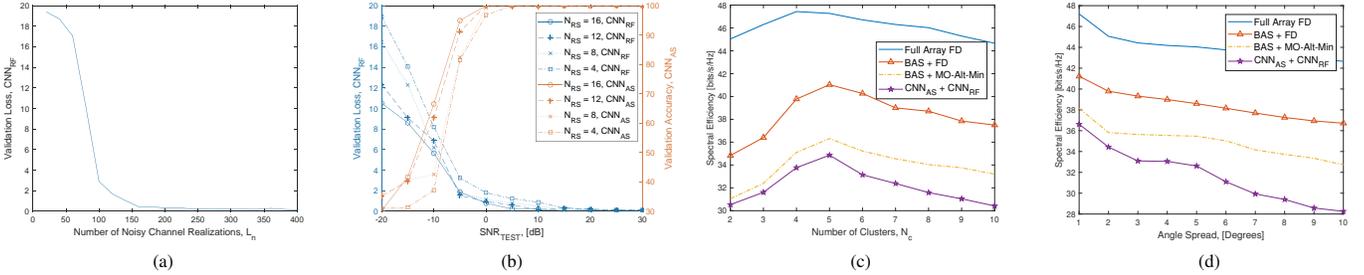Fig. 2. The proposed CNN architecture for antenna selection and RF beamformer design.



Fig. 3. Performance of CNN. (a) Validation loss versus number of channel realizations $L_n$. (b) Validation loss for $CNN_{RF}$ and validation accuracy for $CNN_{AS}$. (c) Spectral efficiency versus number of clusters in when $\sigma_\Theta^2 = 5°$. (d) Spectral efficiency versus angle spread when $N_c = 4$. We set $N_T = N_R = 64$ and $N_S = 4$ for all figures and $N_{RS} = 16$ in (a), (c) and (d).

size $G \times 1$. The proposed network is realized in MATLAB on a PC with 768-core GPU. The proposed network architecture is obtained through an optimization analysis to achieve the best performance and less computational cost.

The use of neural networks in mobile devices is also another issue for practical considerations since a deep neural network is composed of huge number of weights for each layers such as convolutional and fully connected layers [43]. The proposed CNN framework should be applied to the channel matrix in a mobile device to obtain hybrid beamformers where the resolution of the network parameters is of great importance. In common CNN architectures, convolutional layers have relatively less parameter to be optimized as compared to the fully connected layers which have large number of neurons to be updated in each iteration [57]. In order to obtain low resolution CNN structure, each of the weights and bias of all layers are quantized so that the saving in the memory is achieved and the implementation of the CNN is eased.

To train the proposed CNN structure, $N = 100$ different realizations of the channel matrix which is generated for $L = 800$ ($L_c = 4$ for $N_c \in \mathbb{N}_{cluster} = \{3, 4, 5, 6\}$ and $L_n = 200$) noisy realizations with three noise levels, i.e., $SNR_{TRAIN} \in \{15, 20, 25\}$dB. Hence the total size of the training data is $N_R \times N_T \times 3 \times 240000$. In the training process, 70% and 30% of all data generated are selected as the training and validation datasets, respectively. Validation aids in hyperparameter tuning during the training phase to avoid the network simply memorizing the training data rather than learning general features for accurate prediction with new data. The validation data is used to test the performance of the network in the simulations for $J_T = 100$ Monte Carlo trials.

## V. NUMERICAL EXPERIMENTS

We evaluated the performance of the proposed CNN approach via several experiments. In order to prevent the similarity between the

test data and the training data we also add synthetic noise to the test data where the SNR in testing is defined similar to $SNR_{TRAIN}$ as $SNR_{TEST} = 20 \log_{10}(\frac{|[\mathbf{H}]_{i,j}|^2}{\mathbf{\Gamma}_{i,j}})$. We used the stochastic gradient descent algorithm with momentum [58] for updating the network parameters with learning rate 0.01 and mini-batch size of 500 samples for 50 epochs. As a loss function, we use the negative log-likelihood or cross-entropy loss [24]. We select $N_T^{RF} = N_R^{RF} = 4$ for all simulations. For each channel matrix realization, the propagation environment is modeled with $N_c = 4$ clusters and $N_{ray} = 5$ rays for each clusters with the angle spread of $\sigma_\Theta^2 = 5°$ for all transmit and receive azimuth and elevation angles which are uniform randomly selected from the interval $[-60°, 60°]$ and $[-20°, 20°]$ respectively.

### A. Performance of Unquantized CNN

Figure 3 summarizes our assessment of the performance of unquantized CNN. Here, we set $N = 100$, $N_T = N_R = 64$ and $N_S = 4$. Figure 3a shows the validation loss of $CNN_{RF}$ against the number of noisy channel realizations $L$. We observe that the loss is satisfactory for $L \geq 150$; in our simulations, we keep $L = 200$. Note that this is a common result for different number of channel realizations $N$. In fact, we use $N = 100$ in our settings to achieve reasonable network accuracy. To investigate the performance of deep networks against different noise levels in the training data, we demonstrate the validation loss of $CNN_{RF}$ and the classification accuracy of $CNN_{AS}$ in Fig. 3b for different $N_{RS}$ values. It is clear here that both networks attain satisfactory network accuracy for $SNR_{TEST} \geq 0$dB. At low SNR regimes, CNN has poor classification performance due to the deviations between the input and the channel matrices used in the training data. In order to make CNN more robust to noisy inputs, we draw the training data for multiple $SNR_{TRAIN}$ levels. Nevertheless, noise in the training data expectedly limits the performance since the network cannot distinguish the input data if it is corrupted too much. This issue is also reported in [39], [59] for multiple $SNR_{TRAIN}$ case.

The channel statistics are important parameters that change in very short time in mm-Wave channels. Hence, in the training stage, we feed the network with channel realizations of different $N_c$ values. To further investigate the performance with respect to different channel statistics, we compare the algorithms in Fig. 3c and Fig. 3d for different $N_c$ and $\sigma_\Theta^2$, respectively. The proposed CNN framework provides robust performance against different channel statistics. The effectiveness of the proposed techniques can be attributed to training the network with several channel statistics and adding synthetic noise for multiple $\text{SNR}_{\text{TRAIN}}$ values.

As a result, the proposed network provides robust performance against the changes in channel statistics without a need to be re-trained. However, when there is a change in the *full array* system parameters such as $N_T$, $N_R$, $N_{RS}$ and $N_S$, the network does need to be re-trained because these parameters directly dictates the dimensions of the network input and output layers.

### B. Antenna Selection Performance

In this experiment, we present the antenna selection performance for $N_T = N_R = 256$ and $N_{RS} = 16$, and the results are given in Fig. 4a. For fair comparison, the hybrid beamforming is performed by the MO algorithm for all the antenna selection techniques, which are best antenna selection (BAS) with BAB algorithm [23], $\text{CNN}_{\text{AS}}$, Greedy-based antenna selection (GAS) [53], Scheme 1, which is random antenna selection (RAS) and Scheme 2 from Fig. 1. Note that BAS, $\text{CNN}_{\text{AS}}$ and GAS aim to optimize the selected antennas and they refer to the selection Scheme 3 demonstrated in Fig. 1. The performance of the algorithms is also compared with the full array performance ($N_R = N_{RS} = 256$) which is shown for both fully-digital (FD) and hybrid beamforming (HB). Scheme 1 (RAS) and Scheme 2 have much simpler architectures which do not include optimization and we can see that Scheme 2 performs the worst since it has the simplest architecture, i.e., selecting the antennas through the absolute values of the entries of the channel matrix. We observe that the performance $\text{CNN}_{\text{AS}}$ is close to BAS, which employs BAB algorithm to yield an optimum solution. Other algorithms/schemes are suboptimal and $\text{CNN}_{\text{AS}}$ outperforms them. The performance of $\text{CNN}_{\text{AS}}$ is superior in comparison with the other algorithms/schemes which provide sub-optimum performance. We also present the performance for different number of selected antennas as shown in Fig. 4b for the same settings and SNR= 10 dB. From this figure, we obtain similar observations as in Fig. 1, demonstrating the outperformance of $\text{CNN}_{\text{AS}}$.

### C. Hybrid Beamforming Performance

In this experiment, the spectral efficiency of our CNN-based hybrid beamforming approach is evaluated by comparison with MO-Alt-Min [15], SOMP [14] as well as the methods in [9] (called *Wang et.al*) and [41] (called *Sohrabi et.al*). In addition, we compare our work with the multilayer perceptron (MLP) approach proposed in [37] where the MLP architecture is fed and trained with the same training data of $\text{CNN}_{\text{RF}}$. The number of antennas are $N_R = N_T = 256$ and $N_{RS} = 16$ antennas are selected. For antenna selection, BAS is used for all the hybrid beamforming algorithms except $\text{CNN}_{\text{RF}}$, which is fed by the selected channel matrix obtained from $\text{CNN}_{\text{AS}}$. In Fig. 5, the spectral efficiency is presented for $N_S = 1$ (a) and $N_S = 4$ (b) respectively. We can see that $\text{CNN}_{\text{RF}}$ outperforms the other algorithms and it provides very close performance to MO-Alt-Min which is also used in the labeling process. The performance of $\text{CNN}_{\text{RF}}$ is attributed to the extracting features in the input data and matching the data to labels which are obtained through an exhaustive search algorithm.
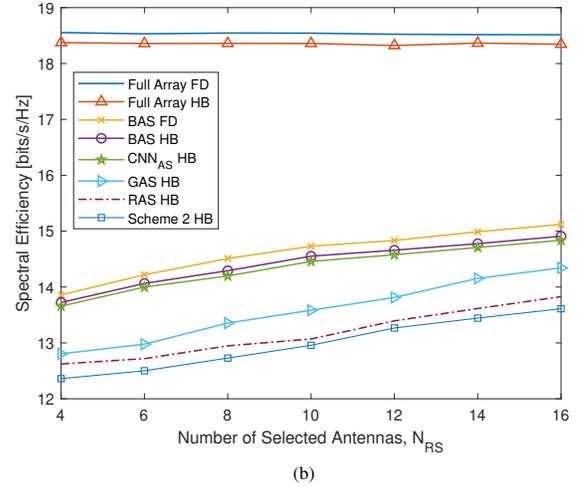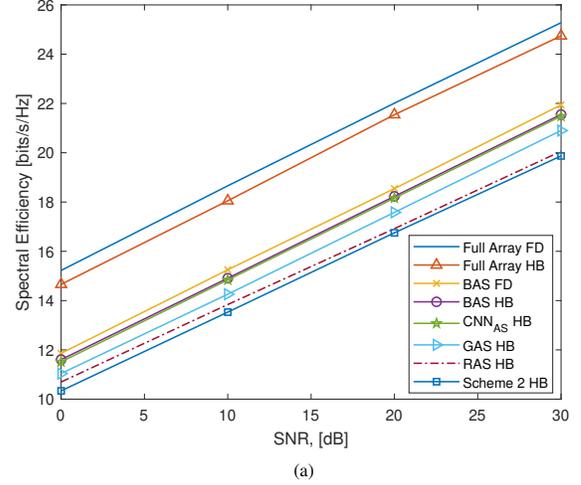


(a)



(b)

Fig. 4. Antenna selection performance for $N_T = N_R = 256$, $N_S = 1$: (a) spectral efficiency vs SNR for fixed $N_{RS} = 16$ (b) spectral efficiency versus $N_{RS}$.

Most of the hybrid beamforming algorithms rely on the perfectness of the channel matrix which is also used in the design process of beamformers. In order to relax this condition, we use a DL approach and train the deep network with several noisy channel data so that it can obtain robust performance due to the changes in the channel data. In Fig. 6, the robustness analysis is conducted for imperfect channel matrix with respect to $\text{SNR}_{\text{TEST}}$. The same simulations settings are used with SNR $= 10$dB and all algorithms are fed with the imperfect channel matrix. As it is seen, proposed CNN approach performs more robust performance as compared to the other algorithms. The advantage of CNN is that it is trained with several noisy channel data so that it can distinguish the labels for noisy inputs.

We further examine the algorithms on the performance of RF beamformer design. In this respect, we define a cost function to measure how close the algorithms are, to the unconstrained beamformers. We define the error between the unconstrained beamformers $\mathbf{F}^{\text{opt}}$, $\mathbf{W}^{\text{opt}}$ and the estimated hybrid beamformers $\hat{\mathbf{F}}_{RF}\hat{\mathbf{F}}_{BB}$, $\hat{\mathbf{W}}_{RF}\hat{\mathbf{W}}_{BB}$ as follows

$$\gamma_F = ||\mathbf{F}^{\text{opt}} - \hat{\mathbf{F}}_{RF}\hat{\mathbf{F}}_{BB}||_F/(N_T N_S), \qquad (18)$$

$$\gamma_W = ||\mathbf{W}^{\text{opt}} - \hat{\mathbf{W}}_{RF}\hat{\mathbf{W}}_{BB}||_{\mathcal{F}}/(N_{RS} N_S). \qquad (19)$$

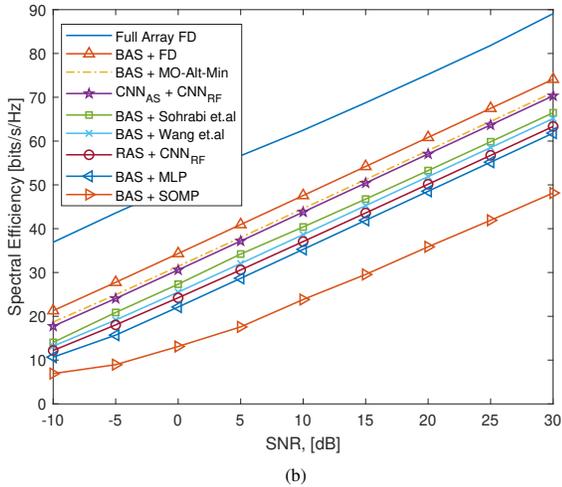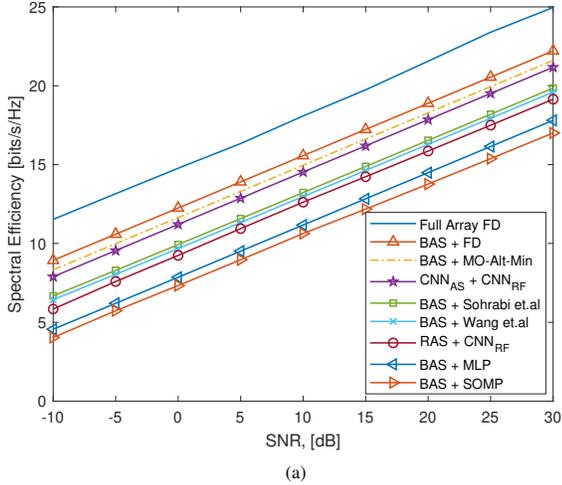Then we present the results in Table I for $N_S = \{1, 2, 3\}$ and

(a)



(b)

Fig. 5. Spectral efficiency of various hybrid beamforming algorithms versus SNR for (a) $N_S = 1$ and (b) $N_S = 4$.



Fig. 6. Performance comparison for corrupted channel data.

TABLE I
PERFORMANCE LOSS FOR HYBRID BEAMFORMER DESIGN.

| | | $\gamma_F$ | $\gamma_W$ |
|---|---|---|---|
| BAS + MO-Alt-min | $N_S = 1$ | 0.010 | 0.0016 |
| | $N_S = 2$ | 0.008 | 0.0014 |
| | $N_S = 3$ | 0.005 | 0.0015 |
| $CNN_{AS} + CNN_{RF}$ | $N_S = 1$ | 0.010 | 0.0060 |
| | $N_S = 2$ | 0.008 | 0.0043 |
| | $N_S = 3$ | 0.005 | 0.0022 |
| BAS + *Sohrabi et.al* | $N_S = 1$ | 0.032 | 0.0062 |
| | $N_S = 2$ | 0.019 | 0.0034 |
| | $N_S = 3$ | 0.008 | 0.0024 |
| BAS + *Wang et.al* | $N_S = 1$ | 0.043 | 0.0069 |
| | $N_S = 2$ | 0.027 | 0.0048 |
| | $N_S = 3$ | 0.013 | 0.0028 |
| BAS + SOMP | $N_S = 1$ | 0.070 | 0.0037 |
| | $N_S = 2$ | 0.032 | 0.0074 |
| | $N_S = 3$ | 0.025 | 0.0165 |

$N_T = N_R = 256$, $N_{RS} = 16$. It can be seen that CNN-based RF chain selection provides less error for both precoder and combiner design as compared to other algorithm. The performance of the CNN is attributed to aiming the highest spectral efficiency with the selection of the best configuration of beamformers. Note also that the algorithms including our CNN approach achieve less $\gamma_F$ and $\gamma_W$ when $N_S = 3$ as compared to $N_S = 1$ except SOMP. In contrast, SOMP performs worse (less $\gamma_F$ and relatively larger $\gamma_W$) when $N_S$ is larger. This explains the performance loss of SOMP observed in Fig. 5b as $N_S$ increases.

### D. Binarized and Quantized CNN

While considering larger CNN that has more layers and nodes, the associated memory and computational cost could be prohibitive for massive MIMO. This hinders deploying large CNNs in mobile devices, which have limited memory and restricted latency to perform tasks such as online learning and incremental learning. In this context, compressing a deep neural network, that has attracted a lot of attention recently [60], is highly desirable. In this paper, we adopt a network quantization to compress and thereby accelerate the CNN. This method compresses the original network by reducing the number of bits required to represent each 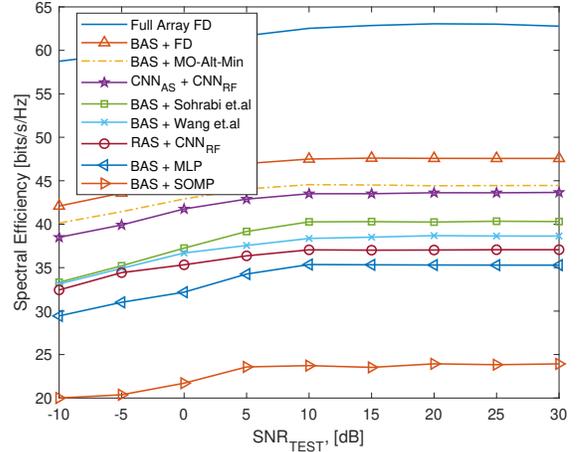weight of convolutional and fully connected layers. It has been observed that this compress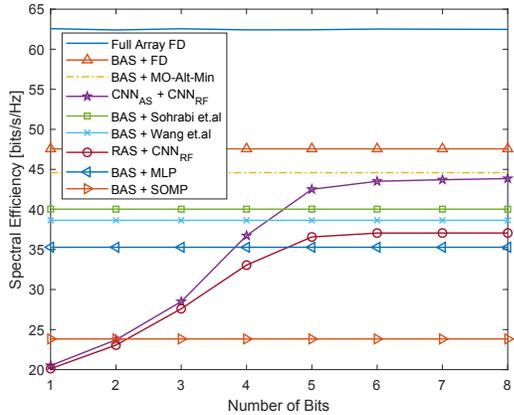ion can support both pretrained and trained-from-scratch models, helps in significantly reducing memory usage and speeds up the computations.

Network quantization in large CNNs can significantly degrade the classification accuracy. We, therefore, investigate the minimum number of bits required to store our proposed twin-CNN network for an acceptable spectral efficiency. Improving the performance of quantized or binarized CNN is an active research area. For example, in the extreme case of binarized or 1-bit CNN, it has been shown that networks trained with back propagation could be resilient to weight distortions introduced by binarization [60].

We examined the performance of the CNN with quantized weights and biases. In Fig. 7a, we present the performance of binarized-CNN where the parameters of the network are either 0 or 1. CNN performance is poorer than SOMP in case of hybrid beamforming. To further investigate the effect of quantization of network parameters, we demonstrate, in Fig. 7b, the performance of CNN versus number of bits used to quantize the weights and biases of all layers of the CNN. As it is seen, at least 5 bits are required for CNN to attain the best subarray and best RF chain performance.

Fig. 7. Performance of quantized-CNN. Spectral efficiency versus SNR given in (a) and versus number of bits in (b).

TABLE II
COMPUTATIONAL COMPLEXITY (IN SECONDS)

| $N_{RS} = 8$ | | | | | |
|---|---|---|---|---|---|
| $N_T$ | CNN | MO-AltMin | Sohrabi et.al | Wang et.al | SOMP |
| 32 | 0.003 | 0.241 | 0.005 | 0.014 | 0.028 |
| 128 | 0.004 | 0.632 | 0.014 | 0.073 | 0.147 |
| 256 | 0.005 | 1.324 | 0.041 | 0.151 | 0.195 |
| $N_{RS} = 16$ | | | | | |
| $N_T$ | CNN | MO-AltMin | Sohrabi et.al | Wang et.al | SOMP |
| 32 | 0.004 | 0.284 | 0.008 | 0.021 | 0.031 |
| 128 | 0.010 | 1.171 | 0.027 | 0.084 | 0.112 |
| 256 | 0.013 | 4.458 | 0.050 | 0.172 | 0.219 |

*E. Computational Complexity*

In this experiment, we measure the computation time of our CNN approach and compare it with other state-of-the-art algorithms. We select $N_R = 64$ and $N_S = 4$. The results are given in Table II with respect to the number of BS antennas $N_T$ when $N_{RS} = \{8, 16\}$. As it is seen, our CNN approach enjoys less computation time where the complexity is due to the classification of the input data. The remaining algorithms have relatively high run times and MO-Alt-Min has the highest computation cost however it has better than the remaining algorithms.

For the sake of completeness, we also calculate the total computation time for the generation of the training data where we select $N_T = N_R = 256$, $N_{RS} = 16$, $L = N = 100$ and use three

SNR$_{\text{TRAIN}}$ levels. In this settings, it takes about two days to generate $256 \times 256 \times 3 \times 30000$ training data. When $N_T = N_R = 25$, $N_{RS} = 16$, it takes only 40 minutes for $25 \times 25 \times 3 \times 30000$. In training, the main challenge is not the time but the memory considerations where large antenna arrays yield higher variables of size $Q_A$ which require large memory allocations to save (even temporarily) the results.

## VI. SUMMARY

We proposed a twin-CNN deep learning approach for joint antenna selection and hybrid beamformer design in mm-Wave communications. Our CNN framework provides significant improvement in the capacity as compared to the conventional beamformer design techniques. This method does not require the precise knowledge of the channel matrix and has significantly better performance than the conventional techniques used in mm-Wave MIMO systems. Instead of computing analog and baseband beamformers, the proposed approach only requires the estimated channel matrix to feed the network and yields the best antenna subarray and analog and baseband beamformers. Hence, it has very low computational complexity. We also investigated the quantized-CNN model when it needs to be applied in a low-memory, low-overhead platform such as a mobile phone. We show that no more than 5 bits are required to save (or access in a cloud-based environment) the CNN in digital form.

## REFERENCES

[1] K. V. Mishra, M. R. Bhavani Shankar, V. Koivunen, B. Ottersten, and S. A. Vorobyov, "Toward millimeter wave joint radar-communications: A signal processing perspective," *IEEE Signal Processing Magazine*, vol. 36, no. 5, pp. 100–114, 2019.

[2] S. H. Dokhanchi, B. S. Mysore, K. V. Mishra, and B. Ottersten, "A mmWave automotive joint radar-communications system," *IEEE Transactions on Aerospace and Electronic Systems*, 2019. in press.

[3] K. V. Mishra, S. S. Ram, S. Vishwakarma, and G. Duggal, "Doppler-resilient 802.11ad-based ultra-short range automotive radar," *arXiv preprint arXiv:1902.01306*, 2019.

[4] J. A. Hodge, K. V. Mishra, and A. I. Zaghloul, "Reconfigurable metasurfaces for index modulation in 5G wireless communications," in *IEEE International Applied Computational Electromagnetics Society Symposium*, pp. 1–2, 2019.

[5] A. Ayyar and K. V. Mishra, "Robust communications-centric coexistence for turbo-coded OFDM with non-traditional radar interference models," in *IEEE Radar Conference*, 2019. in press.

[6] 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, 2013.

[7] M. Alaee-Kerahroodi, K. V. Mishra, M. R. B. Shankar, and B. Ottersten, "Discrete phase sequence design for coexistence of MIMO radar and MIMO communications," in *IEEE International Workshop on Signal Processing Advances in Wireless Communications*, pp. 1–5, 2019.

[8] K. V. Mishra and Y. C. Eldar, "Sub-Nyquist radar: Principles and prototypes," *arXiv preprint arXiv:1803.01819*, 2018.

[9] Z. Wang, M. Li, Q. Liu, and A. L. Swindlehurst, "Hybrid Precoder and Combiner Design With Low-Resolution Phase Shifters in mmWave MIMO Systems," *IEEE J. Sel. Topics Signal Process.*, vol. 12, pp. 256–269, May 2018.

[10] J. B. Tsui, *Digital techniques for wideband receivers*, vol. 2. SciTech Publishing, 2004.

[11] A. Alkhateeb, O. E. Ayach, G. Leus, and R. W. Heath, "Hybrid precoding for millimeter wave cellular systems with partial channel knowledge," in *Information Theory and Applications Workshop*, pp. 1–5, 2013.

[12] A. Alkhateeb, O. E. 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, 2014.

[13] A. Alkhateeb, G. Leus, and R. W. Heath, "Limited Feedback Hybrid Precoding for Multi-User Millimeter Wave Systems," *IEEE Transactions on Wireless Communications*, vol. 14, no. 11, pp. 6481–6494, 2015.

[14] O. E. Ayach, S. Rajagopal, S. Abu-Surra, Z. Pi, and R. W. Heath, "Spatially sparse precoding in millimeter wave MIMO systems," *IEEE Transactions on Wireless Communications*, vol. 13, no. 3, pp. 1499–1513, 2014.

[15] X. Yu, J. Shen, J. Zhang, and K. B. Letaief, "Alternating minimization algorithms for hybrid precoding in millimeter wave MIMO systems," *IEEE Journal of Selected Topics in Signal Processing*, vol. 10, no. 3, pp. 485–500, 2016.

[16] K. V. Mishra, Y. C. Eldar, E. Shoshan, M. Namer, and M. Meltsin, "A cognitive sub-Nyquist MIMO radar prototype," *arXiv preprint arXiv:1807.09126*, 2018.

[17] S. Na, K. V. Mishra, Y. Liu, Y. C. Eldar, and X. Wang, "TenDSuR: Tensor-based 3D sub-Nyquist radar," *IEEE Signal Processing Letters*, vol. 26, no. 2, pp. 237–241, 2019.

[18] O. T. Demir and T. E. Tuncer, "Antenna selection and hybrid beamforming for simultaneous wireless information and power transfer in multi-group multicasting systems," *IEEE Transactions on Wireless Communications*, vol. 15, no. 10, pp. 6948–6962, 2016.

[19] X. Zhai, Q. Shi, Y. Cai, and M. Zhao, "Joint transmit precoding and receive antenna selection for uplink multiuser massive MIMO systems," *IEEE Transactions on Communications*, vol. 66, no. 11, pp. 5249–5260, 2018.

[20] X. Zhai, Y. Cai, Q. Shi, M. Zhao, G. Y. Li, and B. Champagne, "Joint transceiver design with antenna selection for large-scale MU-MIMO mmWave systems," *IEEE Journal on Selected Areas in Communications*, vol. 35, no. 9, pp. 2085–2096, 2017.

[21] H. Li, Q. Liu, Z. Wang, and M. Li, "Joint antenna selection and analog precoder design with low-resolution phase shifters," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 1, pp. 967–971, 2018.

[22] V. V. Ratnam, A. F. Molisch, O. Y. Bursalioglu, and H. C. Papadopoulos, "Hybrid Beamforming With Selection for Multiuser Massive MIMO Systems," *IEEE Transactions on Signal Processing*, vol. 66, no. 15, pp. 4105–4120, 2018.

[23] Y. Gao, H. Vinck, and T. Kaiser, "Massive MIMO antenna selection: Switching architectures, capacity bounds, and optimal antenna selection algorithms," *IEEE Transactions on Signal Processing*, vol. 66, no. 5, pp. 1346–1360, 2018.

[24] Y. Lecun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.

[25] K. V. Mishra, A. Gharanjik, M. R. B. Shankar, and B. Ottersten, "Deep learning framework for precipitation retrievals from communication satellites," in *Euro. Conf. Radar Met. Hydro.*, p. 023, 2018.

[26] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 1798–1828, 2013.

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

[28] Y. Long, Z. Chen, J. Fang, and C. Tellambura, "Data-driven-based analog beam selection for hybrid beamforming under mm-Wave channels," *IEEE Journal of Selected Topics in Signal Processing*, vol. 12, no. 2, pp. 340–352, 2018.

[29] J. A. Hodge, K. V. Mishra, and A. I. Zaghloul, "Joint multi-layer GAN-based design of tensorial RF metasurfaces," in *IEEE International Workshop on Machine Learning for Signal Processing*, 2019. in press.

[30] J. A. Hodge, K. V. Mishra, and A. I. Zaghloul, "RF metasurface array design using deep convolutional generative adversarial networks," in *IEEE International Symposium on Phased Array Systems and Technology*, 2019. in press.

[31] J. A. Hodge, K. V. Mishra, and A. I. Zaghloul, "Multi-discriminator distributed generative model for multi-layer RF metasurface discovery," in *IEEE Global Conference on Signal and Information Processing*, 2019. in press.

[32] P. Zhou, X. Fang, X. Wang, Y. Long, R. He, and X. Han, "Deep Learning-Based Beam Management and Interference Coordination in Dense mmWave Networks," *IEEE Transactions on Vehicular Technology*, vol. 68, pp. 592–603, Jan 2019.

[33] N. Samuel, T. Diskin, and A. Wiesel, "Deep MIMO detection," in *IEEE International Workshop on Signal Processing Advances in Wireless Communications*, pp. 1–5, 2017.

[34] 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, 2018.

[35] V. Raj and S. Kalyani, "Backpropagating through the air: Deep learning at physical layer without channel models," *IEEE Communications Letters*, vol. 22, no. 11, pp. 2278–2281, 2018.

[36] C. Wen, W. Shih, and S. Jin, "Deep learning for massive MIMO CSI feedback," *IEEE Wireless Communications Letters*, vol. 7, no. 5, pp. 748–751, 2018.

[37] H. Huang, Y. Song, J. Yang, G. Gui, and F. Adachi, "Deep-Learning-based Millimeter-Wave Massive MIMO for Hybrid Precoding," *IEEE Trans. Veh. Technol.*, pp. 1–1, 2019.

[38] A. M. Elbir and K. V. Mishra, "Deep Learning Design for Joint Antenna Selection and Hybrid Beamforming in Massive MIMO," in *2019 IEEE International Symposium on Antennas and Propagation USNC/URSI National Radio Science Meeting*, July 2019.

[39] A. M. Elbir, K. V. Mishra, and Y. C. Eldar, "Cognitive radar antenna selection via deep learning," *IET Radar, Sonar & Navigation*, vol. 13, pp. 871–880(9), June 2019.

[40] N. Boumal, B. Mishra, P.-A. Absil, and R. Sepulchre, "Manopt, a Matlab Toolbox for Optimization on Manifolds," *Journal of Machine Learning Research*, vol. 15, pp. 1455–1459, 2014.

[41] F. Sohrabi and W. Yu, "Hybrid Analog and Digital Beamforming for mmWave OFDM Large-Scale Antenna Arrays," *IEEE J. Sel. Areas Commun.*, vol. 35, pp. 1432–1443, July 2017.

[42] A. M. Elbir and K. V. Mishra, "Robust hybrid beamforming with quantized deep neural networks," in *IEEE International Workshop on Machine Learning for Signal Processing*, 2019. in press.

[43] J. Wu, C. Leng, Y. Wang, Q. Hu, and J. Cheng, "Quantized convolutional neural networks for mobile devices," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4820–4828, 2016.

[44] E. Torkildson, C. Sheldon, U. Madhow, and M. Rodwell, "Millimeter-wave spatial multiplexing in an indoor environment," in *IEEE Globecom Workshops*, pp. 1–6, 2009.

[45] R. Méndez-Rial, C. Rusu, A. Alkhateeb, N. González-Prelcic, and R. W. Heath, "Channel estimation and hybrid combining for mmWave: Phase shifters or switches?," in *IEEE Information Theory and Applications Workshop*, pp. 90–97, 2015.

[46] R. Méndez-Rial, C. Rusu, N. González-Prelcic, A. Alkhateeb, and R. W. Heath, "Hybrid MIMO architectures for millimeter wave communications: Phase shifters or switches?," *IEEE Access*, vol. 4, pp. 247–267, 2016.

[47] Z. Marzi, D. Ramasamy, and U. Madhow, "Compressive Channel Estimation and Tracking for Large Arrays in mm-Wave Picocells," *IEEE J. Sel. Topics Signal Process.*, vol. 10, pp. 514–527, April 2016.

[48] J. Wang, Z. Lan, C.-W. Pyo, T. Baykas, C.-W. Sum, M. A. Rahman, J. Gao, R. Funada, F. Kojima, H. Harada, and S. Kato, "Beam codebook based beamforming protocol for multi-Gbps millimeter-wave WPAN systems," *IEEE Journal on Selected Areas in Communications*, vol. 27, no. 8, pp. 1390–1399, 2009.

[49] E. Björnson, L. Van der Perre, S. Buzzi, and E. G. Larsson, "Massive MIMO in Sub-6 GHz and mmWave: Physical, Practical, and Use-Case Differences," *arXiv e-prints*, p. arXiv:1803.11023, Mar 2018.

[50] W. U. Bajwa, J. Haupt, G. Raz, and R. Nowak, "Compressed channel sensing," in *IEEE Conference on Information Sciences and Systems*, pp. 5–10, 2008.

[51] D. Fan, F. Gao, Y. Liu, Y. Deng, G. Wang, Z. Zhong, and A. Nallanathan, "Angle Domain Channel Estimation in Hybrid Millimeter Wave Massive MIMO Systems," *IEEE Transactions on Wireless Communications*, vol. 17, no. 12, pp. 8165–8179, Dec 2018.

[52] S. Joshi and S. Boyd, "Sensor selection via convex optimization," *IEEE Transactions on Signal Processing*, vol. 57, no. 2, pp. 451–462, 2009.

[53] M. Gharavi-Alkhansari and A. B. Gershman, "Fast antenna subset selection in MIMO systems," *IEEE Transactions on Signal Processing*, vol. 52, pp. 339–347, Feb 2004.

[54] D. P. Palomar, J. M. Cioffi, and M. A. Lagunas, "Joint Tx-Rx beamforming design for multicarrier MIMO channels: A unified framework for convex optimization," *IEEE Transactions on Signal Processing*, vol. 51, no. 9, pp. 2381–2401, 2003.

[55] A. Goldsmith, S. A. Jafar, N. Jindal, and S. Vishwanath, "Capacity limits of MIMO channels," *IEEE Journal on Selected Areas in Communications*, vol. 21, no. 5, pp. 684–702, 2003.

[56] T. Kailath, B. Hassibi, and A. H. Sayed, *Linear estimation*. Prentice-Hall, 2000.

[57] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *CoRR*, vol. abs/1409.1556, 2015.

[58] C. M. Bishop, *Pattern recognition and machine learning*. Springer, 2006.

[59] A. M. Elbir, "CNN-Based Precoder and Combiner Design in mmWave MIMO Systems," *IEEE Commun. Lett.*, vol. 23, pp. 1240–1243, July 2019.

[60] Y. Cheng, D. Wang, P. Zhou, and T. Zhang, "A survey of model compression and acceleration for deep neural networks," *arXiv preprint arXiv:1710.09282*, 2017.