

Reinforcement Learning of Beam Codebooks in Millimeter Wave and Terahertz MIMO Systems

Yu Zhang, Muhammad Alrabeiah, and Ahmed Alkhateeb

Abstract

Millimeter wave (mmWave) and terahertz MIMO systems rely on pre-defined beamforming codebooks for both initial access and data transmission. Being pre-defined, however, those codebooks are commonly not optimized for specific environments, user distributions, and/or possible hardware impairments. This leads to large codebook sizes with high beam training overhead which increases the initial access/tracking latency and makes it hard for these systems to support highly mobile applications. To overcome these limitations, this paper develops a deep reinforcement learning framework that learns how to iteratively optimize the codebook beam patterns (shapes) relying only on the receive power measurements and without requiring any explicit channel knowledge. The developed model learns how to autonomously adapt the beam patterns to best match the surrounding environment, user distribution, hardware impairments, and array geometry. Further, this approach does not require any knowledge about the channel, array geometry, RF hardware, or user positions. To reduce the learning time, the proposed model designs a novel *Wolpertinger*-variant architecture that is capable of efficiently searching for an optimal policy in a large discrete action space, which is important for large antenna arrays with quantized phase shifters. This complex-valued neural network architecture design respects the practical RF hardware constraints such as the constant-modulus and quantized phase shifter constraints. Simulation results based on the publicly available DeepMIMO dataset confirm the ability of the developed framework to learn near-optimal beam patterns for both line-of-sight (LOS) and non-LOS scenarios and for arrays with hardware impairments without requiring any channel knowledge.

I. INTRODUCTION

Millimeter wave (mmWave) and terahertz (THz) MIMO systems adopt large antenna arrays to compensate for the significant path loss and ensure sufficient receive signal power. Because

Yu Zhang, Muhammad Alrabeiah and Ahmed Alkhateeb are with Arizona State University (Email: y.zhang, malrabei, alkhateeb@asu.edu). This work is supported by the National Science Foundation under Grant No. 1923676. Part of this work was presented at IEEE Asilomar [1].

of the high cost and power consumption of the mixed-circuit components, however, these systems normally rely either fully or partially on analog beamforming, where transceivers employ networks of phase shifters [2], [3]. This makes the basic MIMO signal processing functions, such as channel estimation and beamforming design, challenging as the channels are seen only through the RF lens. This motivates mmWave/THz massive MIMO systems to rely on pre-defined *beamforming codebooks* for both initial access and data transmission [4]–[6]. The classical pre-defined beamforming/beamsteering codebooks normally consist of a large number of single-lobe beams, each of which can steer the signal towards one direction. These classical codebooks, though, have several drawbacks: (i) To cover all the possible directions, these codebooks consist of a large number of beams, which makes the search over them associated with high beam training overhead. (ii) The second issue is a blowback from the directivity blessing; classical beamsteering codebooks employ single-lobe beams to maximize directivity, which, in many cases, may not be optimal, especially for Non-Line-of-Sight (NLOS) users. (iii) Further, the design of the classical codebooks normally assume that the array is calibrated and its geometry is known, which associates this design processing with high cost (due to the need for expensive calibration) and makes it hard to adapt to systems with unknown or arbitrary array geometries.

Another look at the aforementioned drawbacks reveals that they stem from the lack of environment and hardware adaptability. A mmWave/THz system that has a sense of its environment could discover the most frequent signal directions (for both single and multi-path cases) and accordingly tailor its codebook beam *patterns* (directions, shapes, number of lobes, etc.). Furthermore, the system can also overcome the challenges with intrinsic hardware impairments or unknown/arbitrary array geometries by learning how to calibrate its beams to adapt to the given hardware. All that awareness and adaptability can potentially be achieved if the mmWave/THz system incorporates a data-driven and artificially-intelligent component. Towards this goal, leveraging machine learning tools, especially reinforcement learning, is particularly promising. Reinforcement learning models can efficiently learn from the observed data and responses obtained from both the hardware and environment, which may potentially reduce the channel knowledge requirements. **Hence, the focus of this work is on developing a reinforcement learning based approach that learns how to adapt the codebook beams to the environment and hardware without requiring explicit channel knowledge.**

A. *Prior Work:*

Designing efficient beamforming and combining is essential for realizing the potential of MIMO communications, and it has been an important research topic in the literature of MIMO signal processing [2], [7]–[10]. For MIMO systems with no hardware constraints, i.e., with fully-digital processing and no constraints on the RF hardware, maximum ratio transmission and combining maximize the achievable SNR with single-stream transmission/reception [7]. To realize these solutions, however, the MIMO system should be able to control the magnitude and phase of the signal at each antennas. When only the phase can be controlled, equal-gain transmission solutions have been developed to maximize the SNR or diversity gains [8]. This is particularly interesting for mmWave and terahertz systems where the beamforming/precoding processing is fully or partially done in the RF domain using analog phase shifters [2]. In these systems, however, the phase shifters can normally take only quantized phase shift values. This associates the search over the large space of quantized phase shift values with high complexity (e.g., for a 32-element antenna array with 2-bit phase shifters, there are 4^{32} possible beamforming vectors) [2], [9], [10]. Further, in analog beamforming architectures, the channel is seen through the RF lens, which makes it hard to acquire at the baseband, especially for systems with arbitrary or unknown array geometries. To address these challenges, the first objective of this paper is to design a reinforcement learning based approach to efficiently learn the analog beamforming patterns that adapt to the surrounding environment and the adopted hardware/array geometry without requiring explicit channel knowledge.

Given the high complexity/training overhead associated with the beamforming design for quantized analog/hybrid architectures, these architectures normally rely on using pre-defined beam codebooks [3], [11]–[13]. These classical codebooks, however, are normally designed to have a large number of narrow beams; each beam has a single lobe and points to one direction. This has several limitations: (i) The large number of beams in these classical codebooks leads to large beam training overhead, which makes it hard for these mmWave/terahertz systems to support highly-mobile applications, (ii) single-lobe beams may not be optimal in scenarios with more than one dominant path such as NLOS situations, (iii) classical codebook design approaches normally rely on the knowledge of the channels, which is hard to achieve in architectures with analog beamforming, especially if these systems adopt imperfect hardware or deploy arrays with unknown/arbitrary array geometries, and (iv) the classical codebooks typically assume

fully-calibrated arrays, which is an expensive process. All the aforementioned limitations have motivated the search for adaptive codebook designing methods that leverage machine learning tools to adapt the learned beams based on the environment and hardware. In [14], the beam codebook learning approaches based on neural networks were proposed and shown to achieve good gain compared to the classical beam steering codebooks. The solutions in [14], though, still require full or partial channel knowledge, which is not simple to obtain in mmWave/THz systems. This motivates the development of environment and hardware aware codebook learning approach that does not require explicit channel knowledge, which is the main focus of this paper.

B. Contribution:

Developing environment and hardware awareness using machine learning is not straightforward when the mmWave system constraints are considered, e.g., channels are not available, phase shifters have finite and limited resolution, and hardware envelopes unknown impairments. In this paper, we develop a deep reinforcement learning based framework that can efficiently learn mmWave beam codebooks while addressing all these challenges. The main contributions of this paper can be summarized as follows:

- Designing a deep reinforcement learning based framework that can learn how to optimize the beam pattern for a set of users with similar channels. **The developed framework relies only on receive power measurements and does not require any channel knowledge.** This framework adapts the beam pattern based on the surrounding environment and learns how to compensate for the hardware impairments. This is done by utilizing a novel Wolpertinger architecture [15] which is designed to efficiently explore the large discrete action space. Further, the proposed model accounts for key hardware constraints such as the phase-only, constant-modulus, and quantized-angle constraints [2].
- Developing a reinforcement learning framework that is capable of learning a codebook of beam patterns optimized to serve the users in the surrounding environment. **The proposed framework autonomously optimizes the codebook beam patterns based on the environment, user distribution, hardware impairments, and array geometry.** Further, it relies only on the receive power measurements, does not require any position or channel knowledge (which relaxes the synchronization/coherence requirements), and does not require the users to be stationary during the learning process. This is achieved by developing a novel

pre-processing approach that relies on SNR-based feature matrices to partition/assign the users into clusters based on which parallel neural networks are trained.

- Extensively evaluating the performance of the proposed codebook learning approaches based on the publicly-available DeepMIMO dataset [16]. These experiments adopt both outdoor and indoor wireless communication scenarios and learn codebooks with different sizes. Further, this evaluation is done both for perfect uniform arrays and for arrays with arbitrary geometries and hardware impairments. These experiments provide a comprehensive evaluation of the proposed reinforcement learning based codebook learning approach.

The simulation results show that the proposed approach is capable of learning optimized beam patterns and beam codebooks without the need of providing any channel state information. Instead, based solely on the receive combining gains, the deep reinforcement learning solution adjusts the phases of the beamforming vectors to increase the receive gain and finally yield significant improvements over the classical beamsteering codebooks.

II. SYSTEM AND CHANNEL MODELS

In this section, we introduce in detail our adopted system and channel models. We also describe how the model considers arbitrary arrays with possible hardware impairments.

A. System Model

We consider the system model shown in Fig. 1 where a mmWave massive MIMO base station with M antennas is communicating with a single-antenna user. Further, given the high cost and power consumption of mixed-signal components [3], [17], we consider a practical system where the base station has only one radio frequency (RF) chain and employs analog-only beamforming using a network of r -bit quantized phase shifters. To facilitate the system operation and to respect the hardware constraints, mmWave and massive MIMO systems typically use beamforming codebooks in serving their users. Let \mathcal{W} denote the beam codebook adopted by the base station and assume that it contains N beamforming/combining vectors, with each one of them taking the form

$$\mathbf{w} = \frac{1}{\sqrt{M}} [e^{j\theta_1}, e^{j\theta_2}, \dots, e^{j\theta_M}]^T, \quad (1)$$

where each phase shift θ_m is selected from a finite set Θ with 2^r possible discrete values drawn uniformly from $(-\pi, \pi]$. In the uplink transmission, if a user u transmits a symbol $x \in \mathbb{C}$ to the

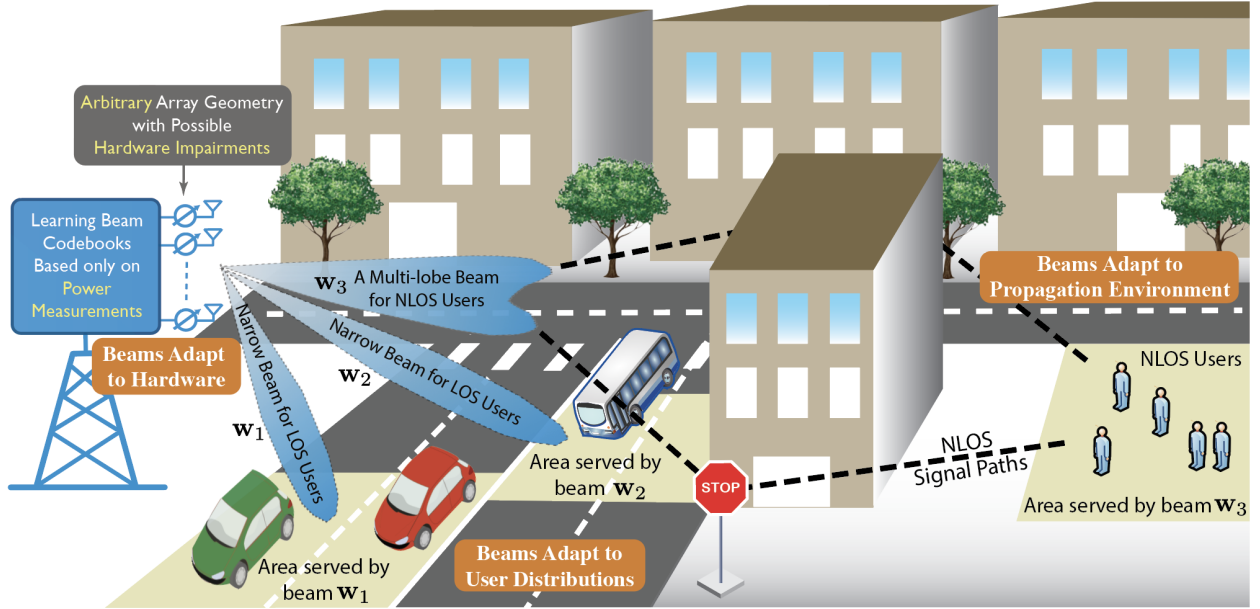


Fig. 1. A mmWave/massive MIMO base station with an arbitrary antenna array serving users with a beam codebook \mathcal{W} . The objective is to develop a learning approach for adapting the codebook \mathcal{W} to match the given hardware and environment based only on SNR measurements (which relaxes the coherence/synchronization requirements).

base station, where the transmitted symbol satisfies the average power constraint $\mathbb{E}[|x|^2] = P_x$, the received signal at the base station after combining can be expressed as

$$y_u = \mathbf{w}^H \mathbf{h}_u x + \mathbf{w}^H \mathbf{n}, \quad (2)$$

where $\mathbf{h}_u \in \mathbb{C}^{M \times 1}$ is the uplink channel vector between the user u and the base station antennas and $\mathbf{n} \sim \mathcal{N}_{\mathbb{C}}(0, \sigma_n^2 \mathbf{I})$ is the receive noise vector at the base station.

B. Channel Model

We adopt a general geometric channel model for \mathbf{h}_u [17], [18]. Assume that the signal propagation between the user u and the base station consists of L paths. Each path ℓ has a complex gain α_ℓ and an angle of arrival ϕ_ℓ . Then, the channel vector can be written as

$$\mathbf{h}_u = \sum_{\ell=1}^L \alpha_\ell \mathbf{a}(\phi_\ell), \quad (3)$$

where $\mathbf{a}(\phi_\ell)$ is the array response vector of the base station. The definition of $\mathbf{a}(\phi_\ell)$ depends on the array geometry and hardware impairments. Next, we discuss that in more detail.

C. Hardware Impairments Model

Most of the prior work on mmWave signal processing has assumed uniform antenna arrays with perfect calibration and ideal hardware [2], [3], [11], [12]. In this paper, we consider a more general antenna array model that accounts for arbitrary geometry and hardware imperfections, and target learning efficient beam codebooks for these systems. This is very important for several reasons: (i) there are scenarios where designing arbitrary arrays is needed, for example, to improve the angular resolution or enhance the direction-of-arrival estimation performance [19], [20], (ii) the fabrication process of large mmWave arrays normally has some imperfections, and (iii) the calibration process of the mmWave phased arrays is an expensive process that requires special high-performance RF circuits [21]. While the codebook learning solutions that we develop in this paper are general for various kinds of arrays and hardware impairments, we evaluate them in Section VII with respect to two main characteristics of interest, namely non-uniform spacing and phase mismatch between the antenna elements. For linear arrays, the array response vector can be modeled to capture these characteristics as follows

$$\mathbf{a}(\phi_\ell) = [e^{j(kd_1 \cos(\phi_\ell) + \Delta\theta_1)}, e^{j(kd_2 \cos(\phi_\ell) + \Delta\theta_2)}, \dots, e^{j(kd_M \cos(\phi_\ell) + \Delta\theta_M)}]^T, \quad (4)$$

where d_m is the position of the m -th antenna, and $\Delta\theta_m$ is the additional phase shift incurred at the m -th antenna (to model the phase mismatch). Without loss of generality, we assume that d_m and $\Delta\theta_m$ are fixed yet unknown random realizations, obtained from the distributions $\mathcal{N}((m-1)d, \sigma_d^2)$ and $\mathcal{N}(0, \sigma_p^2)$, respectively, where σ_d and σ_p model the standard deviations of the random antenna spacing and phase mismatch. Besides, we impose an additional constraint $d_1 < d_2 < \dots < d_M$ to make sure the generated antenna positions physically meaningful.

III. PROBLEM DEFINITION

In this paper, we investigate the design of mmWave beamforming codebooks that are adaptive to the specific deployment (surrounding environment, user distribution, etc.) and the given base station hardware (array geometry, hardware imperfections, etc.). Given the system and channel models described in Section II, the SNR after combining for user u can be written as

$$\text{SNR}_u = \frac{|\mathbf{w}^H \mathbf{h}_u|^2}{|\mathbf{w}|^2} \rho, \quad (5)$$

with $\rho = \frac{P_x}{\sigma_n^2}$. Besides, we define the beamforming/combining gain of adopting \mathbf{w} as a transmit/receive beamformer for user u as

$$g_u = |\mathbf{w}^H \mathbf{h}_u|^2. \quad (6)$$

If the combining vector \mathbf{w} is selected from a codebook \mathcal{W} , with cardinality $|\mathcal{W}| = N$, then, the maximum achievable SNR for user u is obtained by the exhaustive search over the beam codebook as

$$\text{SNR}_u^* = \rho \max_{\mathbf{w} \in \mathcal{W}} |\mathbf{w}^H \mathbf{h}_u|^2, \quad (7)$$

where we set $\|\mathbf{w}\|^2 = 1$ as these combining weights are implemented using only phase shifters with constant magnitudes of $1/\sqrt{M}$, as described in (1). The objective of this paper is to design (learn) the beam codebook \mathcal{W} to maximize the SNR given by (7) averaged over the set of the users that can be served by the base station. Let \mathcal{H} represent the set of channel vectors for all the users that can be served by the considered base station, the beam codebook design problem can be formulated as

$$\mathcal{W}_{\text{opt}} = \arg \max_{\mathcal{W}} \frac{1}{|\mathcal{H}|} \sum_{\mathbf{h}_u \in \mathcal{H}} \left(\max_{\mathbf{w}_n \in \mathcal{W}} |\mathbf{w}_n^H \mathbf{h}_u|^2 \right), \quad (8)$$

$$\text{s. t. } w_{mn} = \frac{1}{\sqrt{M}} e^{j\theta_{mn}}, \quad \forall m = 1, \dots, M, n = 1, \dots, N, \quad (9)$$

$$\theta_{mn} \in \Theta, \quad \forall m = 1, \dots, M, n = 1, \dots, N, \quad (10)$$

where $w_{mn} = [\mathbf{w}_n]_m$ is the m -th element of the n -th beamforming vector in the codebook, $|\mathcal{H}| = K$ is the total number of users, Θ is the set that contains the 2^r possible phase shifts. It is worth mentioning that the constraint in (9) is imposed to uphold the adopted model where the analog beamformer can only perform phase shifts to the received signal, and the constraint in (10) is to respect the quantized phase-shifters hardware constraint.

Due to the unknown array geometry as well as possible hardware impairments, the accurate channel state information is generally hard to acquire. This means that all the channels $\mathbf{h}_u \in \mathcal{H}$ in the objective function are possibly unknown. Instead, the base station may only have access to the beamforming/combining gain g_u (or equivalently, the Received Signal Strength Indicator (RSSI) reported by each user if a downlink setup is considered). Therefore, problem (8) is hard to solve in a general sense for the unknown parameters in the objective function as well as the non-convex constraint (9) and the discrete constraint (10). Given that **this problem is essentially a search problem in a dauntingly huge yet finite and discrete space**, we consider leveraging the powerful exploration capability of deep reinforcement learning to efficiently search over the space to find the optimal or near-optimal solution. Since the number of beams in the codebook is far less than the number of channels in \mathcal{H} , then users sharing similar channels are expected

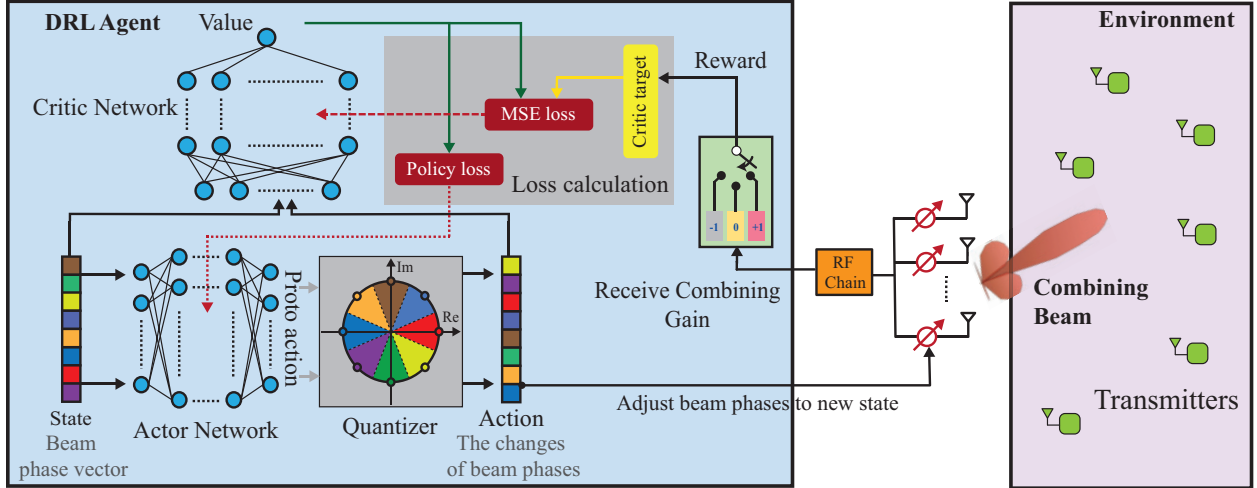


Fig. 2. The proposed beam pattern design framework with deep reinforcement learning. The schematic shows the agent architecture, and the way it interacts with the environment.

to be served by the same beam that achieves the best average beamforming gain compared to the other beams in the codebook. Based on this idea, we consider solving the original problem (8) in two steps. First, we investigate the problem of learning an optimized beam pattern for a single user or a group of users that share similar channels in Section IV, which we refer to as *the beam pattern learning problem* and it can be formulated as

$$\mathbf{w}_{\text{opt}} = \arg \max_{\mathbf{w}} \frac{1}{|\mathcal{H}_s|} \sum_{\mathbf{h}_u \in \mathcal{H}_s} |\mathbf{w}^H \mathbf{h}_u|^2, \quad (11)$$

$$\text{s. t. } w_m = \frac{1}{\sqrt{M}} e^{j\theta_m}, \quad \theta_m \in \Theta, \quad \forall m = 1, \dots, M, \quad (12)$$

where w_m is the m -th element of the beamforming vector and \mathcal{H}_s is the channel set that is supposed to contain a single channel or multiple similar channels. Then, in Section V, we address the codebook design problem (8) by introducing a joint clustering, assignment, and beam pattern learning approach.

IV. BEAM PATTERN LEARNING

In this section, we present our proposed deep reinforcement learning based algorithm for addressing the beam pattern design problem (11), which aims at maximizing the (averaged) beamforming gain of a single user (or a group of users that share the similar channels). Given constraint (12), the design problem is essentially a search problem over a finite yet dauntingly

huge discrete feasible set. For example, for a base station equipped with 32 antennas and 3-bit phase shifters, there are over 7.9×10^{28} legitimate beamforming vectors. With this huge space, finding the optimal beamforming vector by using methods like exhaustive search is definitely infeasible. In order to achieve an efficient search process, we resort to reinforcement learning where the agent (base station) is able to learn from what it has experienced (receive power or feedbacks from users) and try to proceed towards a better direction (beamforming/combining vector). However, when viewing the problem from a reinforcement learning perspective, it features **finite** yet **very high dimensional** action space. This makes the traditional learning frameworks (such as deep Q-learning, deep deterministic policy gradient, etc.) hard to apply. To deal with that, we propose a learning framework based on Wolpertinger architecture [15] to narrow the size of the action space and avoid missing the optimal policy at the same time.

A. Wolpertinger Architecture Overview

It is commonly known that deep Q-networks [22], [23] are very difficult to apply when the number of actions in the action space (which we refer to as the *dimension* of the action space) is huge. This is because the dimension of the output of the deep Q-network relates directly to the number of possible actions, which means that the size of the neural network will keep growing as the number of actions increases. However, for problems approaching real life complexity, it is highly likely to encounter applications that involve a huge action space, different from that in Atari games [22] where only several actions are considered. For example, in our problem, the possible actions in the case given above (where the base station has 32 antennas and adopts 3-bit phase shifters) are in the order of 10^{28} . The number can increase further with more antennas and higher resolution phase shifters. This is definitely intractable for the deep Q-network framework. With this motivation, the Wolpertinger architecture is proposed as a way of reasoning in a space with large number of discrete actions [15]. The Wolpertinger architecture is based on the actor-critic [24] framework and is trained using Deep Deterministic Policy Gradient (DDPG) [25]. This novel architecture utilizes a K-Nearest Neighbor (KNN) classifier to make DDPG suitable for tasks with discrete, finite yet very high dimensional action space. We briefly introduce the basic components of the Wolpertinger architecture as follows.

1) *Actor Network*: We assume an action space $\mathcal{A} \subseteq \mathbb{R}^n$ that is discrete and finite (but possibly with large number of actions), from which the agent selects an action to execute. We also assume a state space $\mathcal{S} \subseteq \mathbb{R}^m$ that contains all the possible states of an environment. We will define

the action and state spaces in the context of the beam pattern learning problem in Section IV-B. The actor network is then constructed as a function approximator parameterized by θ^μ mapping from the state space to the \mathbb{R}^n , that is

$$\mu(\cdot|\theta^\mu) : \mathcal{S} \rightarrow \mathbb{R}^n. \quad (13)$$

Due to the discrete and finite nature of \mathcal{A} , the action predicted by the actor network is probably not within \mathcal{A} . In other words, for any state $\mathbf{s} \in \mathcal{S}$, we can get a predicted *proto-action*

$$\mu(\mathbf{s}|\theta^\mu) = \hat{\mathbf{a}}, \quad (14)$$

where $\hat{\mathbf{a}}$ is highly likely not a “legitimate” action, i.e. $\hat{\mathbf{a}} \notin \mathcal{A}$. Therefore, the proto-action $\hat{\mathbf{a}}$ needs to be transformed (quantized) to a valid action in \mathcal{A} , where the KNN classifier plays a role in.

2) *K-Nearest Neighbor*: Since the predicted proto-action of the actor network is possibly not a valid action, we need to map $\hat{\mathbf{a}}$ to valid actions in \mathcal{A} . One natural solution could be a KNN function, that is, finding k actions in \mathcal{A} that are most closest to $\hat{\mathbf{a}}$ by some distance metric (L_2 distance to name one). More to the point, we assume that there is a function denoted by ξ_k . This function takes in the proto-action $\hat{\mathbf{a}}$ and returns the k nearest neighbors of that proto-action in \mathcal{A} according to L_2 distance, formally

$$\xi_k(\hat{\mathbf{a}}) = \arg \min_{\mathbf{a} \in \mathcal{A}}^{\text{nearest } k} \|\mathbf{a} - \hat{\mathbf{a}}\|_2. \quad (15)$$

The output of $\xi_k(\hat{\mathbf{a}})$ is a set of k actions in \mathcal{A} that are the top k nearest neighbors to $\hat{\mathbf{a}}$, which is denoted by $\mathcal{A}_k = \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k\}$.

3) *Critic Network*: The critic network is constructed as a function approximator parameterized by θ^Q mapping from the joint state space \mathcal{S} and action space \mathcal{A} to \mathbb{R} , that is

$$Q(\cdot, \cdot|\theta^Q) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}. \quad (16)$$

The critic network essentially plays the role of a Q function that takes in the state and action and outputs the predicted Q value of this particular state-action pair. Since k actions are obtained from the KNN function, the critic network then evaluates k state-action pairs (note that they share the same state) and selects the action that achieves the highest Q value

$$\mathbf{a}_t = \arg \max_{\mathbf{a}_t \in \mathcal{A}_k} Q(\mathbf{s}_t, \mathbf{a}_t|\theta^Q). \quad (17)$$

4) *Network Update*: The actor network aims at maximizing the output of the critic network (the predicted Q value) given a particular state, the objective of which can be simply expressed as $J(\theta^\mu) = \mathbb{E} [Q(\mathbf{s}, \mathbf{a}|\mathbf{a}=\mu(\mathbf{s}|\theta^\mu))]$. Thus, the actor policy is updated using the deep deterministic policy gradient, which is given by

$$-\nabla_{\theta^\mu} J(\theta^\mu) \approx -\mathbb{E} [\nabla_{\mathbf{a}} Q(\mathbf{s}, \mathbf{a}) \nabla_{\theta^\mu} \mu(\mathbf{s}|\theta^\mu)] \quad (18)$$

$$\approx -\frac{1}{B} \sum_{b=1}^B \nabla_{\mathbf{a}} Q(\mathbf{s}, \mathbf{a})|_{\mathbf{s}=\mathbf{s}_b, \mathbf{a}=\mu(\mathbf{s}_b|\theta^\mu)} \nabla_{\theta^\mu} \mu(\mathbf{s}|\theta^\mu)|_{\mathbf{s}=\mathbf{s}_b}. \quad (19)$$

The objective of the critic network is to estimate the Q value of the input state-action pair. Thus, the target can be constructed in the exactly same way that is adopted in the deep Q-networks, which is given by

$$y = \mathbb{E} \left[r + \gamma \max_{\mathbf{a}_{t+1}} Q(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}|\theta^Q) \right]. \quad (20)$$

The parameters of the critic network θ^Q is then updated based on the mean squared error over a particular mini-batch, which is given by $\frac{1}{B} \sum_{b=1}^B (y_b - Q(\mathbf{s}_b, \mathbf{a}_b|\theta^Q))^2$.

For the sake of computational stability, the actor and critic networks have duplicates, referred to as the target actor and target critic networks. They are not trainable like the actor and critic networks, but they are utilized for calculating the targets. Despite them being not trainable, the parameters of the target actor and critic networks get updated using the parameters of the critic and actor networks after a certain number of training iterations. Formally, it can be expressed as

$$\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'}, \quad (21)$$

$$\theta^{\mu'} \leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'}, \quad (22)$$

where $\theta^{\mu'}$ and $\theta^{Q'}$ are the parameters of target actor network and target critic network, τ is a non-negative hyper-parameter usually taking a value far less than 1.

B. DRL Based Beam Pattern Design

In this subsection, we describe in detail our proposed DRL based beam pattern design approach. We adopt the Wolpertinger architecture described above as our learning framework.

1) *Reinforcement Learning Setup*: To solve the problem with reinforcement learning, we first specify the corresponding building blocks of the learning algorithm.

- **State**: We define the state \mathbf{s}_t as a vector that consists of the phases of all the phase shifters at the t -th iteration, that is, $\mathbf{s}_t = [\theta_1, \theta_2, \dots, \theta_M]^T$. This phase vector can be converted to the

actual beamforming vector by applying (1). Since all the phases in \mathbf{s}_t are selected from Θ , and all the phase values in Θ are within $(-\pi, \pi]$, (1) essentially defines a bijective mapping from the phase vector to the beamforming vector. Therefore, for simplicity, we will use the term “beamforming vector” to refer to both this phase vector and the actual beamforming vector (the conversion is by (1)), according to the context.

- **Action:** We define the action \mathbf{a}_t as the element-wise changes to all the phases in \mathbf{s}_t . Since the phases can only take values in Θ , a change of a phase means that the phase shifter selects a value from Θ . Therefore, the action is directly specified as the next state, i.e. $\mathbf{s}_{t+1} = \mathbf{a}_t$.
- **Reward:** We define a ternary reward mechanism, i.e. the reward r_t takes values from $\{+1, 0, -1\}$. We compare the beamforming gain achieved by the current beamforming vector, denoted by g_t , with two values: (i) an adaptive threshold β_t , and (ii) the previous beamforming gain g_{t-1} . The reward is computed using the following rule
 - $g_t > \beta_t, r_t = +1$;
 - $g_t \leq \beta_t$ and $g_t > g_{t-1}, r_t = 0$;
 - $g_t \leq \beta_t$ and $g_t \leq g_{t-1}, r_t = -1$.

We adopt an adaptive threshold mechanism that does not rely on any prior knowledge of the channel distribution. The threshold has an initial value of zero. When the base station tries a beam and the resulting beamforming/combining gain surpasses the current threshold, the system updates the threshold by the value of this beamforming/combining gain. Besides, because the update of threshold also marks a successful detection of a new beam that achieves the best beamforming/combining gain so far, the base station also records this beamforming vector. As can be seen in this process, in order to evaluate the quality of a beam (or equivalently, calculate the reward), the system always tracks two quantities, which are the previous beamforming/combining gain and the best beamforming/combining gain achieved so far (i.e. the threshold).

2) *Environment Interaction:* As mentioned in Sections I and III, due to the possible hardware impairments, accurate channel state information is generally unavailable. Therefore, the base station can only resort to the receive power (or beamforming gain feedback reported by the users in a downlink setup) to adjust its beam pattern in order to achieve a better performance. To be more specific, upon forming a new beam $\tilde{\mathbf{w}}$, the base station uses this beam to receive

Algorithm 1 DRL Based Beam Pattern Learning

- 1: Initialize actor network $\mu(\mathbf{s}|\theta^\mu)$ and critic network $Q(\mathbf{s}, \mathbf{a}|\theta^Q)$ with random weights θ^μ and θ^Q
 - 2: Initialize target networks μ' and Q' with the weights of actor and critic networks' $\theta^{\mu'} \leftarrow \theta^\mu$ and $\theta^{Q'} \leftarrow \theta^Q$
 - 3: Initialize the replay memory \mathcal{D} , minibatch size B , discount factor γ
 - 4: Initialize adaptive threshold $\beta = 0$ and the previous average beamforming gain $g_1 = 0$
 - 5: Initialize a random process \mathcal{N} for action exploration
 - 6: Initialize a random beamforming vector \mathbf{w}_1 as the initial state \mathbf{s}_1
 - 7: **for** $t = 1$ to T **do**
 - 8: Receive a proto-action from actor network with exploration noise $\hat{\mathbf{a}}_t = \mu(\mathbf{s}_t|\theta^\mu) + \mathcal{N}_t$
 - 9: Quantize the proto-action to a valid beamforming vector \mathbf{a}_t according to (24)
 - 10: Execute action \mathbf{a}_t , observe reward r_t and update state to $\mathbf{s}_{t+1} = \mathbf{a}_t$
 - 11: Update the threshold β and the previous beamforming gain g_t
 - 12: Store the transition $(\mathbf{s}_t, \mathbf{a}_t, r_t, \mathbf{s}_{t+1})$ in \mathcal{D}
 - 13: Sample a random mini batch of B transitions $(\mathbf{s}_b, \mathbf{a}_b, r_b, \mathbf{s}_{b+1})$ from \mathcal{D}
 - 14: Calculate target $y_b = r_b + \gamma Q'(\mathbf{s}_{b+1}, \mu'(\mathbf{s}_{b+1}|\theta^{\mu'})|\theta^{Q'})$
 - 15: Update the critic network by minimizing the loss $L = \frac{1}{B} \sum_b (y_b - Q(\mathbf{s}_b, \mathbf{a}_b|\theta^Q))^2$
 - 16: Update the actor network using the sampled policy gradient given by (19)
 - 17: Update the target networks every C iterations by (21) and (22)
 - 18: **end for**
-

the symbols transmitted by every user. Then, it averages all the combining gains as follows

$$\bar{g} = \frac{1}{|\mathcal{H}_s|} \sum_{\mathbf{h}_u \in \mathcal{H}_s} |\tilde{\mathbf{w}}^H \mathbf{h}_u|^2, \quad (23)$$

where \mathcal{H}_s represents the targeted user channel set. Recall that (23) is the same as evaluating the objective function of (11) with the current beamforming vector $\tilde{\mathbf{w}}$. Depending on whether or not the new average beamforming/combining gain surpasses the previous beamforming/combining gain as well as the current threshold, the base station gets either reward or penalty, based on which it can judge the “quality” of the current beam and decide how to move.

3) *Exploration*: The exploration happens after the actor network predicts the proto-action $\hat{\mathbf{a}}_{t+1}$ based on the current state (beam) \mathbf{s}_t . Upon obtaining the proto-action, an additive noise is added

element-wisely to $\widehat{\mathbf{a}}_{t+1}$ for the purpose of exploration, which is a customary way in the context of reinforcement learning with continuous action spaces [24], [25]. In our problem, we use temporally correlated noise samples generated by an Ornstein-Uhlenbeck process [26], which is also used in [15], [25]. It is worth mentioning that a proper configuration of the noise generation parameters has significant impact on the learning process. Normally, the extent of exploration (noise power) is set to be a decreasing function with respect to the iteration number, which is commonly known as exploration-exploitation tradeoff [24]. Furthermore, the exact configuration of noise power should relate to specific applications. In our problem, for example, the noise is directly added to the predicted phases. Thus, at the very beginning, the noise should be strong enough to perturb the predicted phase to any other phases in Θ . By contrast, when the learning process approaches to the termination (the learned beam already performs well), the noise power should be decreased to a smaller level that is only capable of perturbing the predicted phase to its adjacent phases in Θ .

4) *Quantization*: The “proto” beam (with exploration noise added) should be quantized in order to be a valid new beam. To this end, we apply a KNN classifier as described in Section IV-A2. Furthermore, we specify $k = 1$ in (15), which is basically a nearest neighbor lookup. Therefore, each quantized phase in the new vector can be simply calculated as

$$[\mathbf{s}_{t+1}]_m = \arg \min_{\theta \in \Theta} |\theta - [\widehat{\mathbf{s}}_{t+1}]_m|, \forall m = 1, 2, \dots, M. \quad (24)$$

5) *Forward Computation and Backward Update*: The current state \mathbf{s}_t and the new state \mathbf{s}_{t+1} (recall that we directly set $\mathbf{s}_{t+1} = \mathbf{a}_t$) are then fed into the critic network to compute the Q value, based on which the targets of both actor and critic networks are calculated. This completes a forward pass. Following that, a backward update is performed to the parameters of the actor and critic networks. A pseudo code of the algorithm can be found in Algorithm 1.

V. BEAM CODEBOOK LEARNING

In this section, we propose a multi-network DRL approach for solving (8) and learning a beam codebook. The solution is built around the beam pattern learning approach described in Section IV. It could be briefly described as a pipeline of three key stages, namely clustering, assignment, and beam pattern learning. The first stage learns to *partition* the users in the environment into clusters based on how similar their channels are (without explicitly estimating those channels). These clusters are, then, *assigned* to different DRL networks in the second

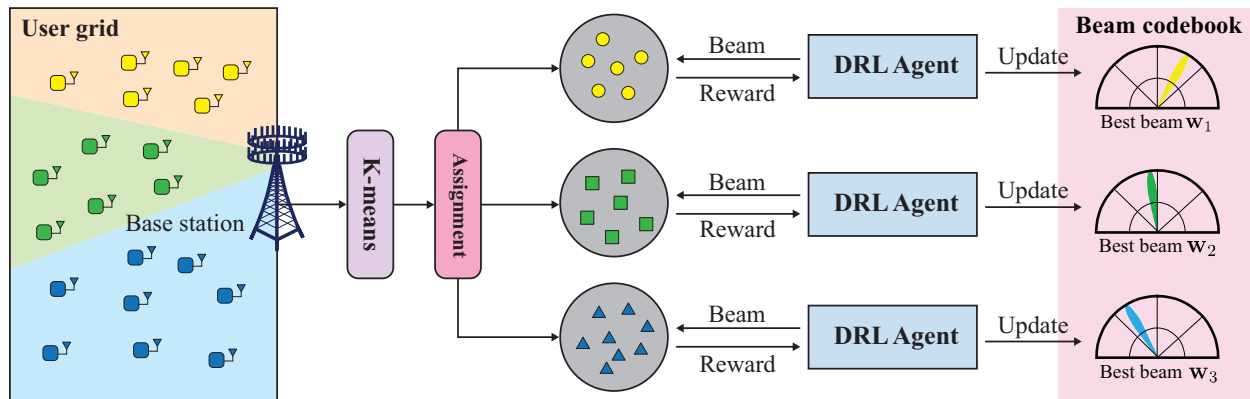


Fig. 3. The proposed beam codebook design framework with deep reinforcement learning. It shows the overall architecture, including the user clustering component.

stage. Such assignment needs to guarantee a form of consistency among the clusters that are assigned to the networks during the learning process. Finally, the third stage is where *the beam pattern learning* happens. Each of the DRL networks is expected to learn a beam pattern, and the collection of those patterns constructs the beam codebook. We detail this approach in the following three subsections.

A. User Clustering

Following into the footsteps of [27], users sharing similar channels are served by the same beam in the codebook. The question then becomes how to cluster the users' channels without knowing them, i.e., without performing expensive channel estimation. As a result of the constant modulus and limited resolution phase shifters, the set of feasible beamforming vectors for (8) forms a huge yet finite subset of \mathbb{C}^M , and all those vectors live on the surface of the M -dimensional unit hypersphere. The proposed clustering method here relies on utilizing a random subset of those vectors, henceforth referred to as *the sensing beams*, for the purpose of gathering sensing information in the form of receive combining gain. This information is used to cluster those users, developing a rough sense of their distribution in the environment.

To perform the clustering, the method starts by constructing a matrix that is comprised of receive combining gains using the sensing beams. Formally, let $\mathcal{F} = \{\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_S\}$ be

a set of S sensing beams that are randomly sampled from the feasible set of (8)¹, where $\mathbf{f}_s \in \mathbb{C}^M$, $\forall s \in \{1, \dots, S\}$. Also, let $\mathcal{H}_{\text{sen}} = \{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_{K'}\}$ denote the channels of the K' users that contribute to the clustering process, where $\mathcal{H}_{\text{sens}} \subseteq \mathcal{H}$. It is worth mentioning that these K' users do not need to be present in the environment at the same time. The receive combining gains used in the clustering algorithm can be collected over a relatively long period of time. This is because essentially the learned clustering is a function of the major elements (e.g. walls, buildings, large trees, etc.) of the environment. Such scatterers/reflectors commonly stay static over long periods of time. As a result, the sensing data can be collected in an accumulative manner and the learned classifier does not need to be updated (re-trained) frequently. The objective then is to collect receive combining gains from the K' users for every beam $\mathbf{f}_s \in \mathcal{F}$. More specifically, the receive combining gains are used to construct the *sensing matrix* \mathbf{P}

$$\mathbf{P} \triangleq \begin{bmatrix} |\mathbf{f}_1^H \mathbf{h}_1|^2 & \cdots & |\mathbf{f}_1^H \mathbf{h}_k|^2 & \cdots & |\mathbf{f}_1^H \mathbf{h}_{K'}|^2 \\ |\mathbf{f}_2^H \mathbf{h}_1|^2 & \cdots & |\mathbf{f}_2^H \mathbf{h}_k|^2 & \cdots & |\mathbf{f}_2^H \mathbf{h}_{K'}|^2 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ |\mathbf{f}_S^H \mathbf{h}_1|^2 & \cdots & |\mathbf{f}_S^H \mathbf{h}_k|^2 & \cdots & |\mathbf{f}_S^H \mathbf{h}_{K'}|^2 \end{bmatrix}, \quad (25)$$

where each column in \mathbf{P} has the receive combining gains from the same user for all sensing beams in \mathcal{F} . It is worth mentioning that since the receive combining gain is the **only** information source to the base station, the sensing matrix \mathbf{P} actually incorporates **all** the information that the base station can leverage from the outside environment.

The sensing matrix is used to extract feature vectors that characterize the user distribution in the environment. Each column in \mathbf{P} represents the receive gains of a single user in the environment. One could cluster the users by directly applying a clustering algorithm (such as k-means) on the columns of \mathbf{P} . However, empirical evidence shows that this clustering does not yield meaningful partitioning of the users (or equivalently the channels). The reason for that could be attributed to the fact that the columns of \mathbf{P} are restricted to the nonnegative orthant of the \mathbb{R}^S vector space; this increases the likelihood of overlapping clusters, which are hard to separate with k-means. As an alternative, we propose to transform the column of \mathbf{P} using

¹To clarify, we use \mathcal{W} and \mathbf{w} to denote the learned codebook and beam. We use \mathcal{F} and \mathbf{f} to denote the sensing beam set and sensing beam.

pair-wise differences. More precisely, the pair-wise differences of the elements of every column are computed, scaled, and stacked in a column vector as follows

$$\mathbf{u}_k = \left(\frac{1}{S} \sum_{s=1}^S |\mathbf{f}_s^H \mathbf{h}_k|^2 \right)^{-1} \begin{bmatrix} |\mathbf{f}_1^H \mathbf{h}_k|^2 - |\mathbf{f}_2^H \mathbf{h}_k|^2 \\ |\mathbf{f}_1^H \mathbf{h}_k|^2 - |\mathbf{f}_3^H \mathbf{h}_k|^2 \\ \vdots \\ |\mathbf{f}_{S-1}^H \mathbf{h}_k|^2 - |\mathbf{f}_S^H \mathbf{h}_k|^2 \end{bmatrix}, \quad \forall k = 1, 2, \dots, K', \quad (26)$$

where \mathbf{u}_k is referred to as the *feature vector* of user k . The column vectors of all K' users are organized in a feature matrix $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{K'}]$. This choice of transformation preserves the relation between the channel vector of a user and the sensing vectors, i.e., the sense of how close a channel vector to each sensing vector. However, it expresses that relation using a feature vector that could fall anywhere in the $\mathbb{R}^{S(S-1)/2}$ vector space (not restricted to the nonnegative orthant). The factor in (26) expresses each elements in the columns of \mathbf{U} as a ratio of a pair-wise difference to the average power of the corresponding column of matrix \mathbf{P} .

The clustering is applied on the columns of the feature matrix \mathbf{U} to produce N clusters. The popular k-means algorithm [28] is adopted to generate those clusters. The algorithm learns to partition the K' users (or equivalently their channels \mathcal{H}_{sen}) into N disjoint subsets

$$\mathcal{H}_{\text{sen}} = \mathcal{H}_1 \cup \mathcal{H}_2 \cup \dots \cup \mathcal{H}_N, \quad (27)$$

where $\mathcal{H}_k \cap \mathcal{H}_l = \emptyset, \forall k \neq l$ and we assume that the subscript of each user group is also the corresponding label of that group. The trained k-means algorithm is used to classify any new user coming into the environment. It is important to note here that the learned clustering is a function of the major elements of the environment not the user distribution, i.e., it is mainly affected by major scatterers and their positions like walls, buildings, large trees, etc. Such scatterers commonly change over long periods of time, and consequently, the learned clusters do not need to be updated frequently.

B. Cluster Assignment

Since the clustering will be frequently repeated whenever there is a change in the environment, an important question arises: how to assign the new clusters to the existing DRL networks, with each of them learning one beam? The answer to this question defines the second stage in our proposed codebook learning approach. For the learning process to be meaningful, a network should consistently be assigned channel clusters that exhibit some form of similarity; the new

cluster should be similar to the previous one in the sense that the network can improve its currently learned beam pattern but not change it completely. To that end, we formulate this cluster assignment task as a linear sum assignment problem, which can be solved efficiently using the *Hungarian algorithm* [29]. In such problem, every pair of new cluster and DRL network is assigned a cost reflecting how suitable this cluster to the network, and the goal is to find N **unique cluster-network assignments** that minimize the total cost sum (total suitability).

To perform the cluster-network assignment, a cost needs to be computed to measure suitability and guide the assignment process. Let $\widehat{\mathcal{H}}_{\text{sen}} = \widehat{\mathcal{H}}_1 \cup \widehat{\mathcal{H}}_2 \cup \dots \cup \widehat{\mathcal{H}}_N$ be the new clusters obtained using the clustering algorithm described in Section V-A. As described in Section IV-B2, the DRL network always tracks the beamforming vectors that achieve the best beamforming gain, which forms a set of “temporarily best” beamforming vectors, denoted by $\Upsilon = \{\widehat{\mathbf{w}}_1, \widehat{\mathbf{w}}_2, \dots, \widehat{\mathbf{w}}_N\}$, where the subscripts stand for the indices of the N DRL networks. We propose to use the average beamforming gain of each beamforming vector in Υ computed on each cluster as the suitability measure. The result of that forms a cost matrix \mathbf{Z} , where the value at the intersection of n -th row and n' -th column of \mathbf{Z} stands for the average beamforming gain of the n -th temporarily best beamforming vector in Υ on the n' -th channel cluster in $\widehat{\mathcal{H}}_{\text{sen}}$. This value is calculated as

$$\mathbf{Z}_{nn'} = \frac{1}{|\widehat{\mathcal{H}}_{n'}|} \sum_{\mathbf{h} \in \widehat{\mathcal{H}}_{n'}} |\widehat{\mathbf{w}}_n^H \mathbf{h}|^2. \quad (28)$$

With the cost matrix, we formulate the cluster assignment task as a linear sum assignment problem, which is given by

$$\min_{\mathbf{X}} - \sum_{n=1}^N \sum_{n'=1}^N \mathbf{X}_{nn'} \mathbf{Z}_{nn'} \quad (29)$$

$$\text{s. t. } \mathbf{X} \text{ is a permutation matrix.} \quad (30)$$

This problem can be efficiently solved by using Hungarian algorithm, the results of which are N association tuples

$$(\widehat{\mathbf{w}}_n, \widehat{\mathcal{H}}_{n'}), \quad n, n' \in \{1, 2, \dots, N\}.$$

In other words, the cluster assignment step forms a bijective mapping from Υ to the set of channel groups

$$\{\widehat{\mathbf{w}}_1, \widehat{\mathbf{w}}_2, \dots, \widehat{\mathbf{w}}_N\} \iff \{\widehat{\mathcal{H}}_1, \widehat{\mathcal{H}}_2, \dots, \widehat{\mathcal{H}}_N\}. \quad (31)$$

Algorithm 2 User Clustering and Cluster Assignment Algorithm

- 1: Initialize a sensing beam set $\mathcal{F} = \{\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_S\}$
 - 2: Initialize the temporarily best beam set $\Upsilon = \{\widehat{\mathbf{w}}_1, \widehat{\mathbf{w}}_2, \dots, \widehat{\mathbf{w}}_N\}$
 - 3: Construct sensing matrix \mathbf{P} by (25)
 - 4: Transform sensing matrix \mathbf{P} to feature matrix \mathbf{U} by applying (26) to the columns of \mathbf{P}
 - 5: Use k-means algorithm to cluster the columns of \mathbf{U} into N clusters
 - 6: **while** environment has not changed **do**
 - 7: Randomly sample a subset of users $\widehat{\mathcal{H}}$ from \mathcal{H}
 - 8: Partition sampled channels using the trained k-means classifier to $\widehat{\mathcal{H}} = \widehat{\mathcal{H}}_1 \cup \widehat{\mathcal{H}}_2 \cup \dots \cup \widehat{\mathcal{H}}_N$
channel clusters
 - 9: Construct the matrix \mathbf{Z} using Υ and the clustering result of $\widehat{\mathcal{H}}$ based on (28)
 - 10: Solve the optimization problem (29) by applying Hungarian algorithm
 - 11: Assign the user clusters to DRL networks based on the association relationship given by
the permutation matrix \mathbf{X}
 - 12: Train the N DRL networks
 - 13: **if** training saturated **then**
 - 14: Fine-tune the learned beam pattern using perturb-and-quantize operations
 - 15: **end if**
 - 16: **end while**
 - 17: Go to line 1
-

C. Neural Network Update and Fine-Tuning

Upon obtaining the clustered channels and their assignment (31), the problem (8) is essentially decomposed into N independent sub-problems which is given by (11). Each DRL network adjusts its own beam based on the assigned user cluster. They only consider the receive combining gains from their designated users. User clustering and cluster assignment are two key stages that enable adaptability and empower the proposed solution with capability of dealing with dynamic environment. Practically speaking, it is impossible to fix all the users until a good beam codebook is learned. Instead, we keep learning cluster and assign the users as they change over time, which partially reflects the dynamics of the environment. Our proposed beam codebook approach accounts for such practical considerations and is able to learn beam codebooks that adapt to the

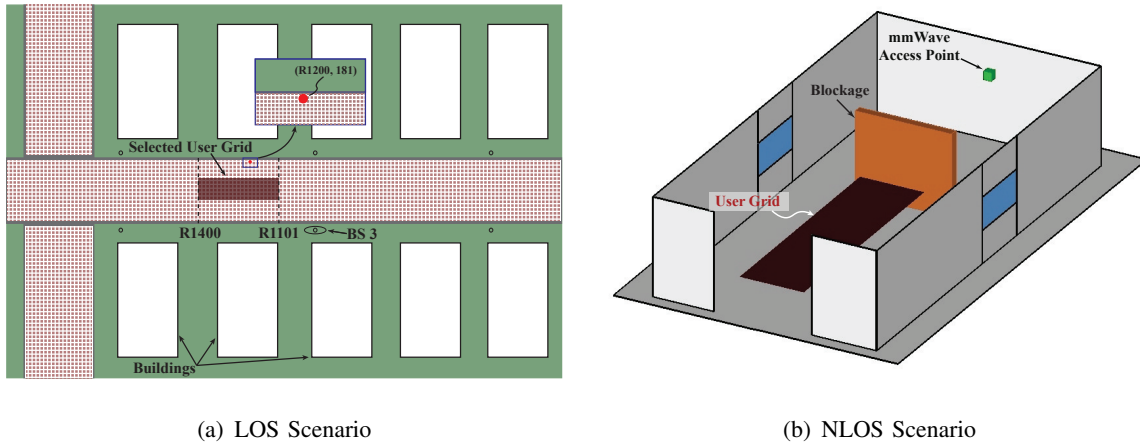


Fig. 4. Two perspective views of the considered communication scenarios. (a) shows the LOS scenario. It is chosen to be outdoor since the likelihood of LOS connection is higher there. (b) shows the NLOS scenario. Similar to (a), this scenario has been chosen for the high likelihood of having NLOS users indoors.

environment. The complete beam codebook learning algorithm is given in Algorithm 2.

The beam pattern learning proceeds as described in Section IV-B with one minor difference, a final perturb-and-quantize fine-tuning step. This step is basically applied after the DRL agent reaches training saturation. It is composed of three simple operations: (i) perturb the beam vector with exploration noise, (ii) quantize the perturb beam vector, and (iii) evaluate the quantized beam vector on the assigned cluster of users. The training algorithm loops over the three operations until the received beamforming gain saturates again. The goal of this last stage is to fine-tune the beam pattern without the relatively expensive agent-training process.

VI. EXPERIMENTS SETUP AND NETWORK TRAINING

To evaluate the performance of the proposed solutions, two scenarios are considered. They are designed to represent two different communication settings. The first has all users experiencing LOS connection with the BS, while the other has them experiencing NLOS connection. The following two subsections provide more details on the scenarios and the training process.

A. Communication Scenarios and Datasets

Two scenarios are used for performance evaluation, as shown in Fig. 4. The first one is an outdoor LOS scenario where all users have LOS connection with the mmWave BS, with an operating frequency of 60 GHz. The second one is chosen to be an indoor NLOS scenario

where all the users have NLOS connection with the mmWave BS, with an operating frequency of 28 GHz. Both scenarios are part of DeepMIMO dataset [16]. Using the DeepMIMO scripts, two sets of channels, namely \mathcal{S}^{LOS} and $\mathcal{S}^{\text{NLOS}}$, are generated, one for each scenario. Table I shows the data generation hyper-parameters. The datasets taking into account the hardware impairments are generated based on the LOS scenario. While our proposed solution can deal with general impairments, we only consider two main sources of impairments, namely antenna spacing and phase mismatches. We generate multiple datasets based on different levels of impairments, measured by the standard deviations of antenna spacing and phase mismatches. Without distinction of them, we denote those datasets with impairments as $\mathcal{S}^{\text{cLOS}}$ ².

B. Machine Learning Model Structure and Pre-processing

While we generate multiple datasets, the learning architecture is the same, which is based on the DDPG framework. It is comprised of two networks, actor and critic. The input of the actor network is the state, i.e. the phases of the phase shifters, hence with a dimension of M . There are two hidden layers, all comprising $16M$ neurons and followed by Rectified Linear Unit (ReLU) activations. The output of the actor network is the predicted action, which also has a dimension of M and is followed by hyperbolic tangent (tanh) activations scaled by π . For the critic network, the input is the concatenation of the state and action, so it has a dimension of $2M$. There are also two hidden layers, all with $32M$ neurons and followed by ReLU activations. The output of the critic network stands for the predicted Q value of the input state-action pair, which is a real scalar (dimension of 1). The hyper-parameters for training can be found in Table II. The training process starts by data pre-processing. The channels in each dataset are normalized to improve the training experience, which is a very common practice in machine learning [30]. As in [27], [31], the channel normalization using the maximum absolute value in the dataset helps the network undergo a stable and efficient training. Formally, the normalization factor is found as follows

$$\Delta = \max_{\mathbf{h}_u \in \mathcal{S}} |[\mathbf{h}_u]_m|, \quad (32)$$

where $\mathcal{S} \in \{\mathcal{S}^{\text{LOS}}, \mathcal{S}^{\text{NLOS}}, \mathcal{S}^{\text{cLOS}}\}$ and $[\mathbf{h}_u]_m$ is the m -th element in the channel vector \mathbf{h}_u .

²“cLOS” is a shorthand of corrupted LOS.

TABLE I
HYPER-PARAMETERS FOR CHANNEL GENERATION

Parameter	value	
Name of scenario	O1_60	I2_28B
Active BS	3	1
Active users	1101 to 1400	201 to 300
Number of antennas (x, y, z)	(1, 32, 1)	(32, 1, 1)
System BW	0.5 GHz	0.5 GHz
Antenna spacing	0.5	0.5
Number of OFDM sub-carriers	1	1
OFDM sampling factor	1	1
OFDM limit	1	1
Number of multi paths	5	5

TABLE II
HYPER-PARAMETERS FOR MODEL TRAINING

Parameter	value	
Models	Actor	Critic
Replay memory size	8192	8192
Mini-batch size	1024	1024
Learning rate	10^{-3}	10^{-3}
Weight decay	10^{-2}	10^{-3}

VII. SIMULATION RESULTS

In this section, we evaluate the performance of the proposed solution using the scenarios described in Section VI. In a nutshell, the numerical results show that our proposed learning solutions can adapt to different environments, user distributions as well as hardware impairments, **without the need to estimate the channels**. We compare the performance of the learned codebook with classical beamsteering codebook, where the beamforming vectors are spatial matched filters for the single-path channels. Therefore, they have the same form of the array response vector and can be parameterized by a simple angle [32]. In our simulation, depending on the adopted size of the classical beamsteering codebook, those angles are evenly spaced in the range of $[0, \pi]$. Next, we will first evaluate the performance of the beam pattern learning solution in Section VII-A, and then evaluate the beam codebook learning solution in Section VII-B.

A. Beam Pattern Learning

We first evaluate our proposed DRL-based beam pattern learning solution on learning a single beam that serves a single user with LOS connection to the BS. The selected target user is highlighted in Fig. 4(a) with a red dot. In Fig. 5, we compare the performance of the learned single beam with a 32-beam classical beamsteering codebook. As only known, classical beamsteering codebook normally performs very well in LOS scenario. However, our proposed method achieves higher beamforming gain than the best beam in the classical beamsteering

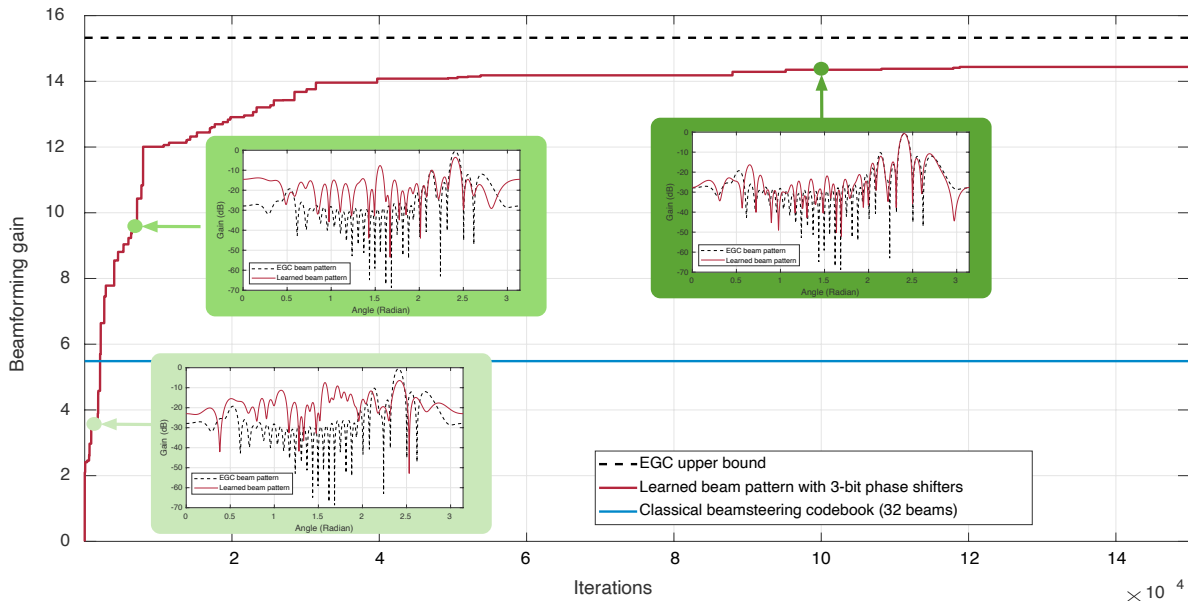


Fig. 5. The beam pattern learning results for a single user with LOS connection to the base station. The base station employs a perfect uniform linear array with 32 antennas and 3-bit phase shifters. In this figure, we show the learning process and the beam patterns learned at three different stages during the iterations. The learned beam patterns are plotted using solid red line, and the equal gain combining/beamforming vector is plotted using dashed black line.

codebook, with negligible iterations. More interestingly, with less than 4×10^4 iterations, the proposed solution can reach more than 90% of the EGC upper bound. It is worth mentioning that the EGC upper bound can only be reached when the user's channel is known and unquantized phase shifters are deployed. By contrast, our proposed solution can finally achieve almost 95% of the EGC upper bound with 3-bit phase shifters and without any channel information. We also plot the learned beam patterns at three different stages (iteration 1000, 5000, and 100000) during the learning process, which helps understand how the beam pattern evolves over time. As shown in Fig. 5, at iteration 1000, the learned beam pattern has very strong side lobes, weakening the main lobe gain to a great extent. At iteration 5000, the gain of the main lobe becomes stronger. However, there are still multiple side lobes with relatively high gains. Finally, at iteration 100000, it can be seen that the main lobe has quite strong gain compared to the other side lobes, having at least 10 dB gain over the second strongest side lobe. And most of the side lobes are below -20 dB. Besides, the learned beam pattern captures the EGC beam pattern very well, which explains the good performance it achieves. The slight mismatching is mainly caused by the use of quantized phase shifters, which is with only 3-bit resolution.

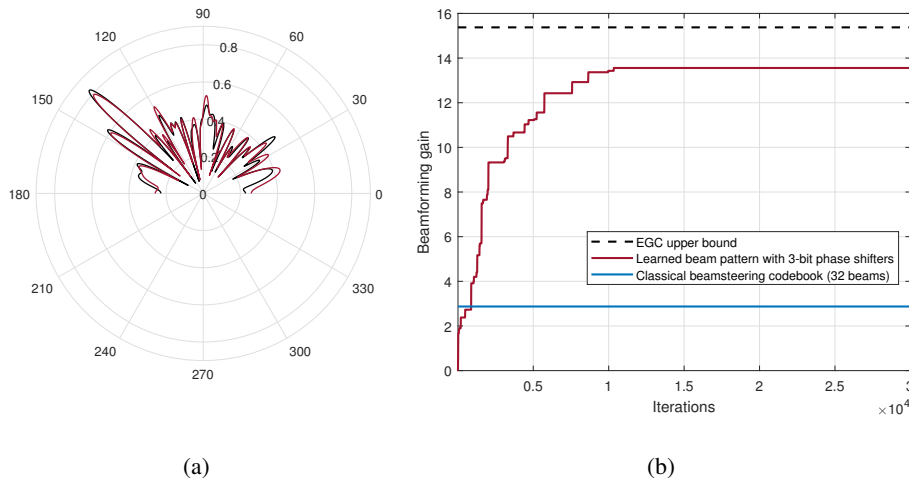


Fig. 6. The beam pattern learned for a single user with LOS connection to the base station. The base station employs a uniform linear array with 32 antennas and 3-bit phase shifters, where hardware impairments exist. The standard deviation of the antenna spacing is 0.1λ and the standard deviation of the phase mismatches is 0.32π . (a) shows the beam patterns for the equal gain combining/beamforming vector (red) and the learned beam (blue). A transformation of $\sqrt[4]{\cdot}$ is used to better show the finer structure of the beams. (b) shows the learning process.

The proposed beam pattern learning solution is also evaluated on a more realistic situation where hardware impairments exist (with the same user considered above). The simulation results confirm that our proposed solution is competent to learn optimized beam pattern that adapts to hardware, showing the capability of compensating the unknown hardware mismatches. Fig. 6 (a) shows the beam patterns for both EGC beam and the learned beam. At the first glance, the learned beam appears distorted and has lots of side-lobes. However, the performance of such beam is excellent, which can be explained by comparing its beam pattern with the EGC beam. **As can be seen from the learned beam pattern, our proposed solution intelligently approximates the optimal beam, where all the dominant lobes are well captured.** By contrast, the classical beamsteering codebook fails when the hardware is not perfect, as depicted in Fig. 6 (b). This is because the distorted array pattern incurred by the hardware impairment makes the pointed classical beamsteering beams only able to capture a small portion of the energy, resulting in a huge degradation in beamforming gain. The learned beam shown in Fig. 6 (a) is capable of achieving more than 90% of the EGC upper bound with approximately only 10^4 iterations, as shown in Fig. 6 (b). This is especially interesting for the fact that the proposed solution does not rely on any channel state information. As is known, the channel estimation in this case relies first on a full calibration of the hardware, which is a hard and expensive process.

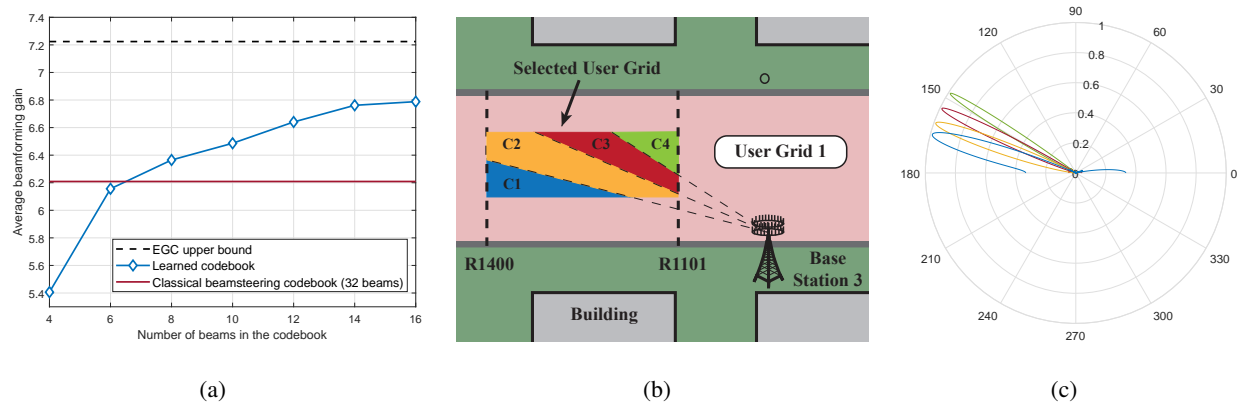


Fig. 7. The learning results of the proposed DRL-based codebook learning solution under a LOS scenario where the base station employs a perfect uniform linear array. (a) shows the average beamforming gain versus the number of beams in the codebook. (b) shows the result of clustering users into 4 groups. (c) shows the beam patterns for the learned 4-beam codebook in (a), which is based on the clustering result in (b).

B. Beam Codebook Learning

In this subsection, we evaluate our proposed DRL-based beam codebook learning solution in several scenarios. The task of learning a beam codebook with multiple beams is significantly different than learning a single beam (pattern) from computational complexity perspective. For example, for a base station with 32 antennas and 4-bit discrete phase shifters, there are 16^{32} possible beamforming vectors, from which a single vector is selected in the beam pattern learning case. However, learning a codebook will further result in finding combinations out of this huge pool. To address this problem, we propose a clustering and assignment approach, given by Algorithm 2, that essentially decomposes the huge task into N independent, parallel and relatively lightweight sub-tasks. This facilitates the problem of learning a codebook with multiple beams. Before we dive into the discussions, it is important to mention that due to the stationarity of our scenario, clustering/assignment is performed only once in our simulations. If the environment is more dynamic, the clustering/assignment is expected to be done more frequently.

Fig. 7 (a) plots the average beamforming gain versus the number of beams in the codebook under the LOS scenario shown in Fig. 4(a), where the BS adopts an ideal uniform linear array. It shows that the average beamforming gain is monotonically increasing as the number of beams increases. Besides, with only 6 beams, the proposed solution has almost the same performance as a 32-beam classical beamsteering codebook. And with 8 beams, it outperforms the 32-beam classical beamsteering codebook. This exhibits how our proposed approach adapts the beams

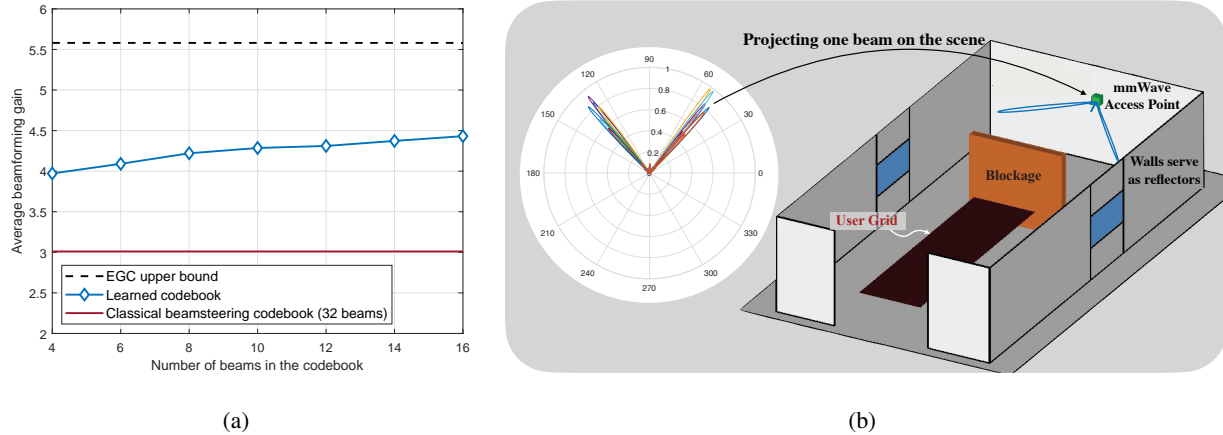


Fig. 8. The learning results of the proposed DRL-based codebook solution under a NLOS scenario. (a) shows the average beamforming gain versus the number of beams in the codebook. (b) shows the beam patterns of the learned 16-beam codebook in (a) and it also shows how one of the learned beams with multi-lobes fits the propagation environment.

based on the user distributions. **As a result, it significantly reduces the training overhead by avoiding scanning directions where there is no user at all.** In Fig. 7 (b), we present the clustering result for the users in this LOS scenario. This is a very important step for learning multiple beams. As stated at the end of Section III, the ultimately optimized codebook should have a collection of beams, where each one of them is optimized to serve a group of users with similar channels. The clustering stage is the first step that our proposed solution takes to attain that objective. Fig. 7 (c) depicts the beam patterns of the learned 4-beam codebook. As shown in the learning result, the proposed solution can cluster users based on the similarity of their channels, and form beams to cover the user grid in order to achieve high beamforming gain.

We also evaluate the proposed solution under a NLOS scenario shown in Fig. 4(b), where all the users experience NLOS connection with an indoor mmWave access point. As can be seen in Fig. 8 (a), the proposed solution surpasses a 32-beam classical beamsteering codebook with only 4 beams. Further, the proposed solution is gradually approaching the EGC upper bound as the size of codebook increases. It should note that in order to achieve the EGC upper bound: (i) The number of beams in the codebook should be equal to the number of users, (ii) continuous phase shifters should be adopted, and most importantly, (iii) accurate channel state information is needed. By contrast, with only 16 beams and 4-bit phase shifters, the proposed solution can reach 80% of the EGC upper bound, relying only on the receive combining gains. **In other words, our approach not only significantly reduces the beam training overhead, but also**

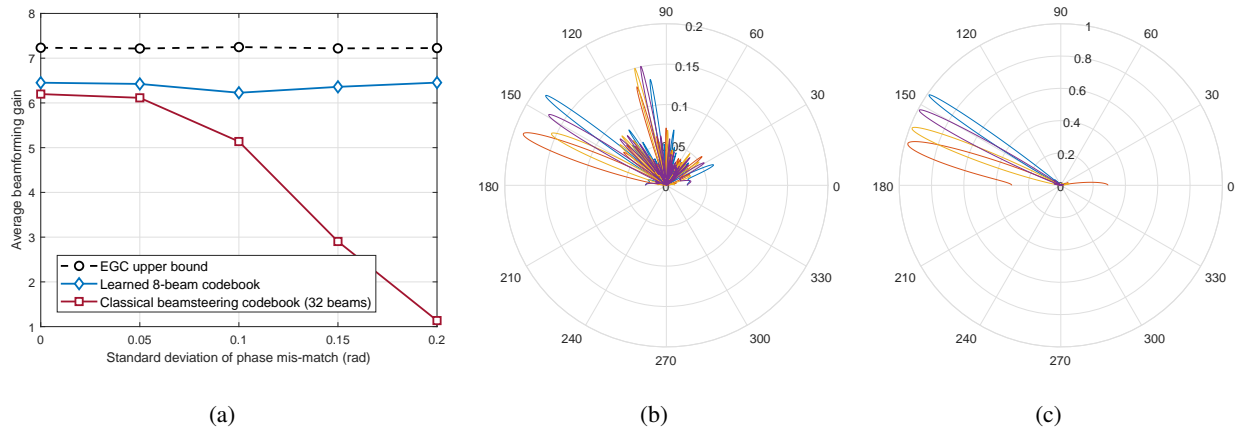


Fig. 9. The learning results of the proposed DRL-based solution under the same LOS scenario with hardware impairments being considered. (a) shows the average beamforming gain versus the standard deviation of phase mismatch, where the antenna spacing mismatch has a fixed standard deviation of 0.1λ . (b) shows the beam patterns of 4 beams in the learned codebook projected onto the “clean” angular space. (c) shows the same beams as in (b) projected onto the “corrupted” angular space.

avoids the prohibitive cost of estimating the channels. To gain more insight, we plot the beam patterns of the learned 16-beam codebook in Fig. 8 (b) and project one of the beams on the adopted scene. It can be seen in Fig. 8 (b) that the learned beams have multi-lobes, different than the pointed beams learned in the LOS scenario. However, such beams achieve better performance compared to the pointed beamsteering beams. The reason becomes clear when we observe that because of the blockage in the considered scenario, the signals transmitted by the users have to resort to reflections to reach the access point, where the walls at both sides of the room serve as major reflectors. This clearly shows how our proposed solution adapts the beam pattern to the propagation environment, gaining more power by receiving signals from multiple directions.

Learning codebooks that overcome the hardware impairments is one of the most important application cases of our DRL-based codebook learning approach. Therefore, we evaluate the proposed solution under the same LOS scenario shown in Fig. 4(a), with hardware impairments being considered. Furthermore, we test our solution under different standard deviations of the phase mismatch, where a fixed antenna spacing mismatch with a standard deviation of 0.1λ is assumed. For each channel dataset, we learn a 8-beam codebook and compare it with a 32-beam classical beamsteering codebook. In Fig. 9 (a), we plot the average beamforming gain versus the standard deviation of the phase mismatch. The result shows that as the standard deviation of the phase mismatch increases, i.e. the hardware impairments become more severe, the proposed DRL based solution keeps a balanced performance. The slight fluctuation is mainly caused by the

uncertainty nature of solving the highly non-convex problem (8). By contrast, the performance of the 32-beam classical beamsteering codebook degrades drastically as the level of hardware impairment increases. This empirically shows the robustness of the proposed codebook learning approach to different levels of hardware impairments. In Fig. 9 (b), we plot the beam patterns of 4 beams in the learned codebook. It can be seen that these beams have quite distorted beam patterns like the single beam case shown in Fig. 6 (a). To show that these distorted beam patterns indeed **match** the hardware impairments, we plot the same beams in Fig. 9 (c), but project them on the “corrupted” angular space. As illustrated in Fig. 9 (c), the learned beams actually appear “clean” and pointy in the corrupted angular space. This empirically verifies the capability of the proposed solution in learning beams that adapt to the flawed hardware.

VIII. CONCLUSIONS AND DISCUSSIONS

In this paper, we considered the problem of designing environment and hardware aware beam codebooks for mmWave/THz MIMO systems with hardware-constrained architectures and without requiring channel knowledge. We developed a DRL based framework that learns how to adapt the beam patterns to the surrounding environment, user distribution and hardware impairments, relying only on receive power measurements. We also introduced a clustering-assignment approach to efficiently learn beam codebooks without requiring any knowledge about the user positions and without requiring the users to be stationary during the learning process. The developed approach is evaluated at various environments, as well as different levels of hardware impairments. Simulation results show promising capability of the proposed solution in learning environment and hardware aware codebooks relying only on the receive power measurements and with finite resolution phase shifters. The learned codebook outperforms the classical beamsteering codebook with much smaller number of beams in all studied scenarios. Further, the beams learned by the proposed solution approach the beam shapes of the ideal beams (characterized by unconstrained beamforming vectors under full channel knowledge) and achieve similar SNR performance. For the future work, it will be interesting to extend the developed framework to other MIMO systems that as those adopting hybrid analog/digital architectures.

REFERENCES

- [1] Y. Zhang, M. Alrabeiah, and A. Alkhateeb, “Reinforcement Learning for Beam Pattern Design in Millimeter Wave and Massive MIMO Systems,” in *Proc. of IEEE Asilomar Conference on Signals, Systems, and Computers*, Nov 2020.
- [2] A. Alkhateeb, J. Mo, N. Gonzalez-Prelcic, and R. W. Heath, “MIMO Precoding and Combining Solutions for Millimeter-Wave Systems,” *IEEE Communications Magazine*, vol. 52, no. 12, pp. 122–131, 2014.

- [3] A. Alkhateeb, O. El Ayach, G. Leus, and R. 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.
- [4] M. Giordani, M. Polese, A. Roy, D. Castor, and M. Zorzi, "A Tutorial on Beam Management for 3GPP NR at mmWave Frequencies," *IEEE Communications Surveys Tutorials*, vol. 21, no. 1, pp. 173–196, 2019.
- [5] IEEE 802.11ad, "IEEE 802.11ad standard draft D0.1." [Online]. Available: www.ieee802.org/11/Reports/tgadupdate.htm
- [6] M. Alrabeiah, A. Hredzak, and A. Alkhateeb, "Millimeter wave base stations with cameras: Vision-aided beam and blockage prediction," in *IEEE Vehicular Technology Conference (VTC2020-Spring)*, 2020, pp. 1–5.
- [7] T. K. Y. Lo, "Maximum ratio transmission," *IEEE Transactions on Communications*, vol. 47, no. 10, pp. 1458–1461, 1999.
- [8] D. Love and R. Heath Jr, "Equal gain transmission in multiple-input multiple-output wireless systems," *IEEE Transactions on Communications*, vol. 51, no. 7, pp. 1102–1110, 2003.
- [9] X. Li, Y. Zhu, and P. Xia, "Enhanced Analog Beamforming for Single Carrier Millimeter Wave MIMO Systems," *IEEE Transactions on Wireless Communications*, vol. 16, no. 7, pp. 4261–4274, 2017.
- [10] O. El Ayach, S. Rajagopal, S. Abu-Surra, Z. Pi, and R. Heath, "Spatially sparse precoding in millimeter wave MIMO systems," *IEEE Transactions on Wireless Communications*, vol. 13, no. 3, pp. 1499–1513, Mar. 2014.
- [11] S. Hur, T. Kim, D. Love, J. Krogmeier, T. Thomas, and A. Ghosh, "Millimeter wave beamforming for wireless backhaul and access in small cell networks," *IEEE Transactions on Communications*, vol. 61, no. 10, pp. 4391–4403, Oct. 2013.
- [12] J. Wang, *et al.*, "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, Nov. 2009.
- [13] A. Alkhateeb and R. W. Heath, "Frequency selective hybrid precoding for limited feedback millimeter wave systems," *IEEE Transactions on Communications*, vol. 64, no. 5, pp. 1801–1818, May 2016.
- [14] M. Alrabeiah, Y. Zhang, and A. Alkhateeb, "Neural Networks Based Beam Codebooks: Learning mmWave Massive MIMO Beams that Adapt to Deployment and Hardware," 2020.
- [15] G. Dulac-Arnold, R. Evans, H. van Hasselt, P. Sunehag, T. Lillicrap, J. Hunt, T. Mann, T. Weber, T. Degris, and B. Coppin, "Deep reinforcement learning in large discrete action spaces," 2015.
- [16] 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.
- [17] R. W. Heath, *et al.* "An overview of signal processing techniques for millimeter wave MIMO systems," *IEEE Journal of Selected Topics in Signal Processing*, vol. 10, no. 3, pp. 436–453, April 2016.
- [18] M. Alrabeiah and A. Alkhateeb, "Deep Learning for TDD and FDD Massive MIMO: Mapping Channels in Space and Frequency," *arXiv e-prints*, p. arXiv:1905.03761, May 2019.
- [19] P. Pal and P. P. Vaidyanathan, "Nested arrays: A novel approach to array processing with enhanced degrees of freedom," *IEEE Transactions on Signal Processing*, vol. 58, no. 8, pp. 4167–4181, Aug 2010.
- [20] M. Rubsamen and A. B. Gershman, "Direction-of-arrival estimation for nonuniform sensor arrays: From manifold separation to fourier domain music methods," *IEEE Transactions on Signal Processing*, vol. 57, no. 2, pp. 588–599, 2009.
- [21] T. Moon, J. Gaun, and H. Hassanieh, "Online millimeter wave phased array calibration based on channel estimation," in *2019 IEEE 37th VLSI Test Symposium (VTS)*, 2019, pp. 1–6.
- [22] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing Atari with Deep Reinforcement Learning," 2013.
- [23] V. Mnih *et al.*, "Human-level Control through Deep Reinforcement Learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [24] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. Cambridge, MA, USA: A Bradford Book, 2018.
- [25] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," 2015.
- [26] G. E. Uhlenbeck and L. S. Ornstein, "On the Theory of the Brownian Motion," *Phys. Rev.*, vol. 36, pp. 823–841, Sep 1930. [Online]. Available: <https://link.aps.org/doi/10.1103/PhysRev.36.823>
- [27] Y. Zhang, M. Alrabeiah, and A. Alkhateeb, "Learning Beam Codebooks with Neural Networks: Towards Environment-Aware mmWave MIMO," 2020.
- [28] C. M. Bishop, *Pattern recognition and machine learning*. springer, 2006.
- [29] H. W. Kuhn, "The Hungarian method for the assignment problem," *Naval research logistics quarterly*, vol. 2, no. 1-2, pp. 83–97, 1955.
- [30] Y. A. LeCun, L. Bottou, G. B. Orr, and K.-R. Müller, "Efficient backprop," in *Neural networks: Tricks of the trade*. Springer, 2012, pp. 9–48.
- [31] Y. Zhang, M. Alrabeiah, and A. Alkhateeb, "Deep Learning for Massive MIMO with 1-Bit ADCs: When More Antennas Need Fewer Pilots," *IEEE Wireless Communications Letters*, 2020.
- [32] 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.