






Article

Adaptive Hybrid Beamforming Codebook Design Using Multi-Agent Reinforcement Learning for Multiuser Multiple-Input–Multiple-Output Systems

Manasjyoti Bhuyan ¹, Kandarpa Kumar Sarma ^{1,*}, Debashis Dev Misra ², Koushik Guha ³
and Jacopo Iannacci ^{4,*}

¹ Department of Electronics and Communication Engineering, Gauhati University, Guwahati 781014, Assam, India; manasjyoti.b@gmail.com

² Department of Computer Science and Engineering, Assam Down Town University, Guwahati 781026 Assam, India; debashish.dm@gmail.com

³ Department of Electronics and Communication Engineering, NIT Silchar, Silchar 788010 Assam, India; koushik@ece.nits.ac.in

⁴ MicroSystems Technology Research Unit, Center for Sensors and Devices (SD), Fondazione Bruno Kessler, Povo, 38123 Trento, Italy

* Correspondence: kandarpaks@gauhati.ac.in (K.K.S.); iannacci@fbk.eu (J.I.)

Abstract: This paper presents a novel approach to designing beam codebooks for downlink multiuser hybrid multiple-input–multiple-output (MIMO) wireless communication systems, leveraging multi-agent reinforcement learning (MARL). The primary objective is to develop an environment-specific beam codebook composed of non-interfering beams, learned by cooperative agents within the MARL framework. Machine learning (ML)-based beam codebook design for downlink communications have been based on channel state information (CSI) feedback or only reference signal received power (RSRP), consisting of an offline training and user clustering phase. In massive MIMO, the full CSI feedback data is of large size and is resource-intensive to process, making it challenging to implement efficiently. RSRP alone for a stand-alone base station is not a good marker of the position of a receiver. Hence, in this work, uplink CSI estimated at the base station along with feedback of RSRP and binary acknowledgment of the accuracy of received data is utilized to design the beamforming codebook at the base station. Simulations using sub-array antenna and ray-tracing channel demonstrate the proposed system's ability to learn topography-aware beam codebook for arbitrary beams serving multiple user groups simultaneously. The proposed method extends beyond mono-lobe and fixed beam architectures by dynamically adapting arbitrary shaped beams to avoid inter-beam interference, enhancing the overall system performance. This work leverages MARL's potential in creating efficient beam codebooks for hybrid MIMO systems, paving the way for enhanced multiuser communication in future wireless networks.



Citation: Bhuyan, M.; Sarma, K.K.; Misra, D.D.; Guha, K.; Iannacci, J. Adaptive Hybrid Beamforming Codebook Design Using Multi-Agent Reinforcement Learning for Multiuser Multiple-Input–Multiple-Output Systems. *Appl. Sci.* **2024**, *14*, 7109. <https://doi.org/10.3390/app14167109>

Academic Editor: Christos Bouras

Received: 10 July 2024

Revised: 3 August 2024

Accepted: 9 August 2024

Published: 13 August 2024

Keywords: multi-agent reinforcement learning; massive MIMO; millimeter wave; hybrid beamforming



Copyright: © 2024 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

MIMO systems employ multiple antennas at both the transmitter and receiver ends to enhance communication performance. MIMO technology leverages spatial diversity and multiplexing to increase data throughput and link reliability without requiring additional bandwidth or transmit power, making it fundamental in modern wireless communication systems, including 4G and 5G networks.

A key component of MIMO systems is beamforming, a signal processing technique used in antenna arrays for directional signal transmission or reception. By combining elements in an antenna array so that signals at particular angles experience constructive interference, beamforming improves signal quality and reduces interference. This spatial filtering capability enhances the communication range and capacity.

Building on the MIMO concept, massive MIMO employs hundreds or even thousands of antennas at the base station. This large-scale antenna system significantly boosts spectral and energy efficiency by spatially multiplexing a large number of users. Massive MIMO is crucial for 5G and future networks, enabling high data rates and improved reliability.

Massive MIMO and hybrid beamforming are interdependent in modern wireless communication systems, especially at millimeter-wave (mmWave) frequencies. While massive MIMO boosts data rates and spectral efficiency, its practical implementation faces challenges such as high hardware complexity and power consumption due to the need for dedicated radio frequency (RF) chains for each antenna. Hybrid beamforming addresses these challenges by combining analog and digital beamforming techniques, reducing the number of RF chains needed while still achieving high array gains and spatial multiplexing benefits. This makes massive MIMO systems more cost-effective and energy-efficient, facilitating widespread deployment in 5G and beyond.

Moreover, hybrid beamforming in mmWave massive MIMO systems effectively addresses practical issues like severe path loss and high power consumption. Recent studies show that hybrid beamforming can achieve a performance close to that of full digital beamforming with significantly lower complexity, enhancing the feasibility of implementing massive MIMO in real-world scenarios [1].

Hybrid beamforming merges analog and digital techniques for efficient signal transmission using compact, cost-effective quantized phase shifters. These phase shifters adjust the signal phase at the antenna level, improving signal quality. Analog weights form beams toward user groups, and digital weights handle MIMO tasks like interference cancellation within the groups. Hybrid setups with fully connected subarray antennas can produce multiple simultaneous beams, with phase-shifting states controlled by digital circuits for real-time adaptation. This integration offers practical, near-optimal performance in modern wireless systems. Advanced MIMO antenna designs integrating phase shifters highlight the potential for enhanced beamforming performance in complex wireless environments [2].

Efficient performance in multi-user mmWave systems involves serving multiple pieces of user equipment (UE) from each base station (BS) simultaneously. Precoding multiplexes different data streams to different users. However, fully digital baseband beamforming is impractical for multi-stream mmWave systems due to high costs and power consumption [3].

In mmWave systems, the large number of antennas and very low signal-to-noise ratio (SNR) before beamforming make it impractical to obtain full CSI for conventional closed-loop precoding matrix calculations [4]. Therefore, alternative beamforming and precoding techniques are necessary to achieve efficient performance while managing cost, power, and CSI availability. This requirement drives research interest in the field of hybrid beamforming in massive MIMO.

Hybrid precoding enables multiplexing multiple data streams by dividing processing between analog and digital domains [5,6]. For example, low-complexity hybrid precoding algorithms exploit the sparse nature of mmWave channels using basis pursuit algorithmic concepts, assuming channel knowledge [5]. Similarly, low-complexity hybrid beamforming algorithms for single-user single-stream MIMO-OFDM systems aim to maximize the received signal strength or sum-rate over different sub-carriers [7]. However, these algorithms were designed for single-user channels, limiting supported streams. In multi-user systems, digital precoding in hybrid setups can design precoders that reduce interference between users, making the development of near-optimal, low-complexity hybrid precoding algorithms for multi-user mmWave systems particularly important.

Large antenna arrays with quantized phase shifters have challenges. These phase shifters, with constant modulus, control only the phase, limiting applications to equal gain transmission schemes to maximize SNR or diversity gain [8].

Both artificial intelligence (AI) and non-AI methods have been explored in the literature to find optimal beamforming codebooks. A benchmark work [9] used deep learning (DL) to find optimal codebooks for transmit beamforming and combining at user terminals,

though it required channel matrix information at both the training and prediction phases. Reinforcement learning (RL) offers a promising solution by eliminating the need for the offline training phases of static deep networks and facilitating adaptive, situation-aware systems capable of learning from the environment. Significant works, such as [10], have implemented RL-based systems with Wolpertinger-variant architectures for beam codebook design, preceded by beam-clustering to reduce the codebook size. Beam clustering in this implementation is actuated through sensing beams. The author in [11] combined the radar-aided direction of arrival (DoA) and direction of departure (DoD) estimation and CS-based position estimation with hybrid beamforming for vehicular communication systems, emphasizing accurate direction estimation to improve the beamforming efficiency. In this work, author has eliminated CSI estimation feedback completely by using the radar-based subsystem. Hybrid beamforming solutions for multi-user millimeter-wave heterogeneous networks is developed in [12] where orthogonal matching pursuit-based analog beamforming and minimum mean square error (MMSE)-based digital beamforming is used.

To support multiple simultaneous users, the proposed work extends [10] by introducing MARL. MARL is an AI research area involving the development of intelligent agents that cooperate or compete to achieve common or individual goals. In a fully cooperative stochastic game, all agents share the same goal and work together to achieve it. Stochastic games involve uncertain action outcomes, with outcome probabilities depending on the current game state. In fully cooperative settings, centralized MARL is generally preferred for effective action coordination among agents. In this work, agent coordination translates to reducing interference between beams, serving different user groups simultaneously. Interference between users within a group served by a single beam is minimized through baseband precoding.

The authors in [10] utilized receiver signal strength (beamforming gain) for user clustering and codebook learning. Although this metric simplifies the design, it underperforms in nLOS environments and with user mobility. The analysis of average spatial autocorrelation functions for individual multipath components in the 28 GHz mmWave band shows that signals lose correlation after about two wavelengths in LOS environments and after about five wavelengths in nLOS environments [13]. This unreliability of signal strength alone challenges the learning agent's ability to gather useful environmental information.

The agent architecture in [10] is adopted for each RL agent in the proposed design. The Wolpertinger-variant structure adapts the continuous action space of deep deterministic policy gradient (DDPG) to work with large discrete action spaces [14]. Multi-agent deep deterministic policy gradient (MADDPG) is used to train the agents. MADDPG is designed for multi-agent systems, maintaining local actor and critic networks for decision making and action evaluation. During training, agents share experience replay buffers and learn from collective experiences to enhance their policies. This combination of centralized training and decentralized execution enables agents to learn effective strategies in complex multi-agent environments.

1.1. Motivation

Quantized phase shifter with fixed numbers of bits makes the search space for beams very large. For example, there will be 8^{64} beams for a 64-element antenna array with 3 quantization bits. In a multiuser case with four RF chains, this number will equate to 4×8^{64} . Finding optimal beams in such a finite but huge space is impractical with exhaustive search or any other traditional technique. Hence, it is a convention to use a beam codebook with large numbers of beams pointing in different directions in an effort to maximize the gain to the users in that direction. This approach is not optimal as this single-lobe beam, which matches filters to the array responses in a particular angle, is not guaranteed to offer maximum possible gain for occluded, non-line of sight (nLOS) users. Also, the large number of beams required in such a codebook renders the beam training inefficient and time consuming, and hence, is inapplicable to mobile users. Additionally,

an accurate array response is required to form such a beamsteering codebook which may not be available for cost-effective systems as calibrating antenna arrays is a sophisticated and costly process.

Motivated by the fact that the environment changes infrequently, the computation power and time are traded for accuracy by utilizing uplink CSI estimated at the base station, with re-clustering required only when the environment changes. CSI provides good autocorrelation properties and is relatively immune to hardware imperfections in the RF stage, making it suitable for beam learning applications. Perfect CSI is not assumed, as the channel is seen through the RF lens in hybrid beamforming. However, it is shown, in this work, that beam learning can be efficiently achieved even with uplink CSI estimates at the BS, by using uplink channel envelope estimates at the BS as fingerprints of specific user locations within the MARL framework.

It is important to note that a learned multibeam codebook differs from simply parallelizing multiple single beams. Unlike fixed codebook beams, learned beams are not matched filters to the antenna array response and can adopt any arbitrary shape suitable for the scattering environment of deployment. Consequently, this may result in beam overlap, leading to inter-user interference, if cooperation between the beams for multiple simultaneous users is not established. The proposed MARL-based approach addresses this issue by cooperatively learning the beamforming vector for each beam per user group.

In summary, RL is employed for beam learning due to its online learning capability, eliminating the need for an offline data collection phase. The Wolpertinger architecture is utilized for its efficiency in exploring vast search spaces for beam configurations. MARL is proposed to facilitate the selection of multiple parallel and diverse beams. Through centralized training, MARL can identify environment-specific, arbitrarily shaped beams that do not interfere with each other. Most existing approaches use RSRP only, which simplifies the design but RSRP does not provide good auto-correlation and is not a reliable indicator of a UE's position. In this work, uplink CSI is utilized, which, although it may not always replace downlink CSI for channel estimation in non-reciprocal systems with hardware imperfections, it offers a better representation of the implicit location of the UE.

1.2. Contribution

An RF codebook design approach for hybrid precoding algorithms in downlink multi-user mmWave systems is presented, demonstrating efficiency and effectiveness in mobile user environments. The proposed method does not require CSI feedback but learns the downlink beam codebook from RSRP feedback and CSI estimates for uplink sounding reference signals (SRSs). By using uplink CSI from SRSs as fingerprints for specific user locations, only sub-band sounding is needed. This increases the SNR at the base station, facilitating cell edge UE recovery. The proposed system replaces traditional fixed codebook beams with learned beams that adapt to the environment in an online process. The goal of this proposed work is to create a robust beamforming codebook for the base station in downlink communication, accommodating uncertain user locations even in mobile scenarios. The contributions from the proposed work are summarized as follows:

1. A multiuser hybrid mmWave MIMO system model designed as a fully cooperative stochastic game under the constraint of the quantized RF phase shifters is proposed. A novel algorithm is developed to realize this model. By employing MADDPG to train DDPG (Wolpertinger variant), this approach effectively minimizes interference among simultaneous users, preventing overlapping beams. Unlike previous methods such as [10], this work uniquely addresses and mitigates potential interference between nearby beams with arbitrary shapes, ensuring the unparalleled performance and reliability in beamforming for multiuser communication systems.
2. A reward function for the RL agent is proposed, considering the comprehensive performance of the end-to-end communication system. The reward for each agent is based on RSRP and binary ARQ status, indicating whether a particular sub-frame scheduled for a specific user is successfully received (ACK) or not (NACK). This

method allows the RL agent to maximize successful transmission rates by improving the beamforming gain and reducing the interference among simultaneous user groups. The cumulative reward function is optimized by MADDPG, enhancing the overall system efficiency and reliability.

3. The proposed system is rigorously evaluated through simulations using a realistic ray tracing channel model. This comprehensive testing spans various SNR and different codebook sizes. The results demonstrate the system’s robustness and efficiency, highlighting its adaptability and performance across diverse conditions.

Simulation results demonstrate that the proposed method can create optimized beam patterns without needing feedback from downlink CSI, relying instead on RSRP, binary ARQ status, and periodic channel estimates from a small sub-band within SRS sub-frames. This deep reinforcement learning-based method efficiently selects beams for the downlink RF codebook, requiring occasional updates, typically when there is a significant change in the operating environment or the base station’s position. The following sections delve into detailed discussions of the proposed systems, methods, algorithms, and results.

2. Multiuser Hybrid Beamforming System Model

The proposed multiuser system model is depicted in Figure 1, wherein a mmWave MIMO base station, equipped with N_{BS} antennas and N_{RF} RF chain, is in communication with M simultaneous users each having N_M antennas and one RF chain through N_t streams. Since each UE is assumed to be served by only one downlink stream and contains only one RF chain, analog combining is applied at the UE. This configuration is similar to works in [4]. The base station utilizes hybrid beamforming, employing a network of r-bit quantized phase shifters.

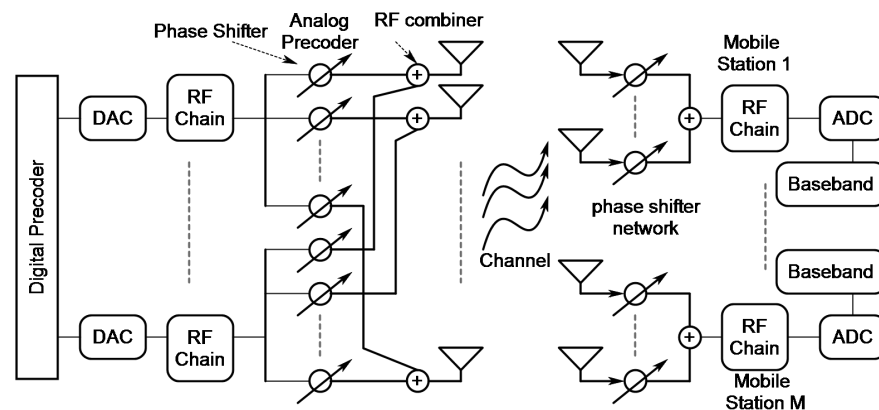


Figure 1. Transceiver architecture for multiuser hybrid beamforming.

The decision to employ a single RF chain per UE is motivated by practical considerations, aiming for lower complexity, cost, and power consumption. Conversely, the BS is equipped with advanced digital signal processing (DSP) capabilities designed to effectively handle multiple data streams.

BS attaches with each UE via one stream. This leads to a total of $N_t = M$ streams, where M represents the maximum number of simultaneous users the BS can serve at once. This aligns with the count of RF chains at the BS ($M \leq N_{RF}$), possible through hybrid schemes enabling spatial multiplexing and multi-user MIMO. This grants the BS the ability to communicate concurrently with multiple UEs using several beams. Design of the end-to-end communication system is shown in Figure 2. Parameter for each processing block in Figure 2 is shown in Table 1.

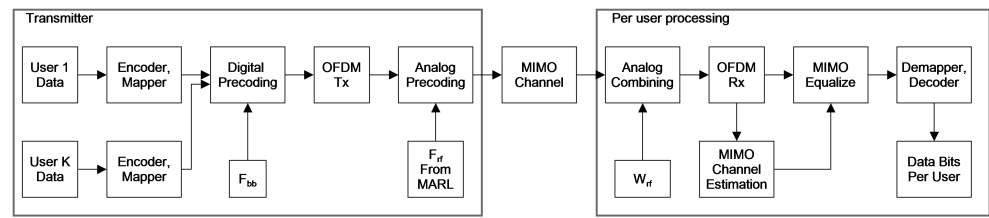


Figure 2. System model for data transmission and reception.

Table 1. Transceiver parameter table.

No. of users	4
Data streams per user	1
No. and type of base station antenna	32, ULA, isotropic, back-baffled
No. and type of receive antenna per user	4, ULA, non-back-baffled
Modulation type	16QAM
Frequency of operation	28 GHz
OFDM FFT length, CP length	256, 64
Encoder type, code rate (fixed)	Convolutional, $\frac{1}{3}$

In consideration of operational efficiency and hardware constraints, beamforming codebooks are commonly resorted to in mmWave and massive MIMO systems to effectively accommodate users. The sum rate achievable across all UEs is optimized using MADDPG in the proposed approach. Through MADDPG, the RF codebook at the base station is estimated. Represented by W , the beam codebook chosen by the base station consists of N beamforming and combining vectors, each crafted in accordance with the structure outlined in Equation (1).

$$w = \frac{1}{\sqrt{N_{BS}}} [e^{j\theta_1}, e^{j\theta_2}, \dots, e^{j\theta_{N_{BS}}}]^T \quad (1)$$

In this context, each phase shift θ_m is chosen from a finite set S containing 2^r discrete values, uniformly selected from the range of $(-\pi, \pi)$. Here, the parameter r represents the number of quantization bit used in phase shifters.

The BS employs baseband precoding denoted by $F_{BB} = [f_{BB_1}, f_{BB_2}, \dots, f_{BB_M}] \in C^{N_{RF} \times M}$ to process the transmit signal $s = [s_1, s_2, \dots, s_M]^T \in C^M$ in compliance with $E\{ss^H\} = \frac{P}{M} I_M$, assuming uniform power distribution among users. Notably, P signifies the average power. RF precoders $F_{RF} \in C^{N_{BS} \times N_{RF}}$, constructed using phase shifters are utilized to direct the signal to N_{BS} transmit antennas. Moreover, considering F_{RF} consists of analog phase shifters, constant equal-norm components in the RF precoder is assumed, i.e., $|[F_{RF}]_{i,j}|^2 = \frac{1}{N_{BS}}$. Furthermore, the power constraint $\|F_{RF}F_{BB}\|_F^2 = M$ is maintained through F_{BB} normalization. Consequently, the transmitted signal comprising $N_{BS} \times 1$ elements is given as

$$x = F_{RF}F_{BB}s. \quad (2)$$

considering $H_k \in C^{N_M \times N_{BS}}$ be the channel matrix between BS and k th user, the received signal for k th user for a narrowband block fading channel is given as

$$\tilde{y}_k = H_k \sum_{n=1}^M F_{RF}F_{BB_n}s_n + n_k, \quad (3)$$

where $n_k \in \mathbb{C}^{N_M}$ is the complex additive white Gaussian noise (AWGN) with $n_k \sim \mathcal{CN}(0, \sigma^2 I_{N_M})$. This signal \tilde{y}_k is received by the k th user and is processed by the combiner $W_{RF_k} \in \mathbb{C}^{N_M}$ to obtain $y_k = W_{RF_k}^H \tilde{y}_k$, i.e.,

$$y_k = W_{RF_k}^H H_k \sum_{n=1}^M F_{RF} F_{BB_n} s_n + W_{RF_k}^H n_k, \quad (4)$$

where RF combiners W_{RF} are designed as quantized phase shifter so that $|[W_{RF_k}]_i|^2 = 1/N_M$.

Achievable rate for the k th user assuming Gaussian symbol transmission through the channel is given as

$$R_k = \log_2 \left| 1 + \frac{\frac{P}{M} |W_{RF_k}^H H_k F_{RF} F_{BB_k}|^2}{\frac{P}{M} \sum_{n \neq k} |W_{RF_k}^H H_k F_{RF} F_{BB_n}|^2 + \sigma^2} \right| \quad (5)$$

Subsequently, the achievable sum rate for the system is given as [9]

$$\bar{R} = \sum_{k=1}^M R_k. \quad (6)$$

In this proposed work, F_{RF} as the beams of the learned codebook is acquired. W_{RF} through conventional beam sweeping is obtained, as detailed in Section 3. The acquisition of F_{BB} occurs in the second step of a two-step procedure, as outlined in [4]. It is important to note that, in this proposed work, the focus is solely on learning F_{RF} . The contribution can also be conceptualized as adaptive beam sectoring that is aware of the environment.

Channel Model

In this work, the ray tracing channel model is employed for simulation purposes. Stochastic channel models lack the spatial detail required for accurate beamforming simulations, making deterministic models like ray tracing preferable for such tasks. Ray tracing channel models treat electromagnetic waves as rays, accounting for interactions like reflection and diffraction with various surfaces in the environment. These models provide detailed insights into signal propagation, aiding in the design and optimization of wireless communication systems.

The ray tracing channel model is applied to an OpenStreetMap (.osm) file corresponding to Canary Wharf in London, UK. The latitude and longitude coordinates (51.50375, -0.01843) specify the BS location. The map is sourced from <https://www.openstreetmap.org> (accessed on 1 August 2023), providing crowd-sourced map data worldwide. Loaded into MATLAB for ray tracing simulation, the map defines transmitter and receiver sites. Multiple receivers are initialized with respective positions, simulating non-stationary users traversing the area. High-rise structures are represented using concrete as the building material.

Figure 3 depicts the ray tracing environment, where one UE experiences LOS conditions while the other encounters nLOS conditions.

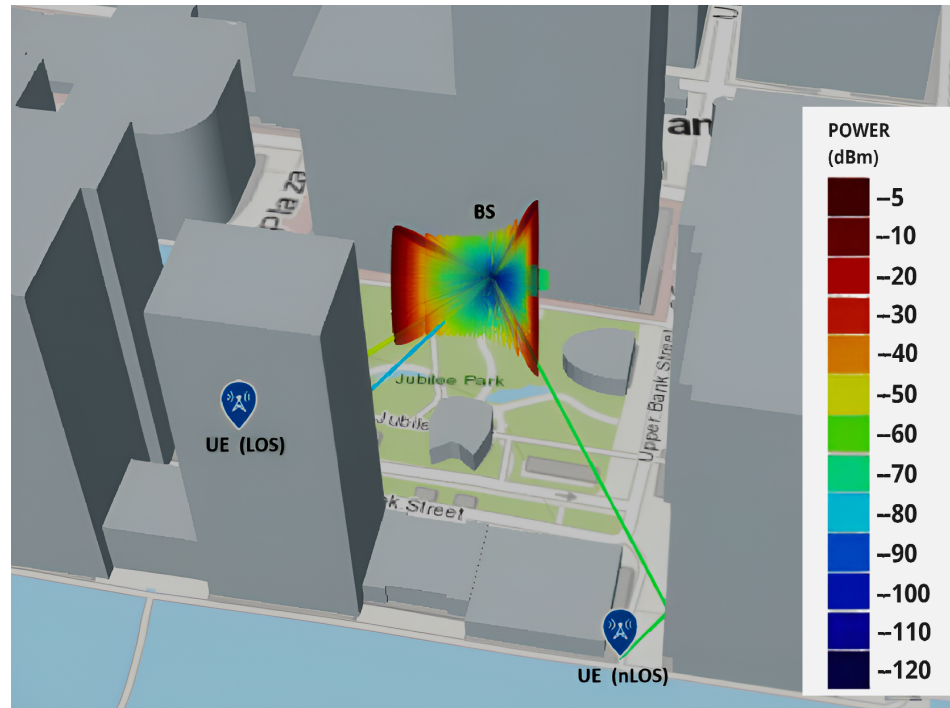


Figure 3. Topographic map of ray tracing environment with one LOS and one nLOS user. Radiation pattern for the 32 antenna BS is also shown forming beams toward each users.

3. Proposed Beam Learning with MARL

The implementation of the MARL can be broken down into sub-tasks, namely data collection through an initial access procedure, data preprocessing, and MARL agent training. Agent training in MARL is performed through continuous interaction with the environment and data collection and preprocessing are executed simultaneously. These processes are explained subsequently.

3.1. Data Collection through Initial Access Procedure

The codebook learning process is initiated by first obtaining a fixed conventional codebook through the initial access procedures and beam management procedures outlined in the 5G New Radio (5G NR) technical report. The steps involved in this fixed codebook learning process are summarized as follows:

Procedure 1: When a connection is established between a transmitter and a receiver, an initial beam alignment is required. This involves finding an optimal transmit–receive beam pair that maximizes the signal strength between the devices. Various methods like synchronization signal blocks (SSBs) and reference signals are used to aid in this process. This is shown in Figure 4.

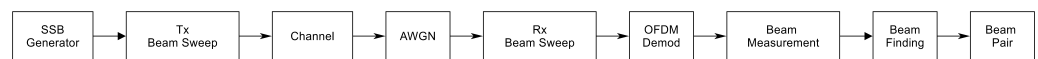


Figure 4. SSB beam search method in initial access procedure.

Procedure 2: Refining transmit-end beam via non-zero-power CSI-RS and SRS. After initial beam acquisition, this beam management aims to refine the beams to improve the communication link further. In this step, reference signals are sent in different directions using finer beams within the initial angular range. UE or BS assesses these beams with the fixed receive beam and selects the best transmit beam.

This proposed system initially employs a standard beamforming procedure and gradually transitions into a more efficient MARL-based system over time. This method essentially substitutes standard codebook beams with learned beams on a one-to-one correspondence

basis. The angular spacing between nearby beams is determined by the number of beams, which corresponds to the number of agents in the MARL framework. This approach simplifies implementation without necessitating alterations to the existing infrastructure. Once the codebook is learned, it can be utilized until the link’s performance deteriorates due to significant changes in the deployment site.

In this research, the procedures for initial beam acquisition and subsequent beam learning are segmented into the following major steps:

1. SSB beam sweeping;
2. Beam measurement and determination at UE;
3. Beam reporting to BS by UE;
4. Send SRS to BS from UE for uplink transmit end beam refinement and also for MARL-based downlink transmit end beam refinement. This procedure differs from method in 5G NR by the fact that the standard used NZ-CSI-RS for downlink transmit end beam refinement. This requires CSI feedback from UE and can work only with traditional matched filter-based codebooks as full-channel estimate feedback from UE which is required for non-codebook based beamforming is unavailable or impractical to achieve and resource intensive;
5. Send NZ-CSI-RP to UE only to obtain RSRP feedback (RSRP consumes very little resource);
6. Decode received SRS and estimate uplink Channel at BS;
7. Send RSRP measurement in SRS to UE for beam refinement at UE;
8. At BS, use RSRP and channel estimate acquired in step 5 and step 6 to learn downlink transmit end beam codebook through the proposed MARL algorithm.

3.2. Components of the Proposed MARL-Based Codebook Learning System and Its Implementation

To make MARL applicable, the environment must be modeled as a Markov process. In [10], this is achieved by incorporating the current beamforming vector as a function of the previous beamforming vector. A similar approach as in [10] is followed, extending this method by also considering the partial and imperfect CSI acquired by the BS during uplink SRS transmission by the UE. The operation of the entire system is illustrated in Figure 5, with each processing block and signal flow explained subsequently.

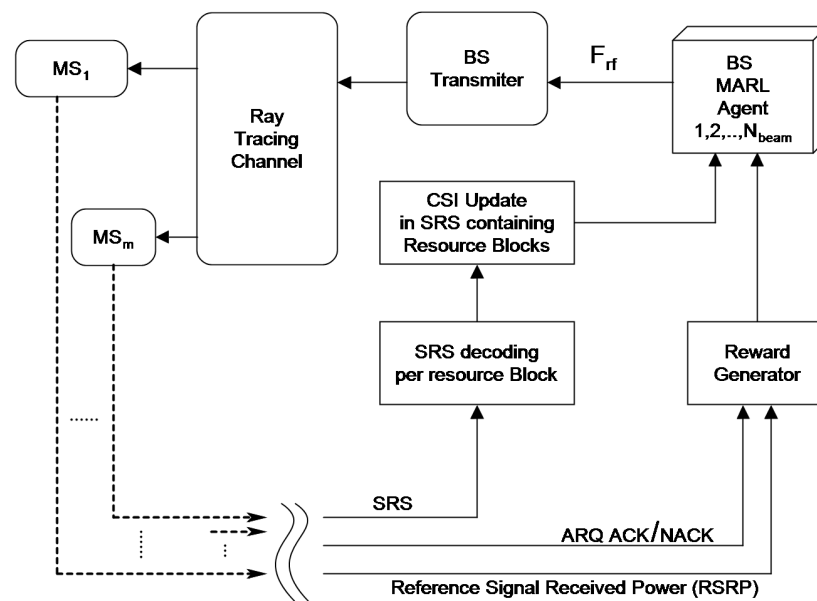


Figure 5. Block diagram of the proposed MARL-based beamforming codebook design.

RL is a type of machine learning where an agent learns to make decisions by performing actions and receiving feedback from the environment in the form of rewards or penalties. This learning process involves trial-and-error, where the agent seeks to maximize cumulative rewards over time. Unlike supervised learning, where the correct actions are provided by a teacher, RL agents must discover optimal actions through interaction with the environment. This approach is formalized through frameworks like Markov decision processes, which provide a mathematical foundation for modeling decision-making scenarios involving uncertainty and delayed rewards. RL has been successfully applied in various domains, including robotics, game playing, and autonomous control systems, due to its ability to handle complex, dynamic environments. Recent advancements in deep reinforcement learning, which combines RL with deep neural networks, have led to significant breakthroughs, such as achieving human-level performance in complex games like Go and Atari. This field continues to evolve, offering promising solutions for real-world applications where adaptive, intelligent behavior is required.

The proposed MARL algorithm builds upon the Wolpertinger Architecture [14], following a similar approach to that described in [10]. The Wolpertinger architecture in reinforcement learning is a sophisticated framework designed to address the challenges associated with large discrete action spaces. This architecture integrates the strengths of deep reinforcement learning with efficient action exploration mechanisms. It operates by embedding the discrete action space into a continuous representation, allowing for more efficient searching and action selection. The Wolpertinger policy effectively reduces the computational complexity by narrowing down potential actions through a k-nearest neighbors approach, ensuring that the most promising actions are evaluated without exhaustive searching. This makes it particularly useful in applications such as content caching at the edge of wireless networks and optimizing beamforming in MIMO systems, where traditional methods struggle with the vast number of possible actions. By efficiently managing large action spaces, the Wolpertinger architecture enhances the scalability and applicability of reinforcement learning in complex, real-world scenarios, ensuring robust and effective decision making.

The Wolpertinger architecture adapts the DDPG, originally crafted for continuous action spaces, to function within a discrete action space through the utilization of a K-nearest neighbor (KNN) classifier. To address non-stationary environment issues in a multi-agent RL system with continuous action spaces, the MADDPG offers a solution. MADDPG achieves this through centralized training and decentralized execution. To accommodate a discrete action space in MARL, an improvisation on MADDPG is made by implementing each agent in MARL using the Wolpertinger architecture. Thus, the proposed MARL essentially embodies MADDPG, with each agent designed to adhere to the Wolpertinger architecture.

The proposed beam learning problem presents a significant challenge due to the large number of possible actions. For instance, considering a base station with 32 antennas, 3-bit phase shifters, and 4 RF channels, each agent faces around 8^{32} potential actions. This complexity is further compounded with additional antennas and higher-resolution phase shifters, rendering conventional deep Q-network frameworks impractical.

Additionally, multi-agent deep Q-networks suffer from instability and renders the environment non-stationary. To overcome these limitations, the Wolpertinger architecture is introduced, offering a solution for navigating spaces with extensive sets of discrete actions [14]. This architecture, rooted in the actor-critic framework, is trained using the DDPG algorithm [15]. Notably, the Wolpertinger architecture incorporates a KNN classifier, enabling DDPG to effectively handle tasks with discrete, finite, yet exceptionally high-dimensional action spaces. Below, a concise overview of the key components of the Wolpertinger architecture is provided.

Actor Networks: The actor maps states from the observation space to actions, serving as a function approximator for this mapping process. Since the actions obtained from the actor fall into a continuous action space, the predicted action may not align perfectly

with the action space of the problem. Therefore, this prediction is referred to as a proto action and is quantized by a KNN classifier to obtain an action available in the discrete action space.

KNN search: KNN search is employed to determine the nearest neighbor of the proto action within the discrete action space. This algorithm utilizes the L_2 distance, also known as a squared Euclidean distance, as a metric to identify the closest vector to the proto action. In essence, the KNN algorithm assesses the spatial proximity of the proto action to the available discrete actions, helping to quantize and align the predicted action with the specific options within the discrete action space.

Exploration noise process: Noise helps agents explore the environment more effectively by injecting randomness into their actions. Exploration is essential in reinforcement learning to discover new states and actions that can lead to better policies. Without exploration, agents might get stuck in suboptimal policies. The noise added to actions is often generated from a stochastic process, such as a Gaussian distribution, Ornstein–Uhlenbeck process, or other types of noise sources. Ornstein–Uhlenbeck process is used in this work to generate noise that is added to the actions of an agent. This noise has the property of being temporally correlated, which means that it tends to stay close to its current value over short time intervals, mimicking the behavior of real-world systems. The peak noise magnitude needs to be such that after adding it to the action in element-wise manner produces a resulting magnitude large enough to cover the full range of a phase shifter array.

Critic Networks: The critic network functions as a Q function, accepting both the state and action inputs and generating the anticipated Q value for the specific state–action combination. Given that the KNN function yields k potential actions, the critic network evaluates k distinct state–action pairs (with a shared state), ultimately pinpointing the action that attains the highest Q value among them.

Target Networks: The target network is a separate neural network that mirrors the actor network. Its parameters are updated less frequently, providing a stable target for the training process. The periodic update of the target network’s parameters enhances the stability and convergence of the learning process, leading to improved training efficiency and more accurate action value estimations.

In this scope, the input (state), outputs (action) and reward process of the MARL algorithm are defined.

State: State comprises of the concatenated vector of the phases of an all-phase shifter at time t and the average normalized envelop of the channel estimate obtained through procedures given in step 1 through step 8 of the major steps mentioned in Section 3.1 of Section 3.

Action: Action comprises element-wise changes of all the phases in the state vector at time t.

Reward: Reward design is pivotal for shaping effective RL policies, efficiently achieving goals, and avoiding unintended behaviors. The reward function provides feedback, guiding the agent to learn optimal behaviors by reinforcing actions that yield higher rewards and discouraging those that do not. Proper reward modeling ensures that the agent learns efficiently and effectively, aligning its actions with the desired outcomes

3.3. Data Preprocessing

The SRS provides the BS with comprehensive channel information across the entire bandwidth. Utilizing this information, the BS optimizes resource allocation, giving preference to areas with superior channel quality over other bandwidth segments. In this proposed work, emphasis was placed on a central cluster consisting of four resource blocks (RBs), each encompassing a bandwidth of 180 kHz. Within each RB, 12 subcarriers are positioned at 15 kHz intervals, resulting in a combined bandwidth of 720 kHz. A frequency-domain vector comprising 48 complex numbers is derived through channel estimation across this contiguous frequency range. Given that only a narrow band of the entire spec-

trum is required for the proposed algorithm, achieving high SNR for SRS transmission is feasible.

For further analysis, this complex vector is transformed into its magnitude and then downscaled by a factor of 2, resulting in a real-valued vector comprising 24 elements. To ensure consistency, in the subsequent preprocessing stage, this 24-element vector is normalized by dividing it by its maximum value. This procedure is iterated for each of the four simultaneous users, producing four channel vectors.

The other part of the state vector input consists of the phases of the phase shifter network for a particular RF chain of length N_{BS} which is 32 in this case. This is also normalized by the maximum absolute value of the phase vector. Here, four such phase shifter vectors are obtained for four RF chains.

3.4. Proposed Algorithms

This section presents two algorithms for the online learning of the beam codebook and a method for user clustering.

3.4.1. Proposed User Clustering through Beam Sweeping

A straightforward clustering approach based on a simple beam sweeping technique is proposed. Users with similar channel characteristics are grouped together and served by a single beam. This method allows the learning of interference-free arbitrary beams using multiple agents within the MARL framework. Additionally, clustering divides the complex task of finding beams across the entire azimuth into parallel sub-tasks, making it more manageable and efficient for developing a multi-user, multi-beam beamforming codebook. This simplification streamlines the acquisition process of a codebook comprising multiple beams. In this work, user clustering is proposed as part of the initial access procedure. The system begins with a traditional beam sweeping initial access procedure, as described in Section 3.1 of Section 3, and updates the codebook as learning progresses. A new UE requesting channel access is assumed to fall into a cluster corresponding to the initial access beam index aligned with that UE. Once sufficient learning is achieved through MARL, the codebook remains fixed in that learned state until further macro changes in the environment occur. The number of UE channel clusters and the number of beams are the same for both the untrained and trained codebooks because, in the proposed MARL algorithm, each beam in a learned codebook corresponds to one beam in the initial codebook. Hence, irrespective of the state of learning of the MARL agents, a new user is assigned to the cluster corresponding to the beam index it obtains through the initial access procedure.

3.4.2. Proposed Reward Function

The proposed reward processing is detailed in Algorithm 1. The reward function is designed to satisfy two goals, namely maximizing the average beamforming gain, in turn maximizing the sum rate of the system and reducing the inter-beam interference. In the proposed approach, the end-to-end system is implemented, where the ARQ signal is sent by the UE to the BS based on whether a frame was received correctly or not, and the received ARQ is also used as an input for reward modeling, addressing concerns that RSRP alone may provide a misleading indication of beamforming gain maximization in a multi-beam system with interference.

Algorithm 1 Proposed reward function

```

1: Initialize dynamic threshold for RSRP,  $Th_t = 0$ .
2: Observe RSRP feedback from UE,  $RSRP_t$ 
3: Observe ACK/NACK (True/False) from UE.  $ARQ_t$ 
4: if  $RSRP_t > Th_t$  and  $ARQ_t = True$  then
5:    $Reward_t = 2$ ;
6:    $Th_t = RSRP_t$ ;
7: else if  $RSRP_t \leq Th_t$  and  $RSRP_t > RSRP_{t-1}$  and  $ARQ_t = True$  then
8:    $Reward_t = 1$ ;
9: else if  $RSRP_t \leq Th_t$  and  $RSRP_t \leq RSRP_{t-1}$  and  $ARQ_t = True$  then
10:   $Reward_t = 0$ ;
11: else if  $RSRP_t > Th_t$  and  $ARQ_t = False$  then
12:   $Reward_t = -1$ ;
13: else
14:   $Reward_t = -2$ ;
15: end if

```

3.4.3. Proposed MARL-Based Codebook Learning Function

Proposed MARL-based beam codebook learning for N agents is given in Algorithm 2. The input to each actor network is the corresponding state. The state is the concatenation of the 24 length channel vector given as N_{CH} and 32 length phase vector which equals N_{BS} . Thus, the length of state vectors are $N_{CH} + N_{BS}$, which is 56 in this case. The output of the actor networks also comprises the predicted phase update vectors which are of length N_{BS} , i.e., 32. The actor network includes a pair of hidden layers, each containing $10 \times (N_{CH} + N_{BS})$ neurons equating to 560. These layers are subsequently activated using rectified linear units (ReLU). The outcome of the actor network stands as the anticipated action. This outcome is then passed through hyperbolic tangent (tanh) activations, which are scaled by π .

Thus, the length of the input of each critic network for a four agents can be given as $(4 \times (N_{CH} + N_{BS}) + 4 \times N_{BS})$, i.e., 336 in this case. The output of the critic network is the predicted Q value, which is a real valued scalar. Hence, output dimension of critic network is 1. The critic network is composed of two hidden layers, each layer containing $5 \times (4 \times (N_{CH} + N_{BS}) + 4 \times N_{BS})$, i.e., 1680 neurons. Following this, ReLU activations are applied to these layers.

Hyper parameter for the MARL is given in Table 2.

Table 2. Hyper parameter table for MARL.

Optimizer	ADAM
Learning rate	0.01
Target soft update parameter	0.95
Replay buffer size	12,288
Batch size	1024
No of samples added to replay buffer before each network update	100

Algorithm 2 Proposed MARL-based beam learning

- 1: Initialize actor networks, critic networks with random weights
- 2: Initialize target networks and with the weights of actor and critic networks
- 3: Initialize the replay memory D , minibatch size B , discount factor γ
- 4: Initialize reward processing algorithm
- 5: Initialize N beams with beam-steering codebook in *procedure 2* of Section 3 for N clusters and N agents.
- 6: Initialize a random process N for action exploration
- 7: For each agent, initialize random initial beamforming vector as state, x .
- 8: **for** $t = 1$ to *max-episode-length* **do**
- 9: For each agent i , select proto action $a_i = \mu_{\theta_i} + N_t$ w.r.t. the current policy and exploration.
- 10: for each agent i , quantize proto action to valid beamforming vector with KNN for $k = 1$.
- 11: Execute action $a = (a_1, a_2, \dots, a_N)$ and observe reward r (with Algorithm 1) and new state x'
- 12: Store (x, a, r, x') in replay buffer D
- 13: $x \leftarrow x'$
- 14: **for** agent $i = 1$ to N **do**
- 15: Sample a random minibatch of S samples (x^j, a^j, r^j, x'^j) from D
- 16: Set $y_j = r_i^j + \gamma Q_i^{\mu'}(x^j, a_1^j, \dots, a_N^j)|_{a_k^j = \mu_k^j(o_i^j)}$
- 17: Update critic by minimizing the loss $L(\theta_i) = \frac{1}{S} \sum_j (y_j - Q_i^{\mu}(x^j, a_1^j, \dots, a_N^j))^2$
- 18: Update actor using the sampled policy gradient:
- 19: $\nabla_{\theta_i} J \approx \frac{1}{S} \sum_j \nabla_{\theta_i} \mu_i(o_i^j) \nabla_{a_i} Q_i^{\mu}(x^j, a_1^j, \dots, a_i^j, \dots, a_N^j)|_{a_i = \mu_i(o_i^j)}$
- 20: **end for**
- 21: Update target network parameter for each agent i :
- 22: $\theta_i' \leftarrow \tau \theta_i + (1 - \tau) \theta_i'$
- 23: **end for**

3.5. Numerical Simulation with MARL

In this proposed work, a 120° sector of a cell for simulation purposes is modeled, restricting transmissions within this azimuth range. Although the 4 RF channels can concurrently serve 4 users within this angular space, real-world scenarios typically involve more than four active users. To address this, users with similar channels are served with a single beam. The assignment of each user to a specific beam, whether before or after the MARL-based codebook learning process, is determined through beam sweeping. Consequently, the number of beams in the learned codebook remains consistent with the initial access codebook, which is adjustable for the performance assessment. Figure 6 illustrates the radiation pattern for one such codebook with nine beams, showcasing variations for different quantization bits.

In the proposed MARL algorithm, the number of agents corresponds to the number of beams utilized in the sector. This configuration effectively breaks down the task of selecting a beam from a large set into finding a single beam within a smaller subset, thereby enhancing the efficiency of the codebook learning process. An additional and significant benefit of employing one agent per beam is the ability to identify optimal non-interfering beams within the sector, even in nLOS scenarios. Each agent in the MARL algorithm strives to maximize the individual beamforming gain while minimizing interference with other agents, as reflected in the reward processing outlined in Algorithm 1.

Upon completion of the learning phase, the acquired codebook becomes readily deployable within the initial access procedure. Users can now be efficiently served using the learned codebook, rendering the traditional matched filter-based beam codebook obsolete. This transition marks a significant advancement in the efficiency and adaptability of beamforming techniques, as the learned codebook optimally caters to the dynamic

needs and complexities of the communication environment without relying on pre-defined beam patterns.

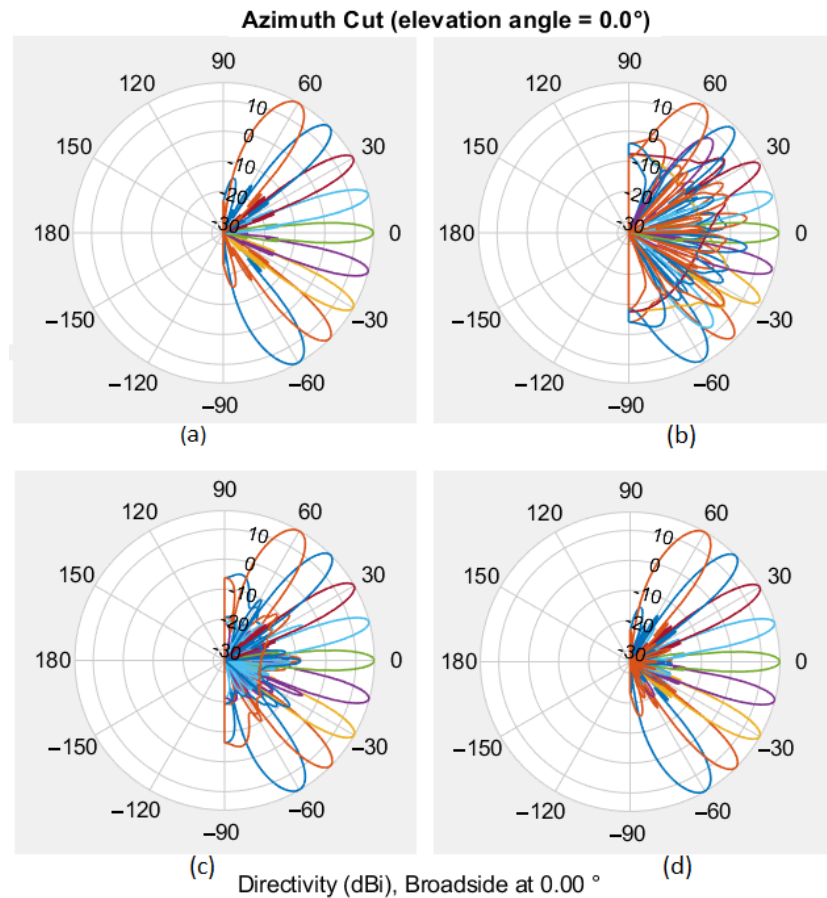


Figure 6. Radiation pattern for 32-element ULA antenna with number of phase shifter quantization bits (a) 0 bits, (b) 2 bits, (c) 4 bits (d) 6 bits. Colored lines represent distinct radiation patterns corresponding to specific steering angles in the azimuth plane.

This learned codebook is valid until there is no significant change in terms of macro structures within the sector. Although such time will only be there occasionally, in the case of such large changes in the structure or the replacement of the BS, learning has to be initiated again for all the beams.

Next, the analog beamforming codebook selection for the UE is carried out. In this work, a conventional beamforming codebook tailored for the UE is employed. The process of selecting beams from the codebook for the UE involves a standard beam search procedure, encompassing steps such as sounding, measurement, and feedback.

In the final step, the baseband beamforming vector (F_{BB}) at the BS is calculated. This computation follows the procedure outlined in [4]. In this process, the BS formulates its zero-forcing digital precoder F_{BB} based on the quantized channel feedback received from the UE. Due to the utilization of RF beamforming and the presence of sparse mmWave channels, it is anticipated that the effective MIMO channel will be well-conditioned [16,17]. This favorable channel condition enables the utilization of a straightforward multi-user digital beamforming strategy like zero-forcing, which can achieve a performance close to the optimal level [18]. The algorithm for obtaining the baseband beamforming vector F_{BB} is detailed in the second stage of the procedure presented in [4].

4. Results and Discussion

This section highlights the performance of the proposed MARL technique. This is demonstrated through a series of experiments where the agents are trained using the

parameters outlined in Table 2. The effectiveness of the proposed MARL-based approach to learning beam codebooks across various scenarios is assessed. Unlike the process of creating a codebook for single-user MIMO systems, as discussed in [10], acquiring a beam codebook with multiple beams for multiuser MIMO involves not only learning the codebook but also identifying optimal combinations from a wide range of potential beamforming vectors.

Figure 7 illustrates the average beamforming gain relative to the number of beams contained within the codebook, specifically in the LOS scenario. In this scenario, the BS employs a uniform linear array with isotropic elements, oriented in a back-baffled configuration. The graph demonstrates a consistent upward trend in average beamforming gain with an increasing number of beams. In line with the observations made in [10], the result in Figure 7 demonstrates nearly equivalent performance to a classical 32-beam beamsteering codebook when employing only 6 beams. Notably, the proposed approach not only matches but also surpasses the performance of [10]. This trend persists as the solution, employing 8 beams, consistently outperforms the 32-beam classical beamsteering codebook while also exceeding the capabilities of [10]. This achievement is particularly significant because it addresses a multiuser scenario with four co-channel users and a nonzero interference probability, representing an improvement over the single-user MIMO configuration in [10]. It is important to note that traditionally, single-user and multiuser codebooks are identical, meaning multiple beams from the same codebook are used for multiuser MIMO. For this comparison, the single-user codebook learned by deep reinforcement learning (DRL) in [10] is utilized and extended it to the multiuser case.

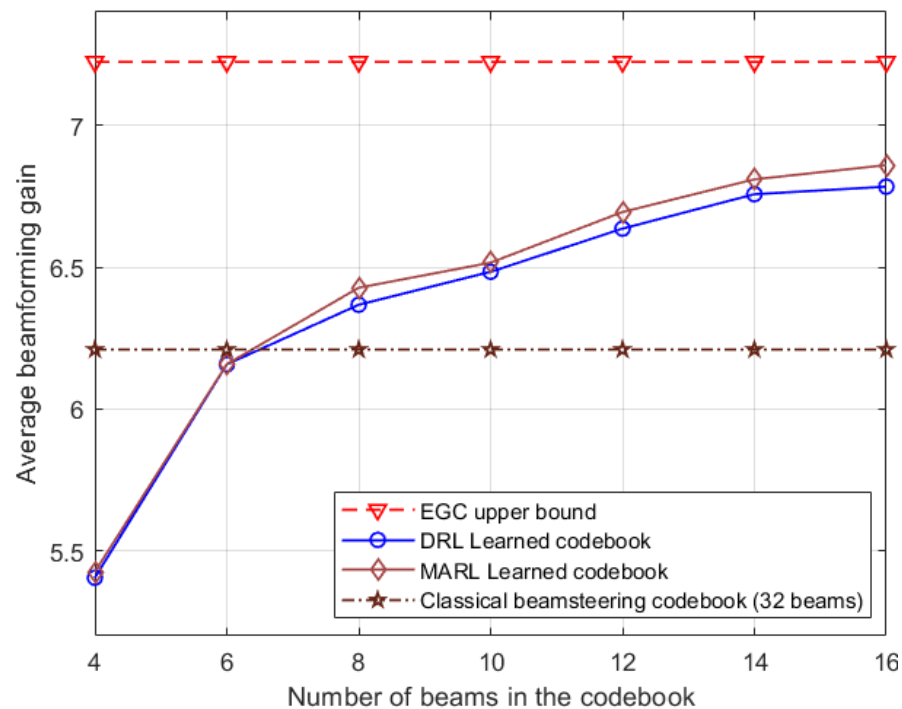


Figure 7. Comparison of average beamforming gain versus number of beams in the learned codebook in LOS with 32-element ULA.

This demonstrates that, even with a simple clustering scheme, which eliminates the need for extensive random beam searches with a large set of matched beams, the proposed MARL approach, utilizing the designed reward function, can effectively understand the wireless environment. The methodology illustrates its capability to dynamically adjust beam configurations based on user distributions and environmental topography, significantly reducing the beam training overhead in a massive MIMO system by minimizing the number of beams required for effective communication in the learned codebook. This ad-

justment effectively mitigates interference within densely populated urban environments, leading to notable performance improvements.

Furthermore, users were strategically placed in nLOS areas within the scenario. The simulation under nLOS conditions highlights the superiority of the proposed MARL system compared to traditional beam codebooks and those proposed in [10]. In this scenario, MARL outperforms the 32-beam classical beamsteering codebook and [10] with just 4 beams. Given that only reflected paths of the channel are available under nLOS conditions, this improvement underscores the adaptability of the MARL system to varying environments. The simulation results for MARL under nLOS conditions are depicted in Figure 8.

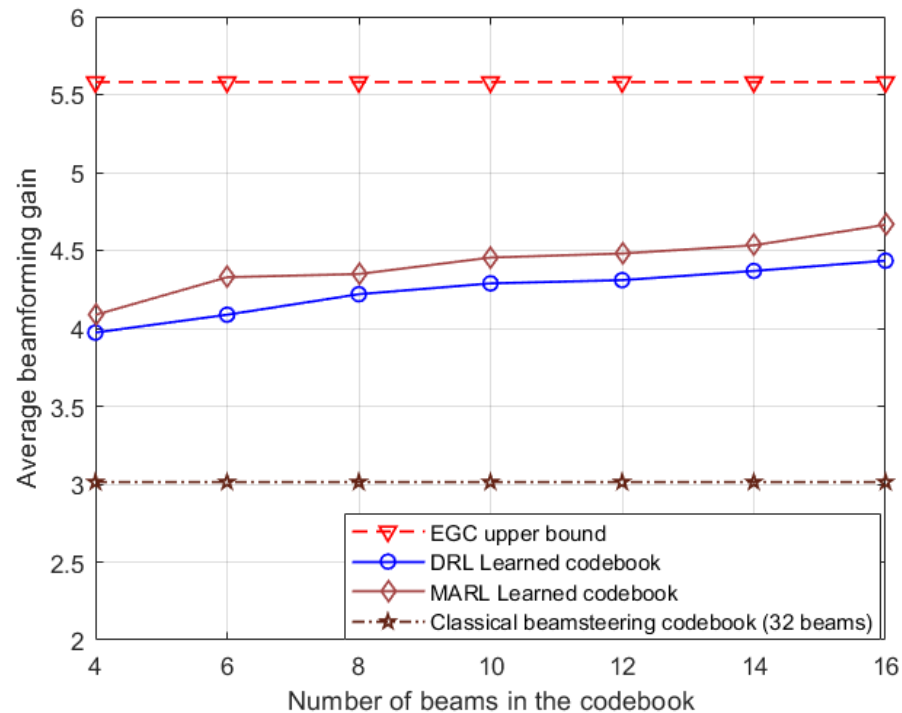


Figure 8. Comparison of average beamforming gain versus number of beams in the learned codebook in nLOS with 32-element ULA.

Figure 8 clearly shows that the learned codebook with only 4 beams in nLOS outperforms the traditional 32-beam beamsteering codebook. Additionally, as the number of beams in the learned codebook increases, the performance of the proposed algorithm approaches that of equal gain combining (EGC). This is notable because the EGC upper bound is typically achievable only in a beamforming setup with perfect CSI and continuous phase shifters. In contrast, the proposed approach uses only a partial, imperfect uplink CSI estimated at the base station as a fingerprint of the channel and a marker of the UE's position, along with RSRP for reward modeling. This demonstrates the practicality and effectiveness of the proposed method, even with limited and imperfect channel information.

This advantage is particularly evident in scenarios with blockages, where user signals rely on reflections to reach the access point. The proposed solution demonstrates its adaptability to the propagation environment by adjusting the beam pattern, thereby capturing signals from multiple directions and gaining more power.

This ability to dynamically adjust and optimize the beam pattern in response to environmental conditions not only improves the performance in nLOS scenarios but also highlights the potential of MARL-based hybrid beamforming in real-world applications where blockages and reflections are common. By leveraging the capabilities of MARL, the system can efficiently manage complex and variable propagation environments, ensuring high data rates and reliable connectivity even under challenging conditions.

The performance of MARL is contrasted against various hybrid precoding techniques including manifold optimization (MO) [19], sparse orthogonal matching pursuit (SOMP) algorithm [5], and the two-stage hybrid beamforming (TS-HB) algorithm [4]. Notably, MO and SOMP were initially proposed for single-user scenarios, but for comparison, these algorithms are adjusted to the multi-user context by adopting the interference cancellation strategy outlined in [9]. In the simulation plot for no interference, the outcomes of fully-digital beamforming and combining are traced. This approach effectively eliminates interference, serving as a reference point in the evaluations.

Figure 9 shows a comparative analysis of the achievable sum-rate performance of the algorithms across various SNR levels, assuming the perfect knowledge of CSI and the available array response. The system parameters are set as follows: each BS has 32 antennas ($N_{BS} = 32$), each UE has 4 antennas ($N_M = 4$), synthetic noise with an SNR of 20 dB, and phase shifter quantization is performed with 3 bits ($r = 3$). The multiuser environment consists of four simultaneous UEs ($M = 4$), each with five paths ($L = 5$). This consistent configuration is maintained across all algorithms to ensure a fair comparison.

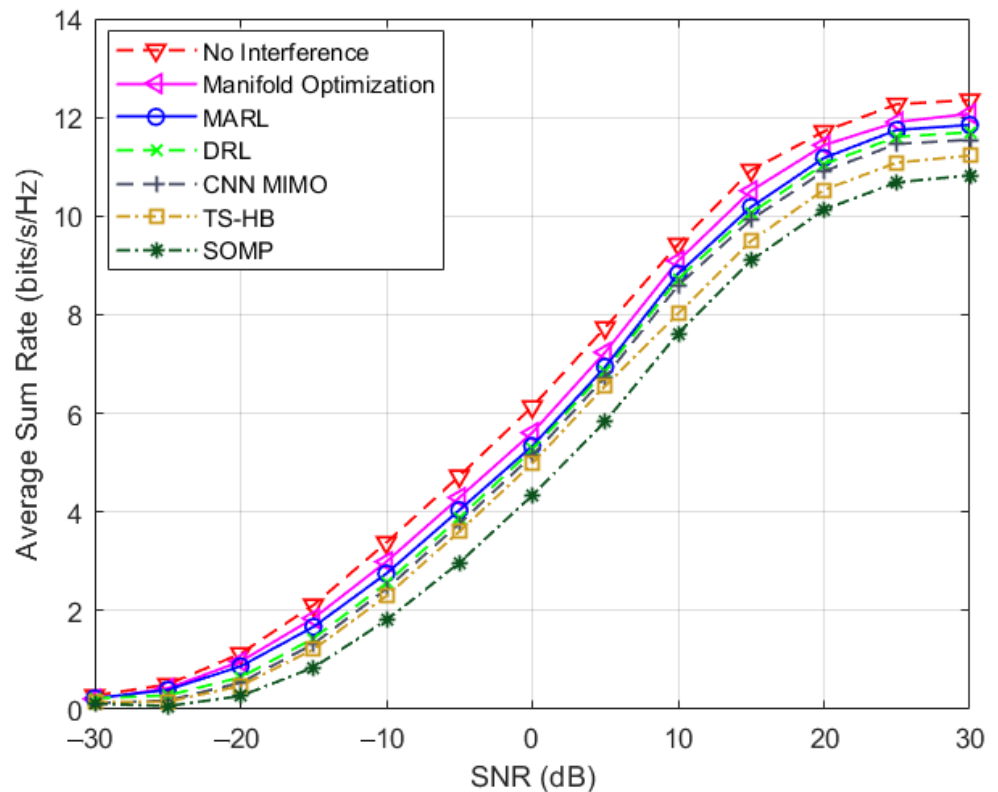


Figure 9. SNR vs. sum-rate comparison for perfect CSI and array response ($N_{BS} = 32$, $N_M = 4$, $r = 3$, $M = 4$, $L = 5$).

For reference, fully digital beamforming and the MO algorithm, known for their near-optimal analog and baseband precoders, have been included. Notably, the performance of the proposed MARL approach closely matches that of the MO algorithm, consistently achieving the highest sum-rate among all the algorithms. As evident from the simulation results with perfect CSI shown in Figure 9, in a typical hybrid beamforming scenario with ideal channel conditions and minimal interference, the sum-rate begins to plateau as the system nears its maximum spatial multiplexing capability. At this point, the effects of quantization noise and hardware imperfections become more pronounced. Additionally, the $\log_2(1 + \text{SNR})$ term grows more slowly, approaching a logarithmic limit. This indicates that the incremental gain in capacity (and thus sum-rate) diminishes as SNR increases.

A perfect CSI and array response data are typically unavailable. To assess the impact of corrupted CSI and array response data on the performance of the proposed system,

complex AWGN noise was added to both the array response and channel matrix. Figure 10 presents a comparative analysis of the sum-rate performance of the algorithms across varying SNR levels, assuming noisy measurements.

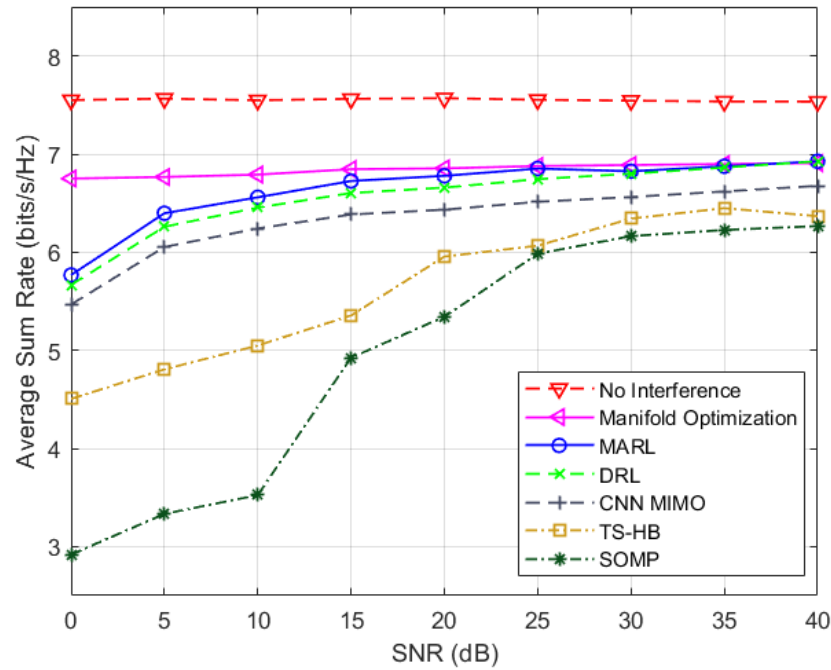


Figure 10. SNR vs. sum-rate comparison with noisy measurement of CSI and array response ($N_{BS} = 32, N_M = 4, r = 3, M = 4, L = 5$).

Compared to DRL, SOMP, TS-HB, and CNN-MIMO [9], MARL demonstrates a superior performance. While SOMP was originally designed for single-user scenarios, it has been adapted for multi-user contexts in this proposed work. Both SOMP and TS-HB require the input of feasible sets F and W , which are the array response sets. Therefore, the precision of these feasible sets significantly affects the performance of SOMP and TS-HB, as it relies on the accuracy of both the channel matrices and array response sets. Corrupt channel and array response data can result in a reduction in the average sum rate at low SNR, as shown in Figure 10. While CNN-MIMO does not require feasible sets of beamforming vectors during the prediction stage, they are necessary during the training stage to acquire labels. Overall, these findings highlight the robustness of the proposed MARL-based approach for downlink RF beamforming codebook design in multi-beam and multi-user MIMO systems.

To highlight the advantages and differences of each of the studied schemes over one another, a table is formed as given in Table 3. As shown in Table 3, the proposed MARL-based beamforming codebook design does not require any difficult-to-obtain or unavailable data during either the training or evaluation phases. The results discussed so far also demonstrate the robustness of MARL against noisy measurements.

Table 3. Comparison table for the different attributes of the studied methods.

	Perfect CSI Not Required	CSI Feedback Not Required	Antenna Array Response Not Required	Multisuser Interference Avoidance in RF Beams	DOA /DOD Not Required	Online Training	Non Iterative Method
MO [19]	×	×	×	✓	×	×	×
TS-HB [4]	×	×	×	×	✓	×	×
SOMP [5]	×	×	×	✓	×	×	×
CNN-MIMO [9]	✓	×	×	✓	✓	×	✓
DRL [10]	✓	✓	✓	×	✓	✓	✓
MARL	✓	✓	✓	✓	✓	✓	✓

5. Conclusions

The proposed MARL-based codebook learning design demonstrated the capability to dynamically adjust beam configurations based on user distributions and environmental topography, effectively mitigating the interference within densely populated urban environments. Leveraging uplink CSI and RSRP for reward modeling, this method reduces the need for extensive random beam searches and minimizes the number of beams required in the learned codebook, significantly reducing the beam training overhead. The proposed system outperformed traditional 32-beam beamsteering codebooks in both LOS and nLOS conditions with fewer beams by using a simple clustering scheme and a designed reward function that facilitated interference-free beam learning. MARL-based approach showcased its adaptability by efficiently managing complex propagation environments, ensuring high data rates and reliable connectivity even in challenging conditions. Comparative analysis highlighted that the proposed MARL approach achieved performance close to the optimal level, providing robust and efficient beamforming solutions for multiuser MIMO systems. This research contributes to the field by offering a practical and efficient solution for beamforming in massive MIMO systems, paving the way for enhanced multiuser communication in future wireless networks. However, the approach is not without limitations. One major limitation is the complexity of the MARL framework, which may pose scalability issues as the number of users and beams increases, potentially requiring higher computational resources. Another limitation is the user clustering method, which requires the complete reclustering and retraining of the MARL agents after significant environmental changes. Additionally, the simulations were conducted in controlled environments using ray-tracing models, and real-world deployments may present unforeseen challenges that could impact the system's performance and adaptability. Further research is needed to address these limitations and enhance the robustness of the proposed method in diverse and dynamic real-world scenarios.

Author Contributions: Data curation, M.B. and D.D.M.; Formal analysis, K.K.S.; Investigation, M.B. and K.G.; Methodology, M.B. and K.K.S.; Project administration, J.I.; Software, D.D.M.; Supervision, K.K.S., K.G. and J.I.; Validation, M.B. and K.K.S.; Visualization, M.B.; Writing—review and editing, D.D.M., K.G. and J.I. All authors have read and agreed to the published version of the manuscript.

Funding: This research received no external funding.

Institutional Review Board Statement: Not applicable.

Informed Consent Statement: Not applicable.

Data Availability Statement: The data presented in this study were generated dynamically using simulation models developed by the authors. The map used for the ray tracing channel simulation is sourced from <https://www.openstreetmap.org>, which provides crowd-sourced map data worldwide. As the data are not derived from any other publicly available dataset and were produced through the methods outlined in the manuscript, they are not available in a public repository.

Conflicts of Interest: The authors declare no conflicts of interest.

References

1. Molisch, A.F.; Ratnam, V.V.; Han, S.; Li, Z.; Nguyen, S.L.H.; Li, L.; Haneda, K. Hybrid beamforming for massive mimo: A survey. *IEEE Commun. Mag.* **2017**, *55*, 134–141. [[CrossRef](#)]
2. Pirapaharan, K.; Prabhashana, W.H.S.C.; Medaranga, S.P.P.; Hoole, P.R.P.; Fernando, X. A New Generation of Fast and Low-Memory Smart Digital/Geometrical Beamforming MIMO Antenna. *Electronics* **2023**, *12*, 1733. [[CrossRef](#)]
3. Singh, J.; Ponnuru, S.; Madhow, U. Multi-gigabit communication: The adc bottleneck. In Proceedings of the 2009 IEEE International Conference on Ultra-Wideband, Vancouver, BC, Canada, 9–11 September 2009; pp. 22–27.
4. Alkhateeb, A.; Leus, G.; Heath, R.W. Limited Feedback Hybrid Precoding for Multi-User Millimeter Wave Systems. *IEEE Trans. Wirel. Commun.* **2015**, *14*, 6481–6494. [[CrossRef](#)]
5. Ayach, O.E.; Rajagopal, S.; Abu-Surra, S.; Pi, Z.; Heath, R.W. Spatially sparse precoding in millimeter wave mimo systems. *IEEE Trans. Wirel. Commun.* **2014**, *13*, 1499–1513. [[CrossRef](#)]

6. Alkhateeb, A.; El Ayach, O.; Leus, G.; Heath, R.W. Hybrid precoding for millimeter wave cellular systems with partial channel knowledge. In Proceedings of the 2013 Information Theory and Applications Workshop (ITA), San Diego, CA, USA, 10–15 February 2013; pp. 1–5.
7. Kim, C.; Kim, T.; Seol, J.-Y. Multi-beam Transmission Diversity with Hybrid Beamforming for MIMO-OFDM Systems. In Proceedings of the 2013 IEEE Globecom Workshops (GC Wkshps), Atlanta, GA, USA, 9–13 December 2013; pp. 61–65.
8. Love, D.J.; Heath, R.W. Equal Gain Transmission in Multiple-input Multiple-output Wireless Systems. *IEEE Trans. Commun.* **2003**, *51*, 1102–1110. [[CrossRef](#)]
9. Elbir, A.M.; Papazafeiropoulos, A.K. Hybrid Precoding for Multiuser Millimeter Wave Massive MIMO Systems: A Deep Learning Approach. *IEEE Trans. Veh. Technol.* **2020**, *69*, 552–563. [[CrossRef](#)]
10. Zhang, Y.; Alrabeiah, M.; Alkhateeb, A. Reinforcement Learning of Beam Codebooks in Millimeter Wave and Terahertz MIMO Systems. *IEEE Trans. Commun.* **2022**, *70*, 904–919. [[CrossRef](#)]
11. Chen, Z.; Cao, Z.; He, X.; Jin, Y.; Li, J.; Chen, P. DoA and DoD Estimation and Hybrid Beamforming for Radar-Aided mmWave MIMO Vehicular Communication Systems. *Electronics* **2018**, *7*, 40. [[CrossRef](#)]
12. Li, Z.; Chen, T. Hybrid Beamforming for Multi-User Millimeter-Wave Heterogeneous Networks. *Electronics* **2022**, *11*, 4221. [[CrossRef](#)]
13. Samimi, M.K.; MacCartney, G.R.; Sun, S.; Rappaport, T.S. 28 Ghz Millimeter-wave Ultrawideband Small-scale Fading Models in Wireless Channels. In Proceedings of the 2016 IEEE 83rd Vehicular Technology Conference (VTC Spring), Nanjing, China, 15–18 May 2016; pp. 1–6.
14. Dulac-Arnold, G.; Evans, R.; van Hasselt, H.; Sunehag, P.; Lillicrap, T.; Hunt, J.; Mann, T.; Weber, T.; Degris, T.; Coppin, B. Deep Reinforcement Learning in Large Discrete Action Spaces. *arXiv* **2015**, arXiv:1512.07679. <https://doi.org/10.48550/arXiv.1512.07679>.
15. Lillicrap, T.P.; Hunt, J.J.; Pritzel, A.; Heess, N.; Erez, T.; Tassa, Y.; Silver, D.; Wierstra, D. Continuous Control with Deep Reinforcement Learning. *arXiv* **2019**, arXiv:1509.02971. <https://doi.org/10.48550/arXiv.1509.02971>.
16. Smith, P.J.; Neil, C.; Shafi, M.; Dmochowski, P.A. On the Convergence of Massive MIMO Systems. In Proceedings of the 2014 IEEE International Conference on Communications (ICC), Sydney, NSW, Australia, 10–14 June 2014; pp. 5191–5196.
17. Yang, H.; Marzetta, T.L. Performance of Conjugate and Zero-forcing Beamforming in Large-scale Antenna Systems. *IEEE J. Sel. Areas Commun.* **2013**, *31*, 172–179. [[CrossRef](#)]
18. Yoo, T.; Goldsmith, A. On the Optimality of Multiantenna Broadcast Scheduling Using Zero-forcing Beamforming. *IEEE J. Sel. Areas Commun.* **2006**, *24*, 528–541.
19. Yu, X.; Shen, J.-C.; Zhang, J.; Letaief, K.B. Alternating Minimization Algorithms for Hybrid Precoding in Millimeter Wave MIMO Systems. *IEEE J. Sel. Top. Signal Process.* **2016**, *10*, 485–500. [[CrossRef](#)]

Disclaimer/Publisher’s Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.