

Machine-Learning Beam Tracking and Weight Optimization for mmWave Multi-UAV Links

Hsiao-Lan Chiang, Kwang-Cheng Chen, *Fellow, IEEE*, Wolfgang Rave, Mostafa Khalili Marandi, and Gerhard Fettweis, *Fellow, IEEE*

Abstract—Millimeter-wave (mmWave) hybrid analog-digital beamforming is a promising approach to satisfy the low-latency constraint in multiple unmanned aerial vehicles (UAVs) systems, which serve as network infrastructure for flexible deployment. However, in highly dynamic multi-UAV environments, analog beam tracking becomes a critical challenge. The overhead of additional pilot transmission at the price of spectral efficiency is shown necessary to achieve high resilience in operation. An efficient method to deal with high dynamics of UAVs applies machine learning, particularly Q-learning, to analog beam tracking. The proposed Q-learning-based beam tracking scheme uses current/past observations to design rewards from environments to facilitate prediction, which significantly increases the efficiency of data transmission and beam switching. Given the selected analog beams, the goal of digital beamforming is to maximize the SINR. The received pilot signals are utilized to approximate the desired signal and interference power, which yield the SINR measurements as well as the optimal digital weights. Since the selected analog beams based on the received power do not guarantee the hybrid beamforming achieving the maximization SINR, we therefore reserve additional analog beams as candidates during the beam tracking. The combination of analog beams with their digital weights achieving the maximum SINR consequently provides the optimal solution to the hybrid beamforming.

Index Terms—UAV communication, mmWave, machine learning, Q-learning, beam tracking, hybrid beamforming, weight optimization, highly dynamic environment.

I. INTRODUCTION

Applications of unmanned aerial vehicles (UAVs) in civil uses become popular in recent years. For example, post-disaster use. A UAV is capable of carrying a network device as an access point that uses an intelligent reflecting surface with beamforming to *reflect* incident signals [1]. A group of UAVs forms an aerial radio access network (aerial-RAN), which serves short-term network infrastructure as an independent wireless network or a long-term extension of existing mobile communication networks [2], [3]. An aerial-RAN can perform tasks such as (i) the UAVs together transmit or receive signals from different directions to detect weak signals from victims and (ii) the UAVs separately serve as independent wireless networks to provide a wide range of services, see Fig. 1. In this example, two followers collect data from ground users and then report the information to the lead UAV, which will pass the data to a remote ground anchor node [4], [5], [6].

Hsiao-Lan Chiang and Kwang-Cheng Chen are with the University of South Florida, Tampa, FL, USA.

Wolfgang Rave, Mostafa Khalili Marandi, and Gerhard Fettweis are with the Vodafone Chair Mobile Communications Systems, Technische Universität Dresden, Dresden, Germany.

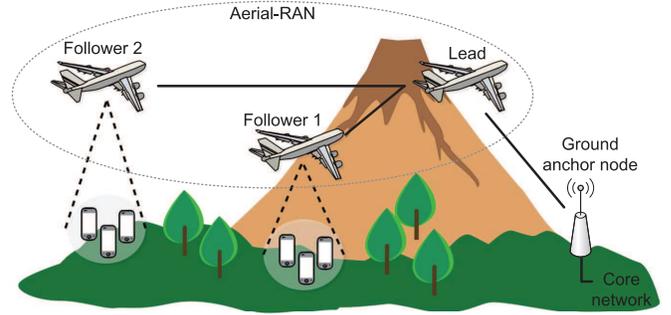


Fig. 1. An example of multi-UAV scenarios. The UAVs are deployed in an area of interest for search and rescue works, where the lead UAV transmits the users' data collected from the followers to the ground anchor node.

Such an operation is often characterized by low-latency and high-resilience constraints. The former is defined as the time to get a response to information sent, while the latter is the ability that provides and maintains an acceptable link quality of services in highly dynamic operations.

Millimeter-wave (mmWave) communication is one of the candidates to satisfy the low-latency requirement due to availability of large chunks of spectrum in unlicensed mmWave frequency bands [7], [8]. Compared with sub-6 GHz communications, mmWave propagation suffers from more severe environmental conditions, such as path loss and a small number of scattering events [9], [10]. In order to improve the data rates and quality of service, beamforming technology for large antenna arrays seems to be a promising approach. At mmWave frequencies, analog beamforming via a passive phased array is taken into account due to cost and power consumption concerns [11], [12], [13]. With more than one analog beamforming vector, linear combinations of multiple analog beamforming vectors with weights of digital beamformers as coefficients provide more degrees of freedom for beamforming designs. Such a beamforming architecture is called hybrid analog-digital beamforming [14], [15].

In hybrid beamforming systems, although both analog and digital beamforming matrices use the same word beamforming, only the former has a specific geometrical meaning in the sense of transmitting or receiving signals towards specific directions in the 3-D space using antenna arrays. In contrast, the digital weights act in the sense of optimum linear combining, given some cost criterion. According to the functions of analog and digital beamforming, hybrid beamforming can be viewed as first converting a MIMO channel matrix (in the spatial domain) into an effective channel (in the angular

domain) using analog beamforming vectors [16], [17]. Then, one can further design the weights of the digital beamformers to linearly combine the analog beamforming vectors based on some optimality criteria. Clearly, the performance of hybrid beamforming is dominated by the analog beam search. In highly dynamic UAV environments with speed up to 100 m/s [18], this challenge (or specifically speaking, analog beam tracking) will be a critical problem.

One of the key performance indicators for dynamic beam tracking could be network resilience [19]. In dynamic environments, the UAVs may have to switch the analog beams rapidly in order to stably provide the acceptable link quality. Given codebooks that consist of candidates for the analog beams, the work in [20] presented a gradient-based algorithm to find a better beam next to the currently used beam, and in [21], the beam tracking problem is formulated as a multi-armed bandit problem. One can also use the extended Kalman filter to recursively track the beams based on the estimated angles of departure and arrival (AoDs/AoAs) [22]. In addition, a conventional object tracking method using reinforcement learning in computer vision [23] has attracted attention and been used in beam tracking [24], [25], [26]. All above-mentioned methods try to find the beam which can achieve an acceptable link quality. However, implementing beam tracking for highly dynamic channels needs a large number of observations (that is, received pilot signals) by sacrificing the spectral efficiency. When we pursue a high-resilient multi-UAV communication, the transmission overhead of pilots is another issue. In this paper, we attempt to strike the balance between the system resilience and efficiency.

To handle the beamforming problem for a time-varying channel, we let the UAVs learn how to interact with the highly dynamic environment during the beam tracking using Q-learning [27], [28]. Q-learning is a model-free reinforcement learning algorithm that uses experience, current measurements, and rewards from the environments to solve the prediction problem without knowing a model of the environment. When applying Q-learning to beam tracking, the crucial problem is to design the reward function based on the noisy observations. Please note that the reward function also influences the experience in Q-learning. Some prior works in [24], [29] used true values of the signal to interference plus noise ratio (SINR) or true values of the received power to define the reward function, which cannot faithfully show the performance of Q-learning-based beam tracking in practical cases. In the proposed method, we use the noisy observations to design the reward function and take current/past observations as arguments in such a way to reduce the pilot overhead.

In the analog beam tracking, the analog beams are selected according to the power of observations.¹ These beams together yield (nearly) the maximum received power. However, the spatial-domain interference from different UAVs could seriously degrade the throughput. Essentially, what really matters to multi-UAV hybrid beamforming is the SINR maximization [30], [31]. To this end, given the selected analog beams, one

can design the corresponding digital weights to maximize the SINR. To obtain the measurements of SINR, we use the received coupling coefficients² (associated with the beams assigned to difference UAVs) to approximate the desired signal and interference power, which facilitates the design of the digital weights. Moreover, it is worth noting that the analog beams leading to the maximum received power may not lead to the maximum SINR [17]. We therefore reserve more candidates for analog beams during the beam tracking. It turns out that the analog beams have to be determined after linear combinations of analog beamforming vectors with the digital weights.

The **contributions** of the proposed method are summarized as follows:

- The proposed method only requires the received coupling coefficients as observations to implement both the analog beam tracking and digital weight optimization. Compared with prior works in the literature which need detailed knowledge, such as channel, we provide a more feasible solution to connect multiple UAVs with low complexity.
- We formulate the beam tracking problem using a Q-learning model and introduce how to use the coupling coefficients to design the rewards. The proposed method can stably track the beams in highly dynamic environments.
- To track the beams in highly dynamic UAV environments, the burden of pilot transmission is inevitable. The proposed beam tracking method uses current and past observations to solve the prediction problem. In such a way, it significantly increases the efficiency of data transmission and beam switching.
- The selected analog beams based on the received power do not ensure that hybrid beamforming achieves the maximization SINR. We manage to reserve additional analog beams as candidates during the beam tracking and then determine which combination of analog beams with their digital weights achieves the maximum SINR. This idea can be simply implemented given the coupling coefficients.

The rest of this paper is organized as follows: Section II describes the multi-UAV beamforming system and time-varying AoDs/AoAs. Section III states the objectives and challenges of the hybrid beamforming problem in highly dynamic environments. To efficiently track the analog beams with limited number of observations, Q-learning is applied to the beam tracking problem for one and multiple links presented in Section IV. Given selected beam pairs, we pursue the corresponding optimal digital weights and the solution is provided in Section V. Simulation results are presented in Section VI, and we conclude our work in Section VII.

We use the following notations throughout this paper.

¹Precisely, the power of observations determines the rewards from environments in Q-learning, and then we use the rewards to find favorable beams.

²A coupling coefficient is a measure of a pair of analog beamforming vectors selected on both sides of the channel [17].

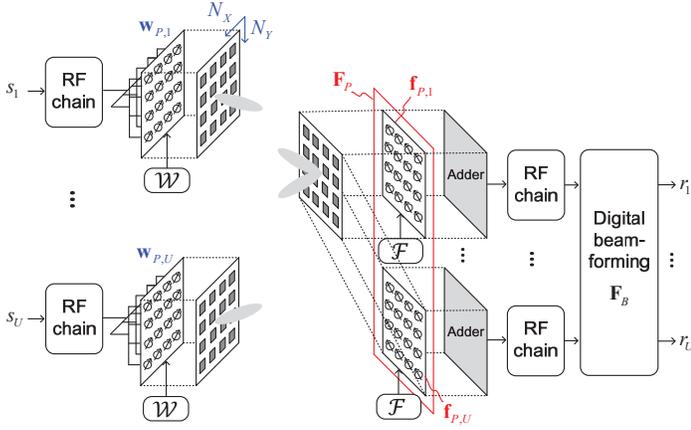


Fig. 2. A multi-UAV hybrid beamforming system has a lead with a hybrid analog-digital beamformer and U followers equipped with analog beamformers.

- a A scalar.
- \mathbf{a} A column vector.
- \mathbf{A} A matrix.
- \mathcal{A} A set.
- $[\mathbf{a}]_n$ The n^{th} entry of \mathbf{a} .
- \mathbf{A}^* The complex conjugate of \mathbf{A} .
- \mathbf{A}^H The Hermitian transpose of \mathbf{A} .
- \mathbf{I}_N The $N \times N$ identity matrix.

II. SYSTEM MODEL

A clustered multi-UAV beamforming system shown in Fig. 2 has one lead and U followers. We assume that these UAVs are perfectly synchronized in time and frequency, and the lead communicates U data streams to U followers at the same time and frequency. That is, we consider space-division multiple access (SDMA) with beamforming to enable data transmission/reception for multiple UAVs [32], [33], and let each UAV be equipped with a uniform rectangular array (URA) of $N = N_X N_Y$ antennas.

The goal of multi-UAV beamforming in a highly dynamic environment is to maximize the system throughput in a discrete time interval $t = 0, \dots, T$. At the cluster lead, the signals are received from specific directions using U analog beamformers at time t , denoted by $\mathbf{f}_{P,u,t} \in \mathbb{C}^{N \times 1}$, $u = 1, \dots, U$. The analog beamformers are implemented in the *passband* as part of the RF front end. Due to the concerns of high implementation costs and power consumption, they have some limitations, e.g., the weights of analog beamformers have unit magnitude because analog beamformers are typically implemented by phase shifters [12]. The U analog beamforming vectors together are denoted by the matrix $\mathbf{F}_{P,t} = [\mathbf{f}_{P,1,t}, \dots, \mathbf{f}_{P,U,t}] \in \mathbb{C}^{N \times U}$, and these vectors can be further combined with the weights of the *baseband* digital beamformer $\mathbf{F}_{B,t} \in \mathbb{C}^{U \times U}$.

Given a pre-defined codebook $\mathcal{F} = \{\tilde{\mathbf{f}}_{n_f} \in \mathbb{C}^{N \times 1}, n_f = 1, \dots, N_F, N_F > U\}$, the U analog beamforming vectors at the lead are selected from the set \mathcal{F} . Beam $\tilde{\mathbf{f}}_{n_f}$ of the URA, i.e., the n_f^{th} member of \mathcal{F} can be represented by the Kronecker product (denoted by \otimes) of the beamforming vectors $\tilde{\mathbf{f}}_{X,n_f} \in$

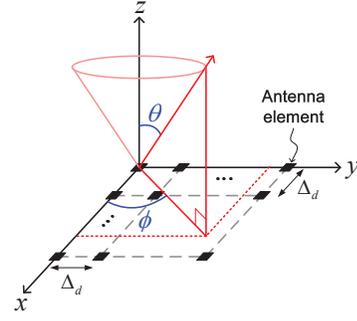


Fig. 3. An array geometry of the URA.

$\mathbb{C}^{N_X \times 1}$ and $\tilde{\mathbf{f}}_{Y,n_f} \in \mathbb{C}^{N_Y \times 1}$ in x - and y -direction respectively [34]:

$$\tilde{\mathbf{f}}_{n_f} = \tilde{\mathbf{f}}_{X,n_f} \otimes \tilde{\mathbf{f}}_{Y,n_f}, \quad (1)$$

and the element of $\tilde{\mathbf{f}}_{X,n_f}$ and $\tilde{\mathbf{f}}_{Y,n_f}$ can be represented by

$$\begin{aligned} [\tilde{\mathbf{f}}_{X,n_f}]_{n_x} &= \frac{\exp\left(-j\frac{2\pi}{\lambda_0} \cos(\phi_{n_f}) \sin(\theta_{n_f})(n_x - 1)\Delta_d\right)}{\sqrt{N_X}}, \\ [\tilde{\mathbf{f}}_{Y,n_f}]_{n_y} &= \frac{\exp\left(-j\frac{2\pi}{\lambda_0} \sin(\phi_{n_f}) \sin(\theta_{n_f})(n_y - 1)\Delta_d\right)}{\sqrt{N_Y}}, \end{aligned} \quad (2)$$

where $n_x = 1, \dots, N_X$ and $n_y = 1, \dots, N_Y$ are the indices of antenna elements in x - and y -direction respectively. Also, ϕ_{n_f} and θ_{n_f} are respectively the n_f^{th} candidate for the azimuth and elevation steering angles at the lead (see Fig. 3), $\Delta_d = \lambda_0/2$ is the distance between neighboring antenna elements, and λ_0 is the wavelength at the carrier frequency.

For the U followers, each only uses a single analog beamformer $\mathbf{w}_{P,u,t} \in \mathbb{C}^{N \times 1}$ with N phase shifters to communicate with the lead.³ Similar to the analog beams at the lead, each follower selects an analog beam from codebook $\mathcal{W} = \{\tilde{\mathbf{w}}_{n_w} \in \mathbb{C}^{N \times 1}, n_w = 1, \dots, N_W, N_W > U\}$.⁴

Via a time-varying channel $\mathbf{H}_{u,t} \in \mathbb{C}^{N \times N}$ between the lead and follower u , the received signal at the lead after the hybrid beamformer is the superposition of the desired signal, interference from other UAVs, and combined noise [30], [31]:

$$\begin{aligned} r_{u,t} &= \underbrace{\mathbf{f}_{B,u,t}^H \mathbf{F}_{P,t}^H \mathbf{H}_{u,t} \mathbf{w}_{P,u,t} s_{u,t}}_{\text{desired signal}} \\ &+ \underbrace{\mathbf{f}_{B,u,t}^H \mathbf{F}_{P,t}^H \sum_{i=1, i \neq u}^U \mathbf{H}_{i,t} \mathbf{w}_{P,i,t} s_{i,t}}_{\text{interference}} + \underbrace{\mathbf{f}_{B,u,t}^H \mathbf{F}_{P,t}^H \mathbf{n}_t}_{\text{combined noise}}, \end{aligned} \quad (3)$$

where $s_{u,t} \in \mathbb{C}$ is the pilot signal satisfying $|s_{u,t}|^2 = 1$ and $\mathbb{E}[s_{u,t} s_{i,t}^*] = 0$, $\mathbf{n}_t \in \mathbb{C}^{N \times 1}$ is an N -dimensional circularly symmetric complex Gaussian (CSCG) random noise vector with mean $\mathbf{0}_{N \times 1}$ and covariance matrix $\sigma_n^2 \mathbf{I}_N$, i.e.,

³We assume that all the UAVs are equipped with a hybrid beamforming architecture since the leading UAV may change over time. The lead is randomly selected from $U + 1$ UAVs at the beginning.

⁴Essentially, these two codebooks are the same, i.e., $\mathcal{W} = \mathcal{F}$. We specify the beamforming problem in terms of two different notations of codebooks for generality.

$\mathbf{n}_t \sim \mathcal{CN}(\mathbf{0}_{N \times 1}, \sigma_n^2 \mathbf{I}_N)$, and $\mathbf{f}_{B,u,t} \in \mathbb{C}^{U \times 1}$ is the u^{th} column of $\mathbf{F}_{B,t}$.

The link between the lead and follower u is modeled as a line-of-sight (LoS) path. According to the relative position and orientation between the transmitter and receiver, the MIMO channel matrix can be determined by the complex path gain $\rho_u \in \mathbb{C}$ and the outer product of two array response vectors $\mathbf{a}_{A,u,t} \in \mathbb{C}^{N \times 1}$ and $\mathbf{a}_{D,u,t} \in \mathbb{C}^{N \times 1}$, which are functions of AoA and AoD [15], [35]. Thus, the channel matrix is expressed by

$$\mathbf{H}_{u,t} = \rho_u \cdot \mathbf{a}_{A,u,t} \cdot \mathbf{a}_{D,u,t}^H. \quad (4)$$

In a manner similar to the steering vector in (1), the array response vectors can be represented by the Kronecker product of the array response vectors in x - and y -direction. Take $\mathbf{a}_{D,u,t}$ as an example:

$$\mathbf{a}_{D,u,t} = \mathbf{a}_{D,X,u,t} \otimes \mathbf{a}_{D,Y,u,t}, \quad (5)$$

and the entries of $\mathbf{a}_{D,X,u,t}$ and $\mathbf{a}_{D,Y,u,t}$ are given by

$$\begin{aligned} [\mathbf{a}_{D,X,u,t}]_{n_x} &= \frac{\exp\left(\frac{-j2\pi}{\lambda_0} \cos(\phi_{D,u,t}) \sin(\theta_{D,u,t})(n_x - 1)\Delta_d\right)}{\sqrt{N_X}}, \\ [\mathbf{a}_{D,Y,u,t}]_{n_y} &= \frac{\exp\left(\frac{-j2\pi}{\lambda_0} \sin(\phi_{D,u,t}) \sin(\theta_{D,u,t})(n_y - 1)\Delta_d\right)}{\sqrt{N_Y}}, \end{aligned} \quad (6)$$

where the random variables $\phi_{D,u,t}$ and $\theta_{D,u,t}$ stand for the azimuth and elevation angles of departure at time t . Given the azimuth and elevation angles of arrival (denoted by $\phi_{A,u,t}$, $\theta_{A,u,t}$), the array response vector at the receiver (i.e., $\mathbf{a}_{A,u,t}$) has a similar form as (5).

To model a highly dynamic environment for the angles under an observed LoS path, a Gaussian random walk is used to generate the time-varying angles $\phi_{A,u,t}$, $\theta_{A,u,t}$, $\phi_{D,u,t}$, and $\theta_{D,u,t}$. For instance, the azimuth angle of arrival $\phi_{A,u,t}$ can be defined by

$$\phi_{A,u,t} = \phi_{A,u,0} + \sum_{i=1}^t \lambda_i, \quad (7)$$

where $\phi_{A,u,0} \sim \mathcal{U}(0, 2\pi)$ is a randomly selected initial angle of $\phi_{A,u,t}$ and follows a uniform distribution, and $\lambda_i \sim \mathcal{N}(0, \sigma_\lambda^2)$ is the disturbance (or white noise) following a normal distribution. The other three time-varying angles are generated in a similar way.

III. PROBLEM STATEMENT

The goal of hybrid beamforming in the multi-UAV system is to maximize the SINR (or system throughput) during the time interval $[0, T]$. Meanwhile, after the combiner $\mathbf{F}_{P,t} \mathbf{F}_{B,t}$, the variance of the combined noise signal is enforced to remain constant, i.e.,

$$\mathbb{E} \left[(\mathbf{f}_{B,u,t}^H \mathbf{F}_{P,t}^H \mathbf{n}_t) (\mathbf{f}_{B,u,t}^H \mathbf{F}_{P,t}^H \mathbf{n}_t)^H \right] = \sigma_n^2 \quad \forall u, t, \quad (8)$$

which leads to a power constraint on the combiner as

$$\mathbf{f}_{B,u,t}^H \mathbf{F}_{P,t}^H \mathbf{F}_{P,t} \mathbf{f}_{B,u,t} = 1 \quad \forall u, t. \quad (9)$$

Then, by introducing two sets $\mathcal{I}_{\mathcal{F},t}$ and $\mathcal{I}_{\mathcal{W},t}$ that include promising candidates for the analog beamforming matrices, we seek $\mathbf{F}_{P,t}$, $\mathbf{F}_{B,t}$, and $\mathbf{W}_{P,t}$ that together achieve the maximum SINR and satisfy the power constraint from $t = 0$ to $t = T$:

$$\begin{aligned} \max_{\mathbf{F}_{P,t} \in \mathcal{I}_{\mathcal{F},t}, \mathbf{W}_{P,t} \in \mathcal{I}_{\mathcal{W},t}} \left\{ \max_{\mathbf{F}_{B,t}} \sum_{u=1}^U \frac{P_{S,u,t}}{P_{I,u,t} + \sigma_n^2} \right\} \quad (10) \\ \text{s.t. } \mathbf{f}_{B,u,t}^H \mathbf{F}_{P,t}^H \mathbf{F}_{P,t} \mathbf{f}_{B,u,t} = 1 \quad \forall u, t, \end{aligned}$$

where $P_{S,u,t}$ and $P_{I,u,t}$ are the power of the desired and interference signals given by

$$P_{S,u,t} = \left| \mathbf{f}_{B,u,t}^H \mathbf{F}_{P,t}^H \mathbf{H}_{u,t} \mathbf{w}_{P,u,t} \right|^2, \quad (11)$$

$$P_{I,u,t} = \sum_{i=1, i \neq u}^U \left| \mathbf{f}_{B,u,t}^H \mathbf{F}_{P,t}^H \mathbf{H}_{i,t} \mathbf{w}_{P,i,t} \right|^2. \quad (12)$$

In the paper, we do not assume the channel state information or any knowledge of AoAs/AoDs is known to the lead. Instead, the required observations are the estimates of *coupling coefficients* associated with a beam pair $(\tilde{\mathbf{f}}_{n_f}, \tilde{\mathbf{w}}_{n_w})$, where $\tilde{\mathbf{f}}_{n_f} \in \mathcal{F}$ and $\tilde{\mathbf{w}}_{n_w} \in \mathcal{W}$. By correlating the received pilot signals with the known transmitted ones, we can obtain such observations given by⁵

$$\begin{aligned} y_{u,t}(n_f, n_w) &= s_{u,t}^* \underbrace{\left(\tilde{\mathbf{f}}_{n_f}^H \mathbf{H}_{u,t} \tilde{\mathbf{w}}_{n_w} s_{u,t} + \tilde{\mathbf{f}}_{n_f}^H \sum_{i=1, i \neq u}^U \mathbf{H}_{i,t} \tilde{\mathbf{w}}_{n_w} s_{i,t} + \tilde{\mathbf{f}}_{n_f}^H \mathbf{n}_t \right)}_{\text{received pilot signal}} \\ &= \tilde{\mathbf{f}}_{n_f}^H \mathbf{H}_{u,t} \tilde{\mathbf{w}}_{n_w} + \underbrace{\left(\tilde{\mathbf{f}}_{n_f}^H \sum_{i=1, i \neq u}^U \mathbf{H}_{i,t} \tilde{\mathbf{w}}_{n_w} s_{i,t}^* + s_{u,t}^* \tilde{\mathbf{f}}_{n_f}^H \mathbf{n}_t \right)}_{\triangleq z_t} \\ &= \underbrace{\tilde{\mathbf{f}}_{n_f}^H \mathbf{H}_{u,t} \tilde{\mathbf{w}}_{n_w}}_{\text{coupling coefficient}} + z_t, \end{aligned} \quad (13)$$

where z_t denotes the superposition of the combined interference and noise, and we assume that it follows a complex normal distribution, i.e., $z_t \sim \mathcal{CN}(0, \sigma_z^2)$.

Given the observations $\{y_{u,t}(n_f, n_w) \forall u, t\}$, the strategy of solving the problem (10) could be, first, using the observations to find the sets $\mathcal{I}_{\mathcal{F},t}$ and $\mathcal{I}_{\mathcal{W},t}$ that ideally consist of the optimal analog beamforming matrices. However, due to the hardware constraint on the analog beamformer, the beam probing is time-consuming. When the channel is highly dynamic, the observations acquired early may become unreliable. How to use the observations to interact with the highly dynamic environment during the beam probing becomes a crucial problem. As a result, the idea of *Q-learning* algorithm [28] is borrowed to find appropriate beams (i.e., the members of $\mathcal{I}_{\mathcal{F},t}$ and $\mathcal{I}_{\mathcal{W},t}$) for time-varying channels. The concept of Q-learning is to let the UAVs learn the optimal behavior directly from the interaction with the environment. Once we determine the candidate sets $\mathcal{I}_{\mathcal{F},t}$ and $\mathcal{I}_{\mathcal{W},t}$, the observations associated

⁵The notation of observation $y_{u,t}(n_f, n_w)$ is simplified from its formal expression given by $y_{u,t}(n_f = n_f(u, t) \in \{1, \dots, N_F\}, n_w = n_w(u, t) \in \{1, \dots, N_W\})$.

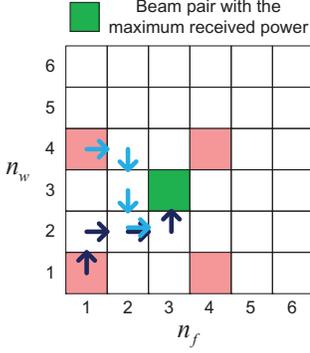


Fig. 4. All the candidates for the beam pairs $\{(\tilde{\mathbf{f}}_{n_f}, \tilde{\mathbf{w}}_{n_w}) \forall n_f, n_w\}$ are represented by the grid map, where the red ones are trained during the initial beam search. An example of the Q-learning-based beam selection is given in Example 1. According to the updated Q-values, see Table I, it will converge to beam pair $(\tilde{\mathbf{f}}_3, \tilde{\mathbf{w}}_3)$ after few iterations.

with the members of $\mathcal{I}_{\mathcal{F},t}$ and $\mathcal{I}_{\mathcal{W},t}$ are used to generate the corresponding digital weights and the SINR measurement.

IV. ANALOG BEAM TRACKING USING Q-LEARNING

In this section, we introduce an analog beam tracking algorithm for highly dynamic environments. Starting from a single link between the lead and a follower, we adopt Q-learning to deal with the beam tracking problem. The idea can be easily extended to multiple links with additional constraints.

A. Beam Selection Using Q-Learning for One Link

To begin with, let us focus on the link between the lead and follower u . That is, we seek the candidates for $\mathbf{f}_{P,u,t}$ and $\mathbf{w}_{P,u,t}$. When the codebook size is large, the efficient way of beam tracking is to start from some specific directions that cover the 3-D environment. This phase is called *initial beam search*. For example, Fig. 4 shows $N_F N_W = 6 \times 6$ candidates for the analog beam pair, where N_F and N_W are the numbers of elements in codebooks \mathcal{F} and \mathcal{W} respectively. In the example, the four beam pairs highlighted in red are initially explored. To be formal, we define two sets that consist of the beams used in the initial search by $\mathcal{F}_{\mathcal{I}} = \{\tilde{\mathbf{f}}_1, \tilde{\mathbf{f}}_4\}$ and $\mathcal{W}_{\mathcal{I}} = \{\tilde{\mathbf{w}}_1, \tilde{\mathbf{w}}_4\}$ and assume that both the lead and follower have the same initial beam search pattern. After the beam probing using these four beam pairs, the one having the maximum received power will be selected as a starting point of beam tracking in the next phase.

The beam tracking is conventionally implemented by searching a better choice next to the currently used beam pair [20], [36]. Both the initial beam search and beam tracking in the above-mentioned work only explore the environment rather than interact with the environment. The concept of “interaction with the environment” can be viewed as a beam selection algorithm that can *explore* uncharted territory and, meanwhile, *exploit* the searching experience. Concerning a highly dynamic environment, the exploration-exploitation balance becomes more important to the beam tracking. The idea of Q-learning is to let an agent (e.g., a UAV) learn to strike the balance between exploration and exploitation.

TABLE I
THE Q-VALUES ARE UPDATED ACCORDING TO THE STATES AND ACTIONS GIVEN IN EXAMPLE 1 AND FIG. 4. HERE WE LET THE Q-VALUES BE UPDATED BY EITHER 0 OR 1 FOR SIMPLICITY.

Time (t)	Episode	Step ($N_S = 4$)	State S_t	Action A_t			
				\uparrow	\downarrow	\rightarrow	\leftarrow
0	0	0	$(\tilde{\mathbf{f}}_1, \tilde{\mathbf{w}}_1)$	1	0	0	0
1		1	$(\tilde{\mathbf{f}}_1, \tilde{\mathbf{w}}_2)$	0	0	1	0
2		2	$(\tilde{\mathbf{f}}_2, \tilde{\mathbf{w}}_2)$	0	0	1	0
3		3	$(\tilde{\mathbf{f}}_3, \tilde{\mathbf{w}}_2)$	1	0	0	0
4	1	0	$(\tilde{\mathbf{f}}_1, \tilde{\mathbf{w}}_4)$	0	0	1	0
5		1	$(\tilde{\mathbf{f}}_2, \tilde{\mathbf{w}}_4)$	0	1	0	0
6		2	$(\tilde{\mathbf{f}}_2, \tilde{\mathbf{w}}_3)$	0	1	0	0
7		3	$(\tilde{\mathbf{f}}_2, \tilde{\mathbf{w}}_2)$	0	0	2	0

⋮

In Q-learning, the experience is recorded in a Q-learning table (or Q-table), see Table I, which is updated according to the current measurements. The Q-table is constructed according to three components: states, actions, and state-action values (also known as Q-values). Before the learning begins, the state-action values in the Q-table are initialized to zero. In a state S_t at time t , the UAV always implements the following four steps: select an action A_t from the action set $\mathcal{A} = \{\text{up, down, right, left}\}$, go to the next state S_{t+1} , observe a reward R_{t+1} , and update the Q-value, given by [28, Ch. 6]

$$Q(S_t, A_t) \leftarrow \underbrace{(1-\alpha)Q(S_t, A_t)}_{\text{old value}} + \alpha \underbrace{\left[R_{t+1} + \gamma \max_{a \in \mathcal{A}} Q(S_{t+1}, a) \right]}_{\text{new information}}, \quad (14)$$

where $0 < \alpha < 1$ is the learning rate (or step size), $0 < \gamma < 1$ is the discount factor determining the importance of future rewards. The Q-value update can be described as a weighted average between the old value and new information.

The reward can be regarded as the feedback from the environment given an action. In terms of maximizing the SINR, the reward is supposed to be a function of SINR. Nevertheless, we only have the coupling coefficients as measurements which suffer from noise and interference. We therefore define the reward function as follows. According to the received power of the coupling coefficients corresponding to the trained beam pairs at time t and $t+1$, the reward is defined, in terms of thresholds, by functions of the received power

$$R_{t+1} = \begin{cases} 1, & \text{if } \frac{|y_{u,t+1}(n'_f, n'_w)|^2}{|y_{u,t}(n_f, n_w)|^2} > c_u \\ 0, & \text{if } c_l < \frac{|y_{u,t+1}(n'_f, n'_w)|^2}{|y_{u,t}(n_f, n_w)|^2} \leq c_u \\ -1, & \text{otherwise} \end{cases} \quad (15)$$

where (n'_f, n'_w) is the beam index pair used at time $t+1$. Due to the noise and interference, the observations, $y_{u,t+1}(n'_f, n'_w)$ and $y_{u,t}(n_f, n_w)$, may be unreliable for determining the reward. To reduce the uncertainty, we define a lower threshold c_l and an upper threshold c_u . If the ratio of $|y_{u,t+1}(n'_f, n'_w)|^2$ to $|y_{u,t}(n_f, n_w)|^2$ is between c_l and c_u , the measurement is treated as ambiguity so that the reward is equal to zero. A more detailed discussion about the upper and lower thresholds is provided in Appendix A.

To elaborate the Q-learning-based beam selection, let us take an example by Fig. 4 and Table I.

Example 1. When starting from a state $S_0 = (\tilde{\mathbf{f}}_1, \tilde{\mathbf{w}}_1)$, one of the neighboring beam pairs $\{(\tilde{\mathbf{f}}_1, \tilde{\mathbf{w}}_2), (\tilde{\mathbf{f}}_1, \tilde{\mathbf{w}}_6), (\tilde{\mathbf{f}}_2, \tilde{\mathbf{w}}_1), (\tilde{\mathbf{f}}_6, \tilde{\mathbf{w}}_1)\}$ will be explored by choosing an action from \mathcal{A} according to the state-action values, i.e., $\max_{a \in \mathcal{A}} Q(S_0, a)$. Since all the Q-values at S_0 are initialized to zero, an action will be selected randomly (or according to some predefined criteria). We assume that the action “up” is selected so that the next state becomes $S_1 = (\tilde{\mathbf{f}}_1, \tilde{\mathbf{w}}_2)$. The corresponding reward and Q-value $Q(S_0, A_0 = \text{up})$ will be updated accordingly, see Table I. In the example, we simply let the Q-values be updated by either 0 or 1, where a value of 1 implies that the agent chooses the action and gets a positive reward. In Q-learning, a sequence of $N_S = 4$ time slots (also called steps) is defined as an *episode*. Each episode starts from a state, which could be pre-defined or determined by the received power. Fig. 4 shows that the initial beam search needs in total four episodes with starting states at $(\tilde{\mathbf{f}}_1, \tilde{\mathbf{w}}_1)$, $(\tilde{\mathbf{f}}_1, \tilde{\mathbf{w}}_4)$, $(\tilde{\mathbf{f}}_4, \tilde{\mathbf{w}}_1)$, and $(\tilde{\mathbf{f}}_4, \tilde{\mathbf{w}}_4)$ respectively. In each episode, the beam probing takes N_S time slots to update the Q-values. When finishing the first episode, the agent starts the next episode using beam pair $(\tilde{\mathbf{f}}_1, \tilde{\mathbf{w}}_4)$. With a sufficiently large number of significant Q-values, Q-learning will converge to the beam pair $(\tilde{\mathbf{f}}_3, \tilde{\mathbf{w}}_3)$ corresponding to the maximum received power. \square

After the initial beam search, some beam pairs have been explored and the beam tracking will start from the beam pair with the maximum received power during the initial beam search, which is denoted by S_{MP} (i.e., the state or beam pair with respect to the maximum power).

According to the updated Q-values, an agent exploits what it has already experienced in order to obtain a positive reward, but it also has to explore the uncharted or changed environment to see if it can make better action selections in the future. One of the challenges in reinforcement learning is the trade-off between the exploration and exploitation. By introducing a parameter $0 < \varepsilon < 1$, an ε -greedy action is obtained to better balance the exploration and exploitation:

$$A_t = \begin{cases} \operatorname{argmax}_{a \in \mathcal{A}} Q(S_t, a), & \text{with prob. } 1 - \varepsilon \\ \text{a random action,} & \text{with prob. } \varepsilon \end{cases} \quad (16)$$

The agent chooses the action as it believes that the action yields the best long-term effect with probability $1 - \varepsilon$. Or the agent chooses an action uniformly at random with probability ε .

The pseudocode of the Q-learning-based beam tracking algorithm is shown in **Algorithm 1**, which includes two phases: the initial beam search and beam tracking. The difference between these two phases is the decision of the starting state of each episode. During the initial beam search, the starting state is selected from the pre-defined sets $\mathcal{F}_{\mathcal{I}}$ and $\mathcal{W}_{\mathcal{I}}$. In Example 1, $\mathcal{F}_{\mathcal{I}} = \{\tilde{\mathbf{f}}_1, \tilde{\mathbf{f}}_4\}$ and $\mathcal{W}_{\mathcal{I}} = \{\tilde{\mathbf{w}}_1, \tilde{\mathbf{w}}_4\}$. During the beam tracking, the starting state is selected according to the maximum received power. Moreover, the selected beam pair at time t is denoted by $(\hat{\mathbf{f}}_{P,u,t}, \hat{\mathbf{w}}_{P,u,t})$. We assume that the

Algorithm 1 Q-learning beam tracking for a single link.

Input: Observations $\{y_{u,t}(n_f, n_w), t = 0, \dots, T\}$

Output: Selected beam pairs $\{(\hat{\mathbf{f}}_{P,u,t}, \hat{\mathbf{w}}_{P,u,t}), t = 0, \dots, T\}$

```

1: Initialize Q-table
2:  $t = 0$ 
3: for  $i = 1$  : number of episodes
4:   if initial beam search
5:      $S_t \in \{(\tilde{\mathbf{f}}_{n_f}, \tilde{\mathbf{w}}_{n_w}) \mid \tilde{\mathbf{f}}_{n_f} \in \mathcal{F}_{\mathcal{I}}, \tilde{\mathbf{w}}_{n_w} \in \mathcal{W}_{\mathcal{I}}\}$ 
6:   else if beam tracking
7:      $S_t = S_{MP}$ 
8:   end if
9:   for  $j = 1$  : number  $N_S$  of steps
10:    choose  $A_t$  and go to  $S_{t+1} \equiv (\hat{\mathbf{f}}_{P,u,t+1}, \hat{\mathbf{w}}_{P,u,t+1})$ 
11:    obtain  $R_{t+1}$  according to observations
12:    update  $Q(S_t, A_t)$ 
13:    update  $S_{MP}$  according to observations
14:     $t = t + 1$ 
15:   end step
16: end episode

```

analog beam pairs are determined at the UAV lead, and time division duplex (TDD) technique that separates the transmit and receive signals in the time domain can be used to inform the followers to update their beams.

B. Overhead Reduction Using Offline Q-Learning

In Algorithm 1, the observations are available at each time slot t . This implies that the beam switching and pilot transmission/reception are executed in every time slot, which is not a well-designed manner in the sense of system efficiency. To reduce the overhead, we reserve all observations so that the Q-learning can execute offline. When using past observations to obtain the rewards and update the Q-values, we name the Q-learning algorithm *offline* Q-learning. Otherwise, it is called *online* Q-learning.

For the offline Q-learning, only the observations associated with large received power have to be updated regularly. Therefore, at the end of each episode, the beam pairs with respect to the maximum received power (i.e., S_{MP}) will be chosen and employed at the beginning of each episode in order to update the corresponding observations. For other steps in an episode, the pilot transmission and beam switching are not necessary unless a specific state has not been explored.

C. Beam Selection Using Q-Learning for Multiple Links

The idea of Q-learning-based beam tracking for one link can be easily extended to the case of multiple links, similar to multi-agent systems [37], [38]. For multi-UAV beam probing, the lead receives the observations from different followers simultaneously in an SDMA manner. In this case, the members of \mathcal{F} at the lead UAV’s side should not be selected repeatedly. As a result, the action set in (16) has to be updated in real time.

In each beam probing, which could be in the stage of initial beam search or beam tracking, the Q-learning-based beam selection starts from a follower corresponding to the maximum

received power at the moment. We further define a set \mathcal{A}' that includes the actions which will make different followers go to the same states. Thus, the action selection given in (16) can be reformulated as

$$A_t = \begin{cases} \operatorname{argmax}_{a \in \mathcal{A} \setminus \mathcal{A}'} Q(S_t, a), & \text{with prob. } 1 - \varepsilon \\ \text{randomly selected from } \mathcal{A} \setminus \mathcal{A}', & \text{with prob. } \varepsilon \end{cases} \quad (17)$$

After making the decision about the next state for a follower, the lead has to update \mathcal{A}' accordingly.

V. DIGITAL BEAMFORMING

In the previous section, we use Q-learning to find the members of sets $\mathcal{I}_{\mathcal{F},t}$ and $\mathcal{I}_{\mathcal{W},t}$ in the problem (10). However, the selected beam pairs may not be the optimal solution to the problem for the reasons that (i) Q-learning usually only provides a good enough solution⁶ and (ii) the digital beamformer weights are not taken into account during the procedure of analog beam selection. In the sense of hybrid beamforming, a better solution should be the one whose linear combination with the digital weights leading to the maximum SINR. This issue can be solved by keeping more than one promising members with large received power in $\mathcal{I}_{\mathcal{F},t}$ and $\mathcal{I}_{\mathcal{W},t}$ [17]. We use Example 2 to explain the idea.

Example 2. Two selected beam pairs with large received power for each follower are collected in the following two sets:

$$\{[\tilde{\mathbf{f}}_{P,1}, \tilde{\mathbf{f}}_{P,2}, \tilde{\mathbf{f}}_{P,3}], [\tilde{\mathbf{f}}_{P,1}, \tilde{\mathbf{f}}_{P,3}, \tilde{\mathbf{f}}_{P,4}]\}$$

and

$$\{[\tilde{\mathbf{w}}_{P,1}, \tilde{\mathbf{w}}_{P,1}, \tilde{\mathbf{w}}_{P,2}], [\tilde{\mathbf{w}}_{P,2}, \tilde{\mathbf{w}}_{P,3}, \tilde{\mathbf{w}}_{P,4}]\}.$$

Given these two sets, we can generate all the members of $\mathcal{I}_{\mathcal{F},t}$ and $\mathcal{I}_{\mathcal{W},t}$, given by

$$\mathcal{I}_{\mathcal{F},t} = \underbrace{\{[\tilde{\mathbf{f}}_{P,1}, \tilde{\mathbf{f}}_{P,2}, \tilde{\mathbf{f}}_{P,3}], [\tilde{\mathbf{f}}_{P,1}, \tilde{\mathbf{f}}_{P,2}, \tilde{\mathbf{f}}_{P,4}], [\tilde{\mathbf{f}}_{P,1}, \tilde{\mathbf{f}}_{P,3}, \tilde{\mathbf{f}}_{P,4}]\}}_{\substack{\text{the 1st candidate} \\ \text{for } \mathbf{F}_{P,t}}}$$

which has a cardinality of 3 because the members of \mathcal{F} at lead UAV should not be selected repeatedly, and the other set can be represented by

$$\mathcal{I}_{\mathcal{W},t} = \underbrace{\{[\tilde{\mathbf{w}}_{P,1}, \tilde{\mathbf{w}}_{P,1}, \tilde{\mathbf{w}}_{P,2}], [\tilde{\mathbf{w}}_{P,1}, \tilde{\mathbf{w}}_{P,1}, \tilde{\mathbf{w}}_{P,4}], \dots, [\tilde{\mathbf{w}}_{P,2}, \tilde{\mathbf{w}}_{P,3}, \tilde{\mathbf{w}}_{P,4}]\}}_{\substack{\text{the 1st candidate} \\ \text{for } \mathbf{W}_{P,t}}}$$

which has a cardinality of 8. In this example, given the above $\mathcal{I}_{\mathcal{F},t}$ and $\mathcal{I}_{\mathcal{W},t}$, we have to evaluate a total of 24 combinations with their digital weights to maximize the SINR. \square

The above-mentioned idea is different from the work represented in [25] that keeps candidates in subspace. In our opinion, the better solution is supposed to keep candidates with large received power because the idea in [25] only takes into account the main lobes of analog beams, while the proposed method considers both the main and side lobes.

⁶Q-learning uses experience to solve a prediction problem, which can be viewed as a Monte Carlo method.

A. Digital Weight Optimization

To simplify the following descriptions of digital beamforming, we assume that $\mathcal{I}_{\mathcal{F},t}$ and $\mathcal{I}_{\mathcal{W},t}$ only include one member respectively, i.e.,

$$\begin{aligned} \mathcal{I}_{\mathcal{F},t} &= \{\hat{\mathbf{F}}_{P,t} = [\hat{\mathbf{f}}_{P,1,t}, \dots, \hat{\mathbf{f}}_{P,U,t}]\} \\ \mathcal{I}_{\mathcal{W},t} &= \{\hat{\mathbf{W}}_{P,t} = [\hat{\mathbf{w}}_{P,1,t}, \dots, \hat{\mathbf{w}}_{P,U,t}]\}. \end{aligned} \quad (18)$$

In the numerical results, we will provide more discussion about the idea. Given $\hat{\mathbf{F}}_{P,t}$ and $\hat{\mathbf{W}}_{P,t}$, the hybrid beamforming problem (10) becomes a digital beamforming problem subject to the power constraint, which can be formulate as

$$\begin{aligned} &\sum_{u=1}^U \max_{\mathbf{f}_{B,u,t}} \frac{\hat{P}_{S,u,t}}{\hat{P}_{I,u,t} + \sigma_n^2} \\ \text{s.t. } &\mathbf{f}_{B,u,t}^H \hat{\mathbf{F}}_{P,t}^H \hat{\mathbf{F}}_{P,t} \mathbf{f}_{B,u,t} = 1 \quad \forall u \end{aligned} \quad (19)$$

where $t = 1, \dots, T$. The signal and interference power are subject to the selected analog beams

$$\hat{P}_{S,u,t} \triangleq P_{S,u,t} \Big|_{\mathbf{F}_{P,t}=\hat{\mathbf{F}}_{P,t}, \mathbf{W}_{P,t}=\hat{\mathbf{W}}_{P,t}}, \quad (20)$$

$$\hat{P}_{I,u,t} \triangleq P_{I,u,t} \Big|_{\mathbf{F}_{P,t}=\hat{\mathbf{F}}_{P,t}, \mathbf{W}_{P,t}=\hat{\mathbf{W}}_{P,t}}. \quad (21)$$

To satisfy the power constraint on the combiner, one can define U unit vectors $\{\mathbf{x}_u \mid \|\mathbf{x}_u\|_2 = 1, u = 1, \dots, U\}$ that obey the relation [17]

$$\mathbf{f}_{B,u,t} = (\hat{\mathbf{F}}_{P,t}^H \hat{\mathbf{F}}_{P,t})^{-0.5} \mathbf{x}_u. \quad (22)$$

Upon replacing $\mathbf{f}_{B,u,t}$ with $(\hat{\mathbf{F}}_{P,t}^H \hat{\mathbf{F}}_{P,t})^{-0.5} \mathbf{x}_u$ in the problem, the received signal and interference power can be written by

$$\hat{P}_{S,u,t} = \left| \mathbf{x}_u^H (\hat{\mathbf{F}}_{P,t}^H \hat{\mathbf{F}}_{P,t})^{-0.5} \hat{\mathbf{F}}_{P,t}^H \mathbf{H}_{u,t} \hat{\mathbf{w}}_{P,u,t} \right|^2, \quad (23)$$

$$\hat{P}_{I,u,t} = \sum_{i=1, i \neq u}^U \left| \mathbf{x}_u^H (\hat{\mathbf{F}}_{P,t}^H \hat{\mathbf{F}}_{P,t})^{-0.5} \hat{\mathbf{F}}_{P,t}^H \mathbf{H}_{i,t} \hat{\mathbf{w}}_{P,i,t} \right|^2. \quad (24)$$

Then, we can find that the problem (19) is equivalent to seeking vectors $\mathbf{x}_1, \dots, \mathbf{x}_U$ that maximize the SINR for U followers. As a result, the maximization problem (19) can be reformulated as

$$\sum_{u=1}^U \max_{\mathbf{x}_u} \frac{\hat{P}_{S,u,t}}{\hat{P}_{I,u,t} + \sigma_n^2}. \quad (25)$$

B. SINR Approximation Using Coupling Coefficients

In (23) and (24), the couplings of the channel and analog beams, such as $\hat{\mathbf{F}}_{P,t}^H \mathbf{H}_{u,t} \hat{\mathbf{w}}_{P,u,t}$ and $\hat{\mathbf{F}}_{P,t}^H \mathbf{H}_{i,t} \hat{\mathbf{w}}_{P,i,t}$, can be viewed as effective channel vectors. Since the observations, given in (13), are the coupling of the channel and one analog beam pair, we can use them to construct the estimates of effective channel vectors, defined by

$$\begin{aligned} \hat{\mathbf{h}}_{E,u,t} &= \hat{\mathbf{F}}_{P,t}^H \mathbf{H}_{u,t} \hat{\mathbf{w}}_{P,u,t} + z_t \\ &= \underbrace{\begin{bmatrix} \hat{\mathbf{f}}_{P,1,t}^H \mathbf{H}_{u,t} \hat{\mathbf{w}}_{P,u,t} + z_t \\ \vdots \\ \hat{\mathbf{f}}_{P,U,t}^H \mathbf{H}_{u,t} \hat{\mathbf{w}}_{P,u,t} + z_t \end{bmatrix}}_{\hat{\mathbf{h}}_{E,u,t}}, \end{aligned} \quad (26)$$

The entries of $\hat{\mathbf{h}}_{E,u,t}$ can be obtained from $\{y_{u,t}(n_f, n_u)\} \forall u$

and

$$\hat{\mathbf{h}}_{E,i,t} = \hat{\mathbf{F}}_{P,t}^H \mathbf{H}_{i,t} \hat{\mathbf{w}}_{P,i,t} + z_t, \quad (27)$$

where the entries of $\hat{\mathbf{h}}_{E,i,t}$ can be obtained from $\{y_{u,t}(n_f, n_w) \forall u\}$ as well. The collected observations suffice to generate the estimates of $\hat{P}_{S,u,t}$ and $\hat{P}_{I,u,t} + \sigma_n^2$ represented by

$$\begin{aligned} \hat{P}_{S,u,t} &\approx \left| \mathbf{x}_u^H (\hat{\mathbf{F}}_{P,t}^H \hat{\mathbf{F}}_{P,t})^{-0.5} \hat{\mathbf{h}}_{E,u,t} \right|^2 \\ &= \mathbf{x}_u^H \underbrace{(\hat{\mathbf{F}}_{P,t}^H \hat{\mathbf{F}}_{P,t})^{-0.5} \hat{\mathbf{h}}_{E,u,t} \hat{\mathbf{h}}_{E,u,t}^H (\hat{\mathbf{F}}_{P,t}^H \hat{\mathbf{F}}_{P,t})^{-0.5}}_{\triangleq \mathbf{A}_{u,t}} \mathbf{x}_u \\ &= \mathbf{x}_u^H \mathbf{A}_{u,t} \mathbf{x}_u \end{aligned} \quad (28)$$

and

$$\begin{aligned} \hat{P}_{I,u,t} + \sigma_n^2 &\approx \sum_{i=1, i \neq u}^U \left| \mathbf{x}_u^H (\hat{\mathbf{F}}_{P,t}^H \hat{\mathbf{F}}_{P,t})^{-0.5} \hat{\mathbf{h}}_{E,i,t} \right|^2 + \sigma_n^2 \\ &= \mathbf{x}_u^H \left(\sum_{i=1, i \neq u}^U (\hat{\mathbf{F}}_{P,t}^H \hat{\mathbf{F}}_{P,t})^{-0.5} \hat{\mathbf{h}}_{E,i,t} \hat{\mathbf{h}}_{E,i,t}^H (\hat{\mathbf{F}}_{P,t}^H \hat{\mathbf{F}}_{P,t})^{-0.5} + \sigma_n^2 \mathbf{I}_U \right) \mathbf{x}_u \\ &= \mathbf{x}_u^H \underbrace{\mathbf{B}_{u,t}}_{\triangleq \mathbf{B}_{u,t}} \mathbf{x}_u \end{aligned} \quad (29)$$

Using (28) and (29), the SINR for follower u conditional on $\mathbf{W}_{P,t} = \hat{\mathbf{W}}_{P,t}$ and $\mathbf{F}_{P,t} = \hat{\mathbf{F}}_{P,t}$ can be approximated by the following equation

$$\frac{\hat{P}_{S,u,t}}{\hat{P}_{I,u,t} + \sigma_n^2} \approx \frac{\mathbf{x}_u^H \mathbf{A}_{u,t} \mathbf{x}_u}{\mathbf{x}_u^H \mathbf{B}_{u,t} \mathbf{x}_u}. \quad (30)$$

Using the property that $\mathbf{B}_{u,t}$ is a positive definite matrix, the optimal solution of \mathbf{x}_u that attains the maximum SINR can be stated as follows (also see Appendix B):

$$\begin{aligned} \mathbf{x}_u^* &= \arg \max_{\mathbf{x}_u} \frac{\mathbf{x}_u^H \mathbf{A}_{u,t} \mathbf{x}_u}{\mathbf{x}_u^H \mathbf{B}_{u,t} \mathbf{x}_u} \\ &= \frac{\mathbf{B}_{u,t}^{-0.5} \mathbf{e}_{\max}(\mathbf{B}_{u,t}^{-0.5} \mathbf{A}_{u,t} \mathbf{B}_{u,t}^{-0.5})}{\left\| \mathbf{B}_{u,t}^{-0.5} \mathbf{e}_{\max}(\mathbf{B}_{u,t}^{-0.5} \mathbf{A}_{u,t} \mathbf{B}_{u,t}^{-0.5}) \right\|_2}, \end{aligned} \quad (31)$$

where $\mathbf{e}_{\max}(\mathbf{B}_{u,t}^{-0.5} \mathbf{A}_{u,t} \mathbf{B}_{u,t}^{-0.5})$ is the eigenvector of $\mathbf{B}_{u,t}^{-0.5} \mathbf{A}_{u,t} \mathbf{B}_{u,t}^{-0.5}$ corresponding to the maximum eigenvalue. In the same manner, given $\mathbf{A}_{u,t}$ and $\mathbf{B}_{u,t}$ for all u , we have the optimal solution of \mathbf{x}_u , $u = 1, \dots, U$. The corresponding estimated digital beamformer weights are therefore given by

$$\begin{aligned} \hat{\mathbf{F}}_{B,t} &= [\hat{\mathbf{f}}_{B,1,t}, \dots, \hat{\mathbf{f}}_{B,U,t}] \\ &= (\hat{\mathbf{F}}_{P,t}^H \hat{\mathbf{F}}_{P,t})^{-0.5} [\mathbf{x}_1^*, \dots, \mathbf{x}_U^*]. \end{aligned} \quad (32)$$

The digital weights represented in (32) are derived from the constraint that the variance of the combined noise signal is still AWGN. When concerning $\mathbf{F}_{B,t}$ acting as part of the precoder for data transmission (i.e., sending signals from the lead to followers), the power constraint could be $\|\mathbf{F}_{P,t} \mathbf{F}_{B,t}\|_F = U$ [30], which leads to zero-forcing (ZF) digital beamforming.

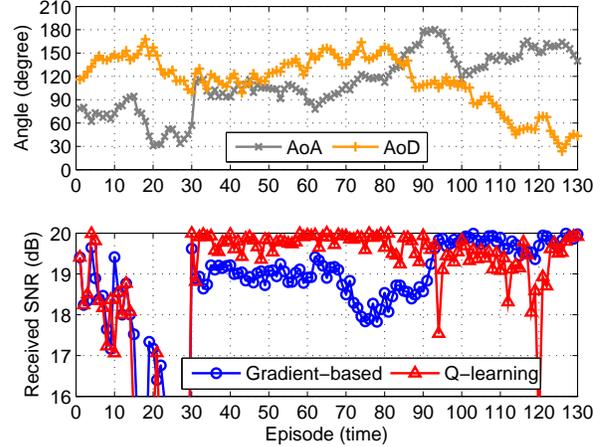


Fig. 5. An typical example of the received SNR using the Q-learning-based and gradient-based beam tracking methods for one link ($U = 1$) with $\sigma_\lambda^2 = 16$. In this example, we assume fixed elevation angles $\theta_{A,u,t} = \theta_{D,u,t} = 15^\circ \forall t$. Compared with the gradient-based method, the Q-learning-based beam tracking is robust to the large variance of angle.

VI. SIMULATION RESULTS

In this section, we numerically illustrate the multi-UAV beamforming performance in highly dynamic environments, while each result at a time slot averages 1000 trials. The system parameters in the simulations are listed as follows.

- The lead connects to $U = 3$ followers using SDMA at the same time and same frequency. The number of antennas $N = 16$ (4×4), and the SINR = 20 dB in the simulations for each follower is given by $\frac{|\rho_u|^2}{\sigma_z^2}$, where $|\rho_u|^2$ is the average receive power for follower u and $\sum_{u=1}^U |\rho_u|^2 = 1$.
- In the codebooks, the candidates for azimuth angles ϕ_{n_f} and ϕ_{n_w} are $\{15^\circ + n \cdot 30^\circ\}_{n=0}^{11}$, and the candidates for elevation angles θ_{n_f} and θ_{n_w} are $\{15^\circ + n \cdot 30^\circ\}_{n=0}^2$.
- The Q-learning parameters include the learning rate $\alpha = 0.5$, discount factor $\gamma = 0.5$, probability of ε -greedy action $\varepsilon = 0.1$, upper threshold $c_u = 1.1$, and lower threshold $c_l = 0.9$.
- The random walk channel model has normally distributed disturbance $\lambda_i \sim \mathcal{N}(0, \sigma_\lambda^2)$, where $\sigma_\lambda^2 = 4, 16$.

According to the number of all potential steering angles, the size of codebook \mathcal{F} at the lead should be 36. To alleviate the loading at the lead and speed up the convergence and learning rate, we group the followers into three zones in elevation angle (i.e., $0^\circ - 30^\circ$, $30^\circ - 60^\circ$, and $60^\circ - 90^\circ$), and each zone has three followers. Due to the space limitation, we only show the simulation results with three followers in the zone of elevation angle between 0° and 30° , and the codebook size of \mathcal{F} becomes $N_F = 12$, where the 12 candidates all have the same elevation angle $\theta_{n_f} = 15^\circ$.

A. Q-Learning and Gradient-Based Beam Tracking Methods

The first numerical result of the beam tracking in Fig. 5 is described by an example of the performance comparison of the proposed Q-learning and reference gradient-based tracking methods [20]. We use one realization of the time-varying AoA

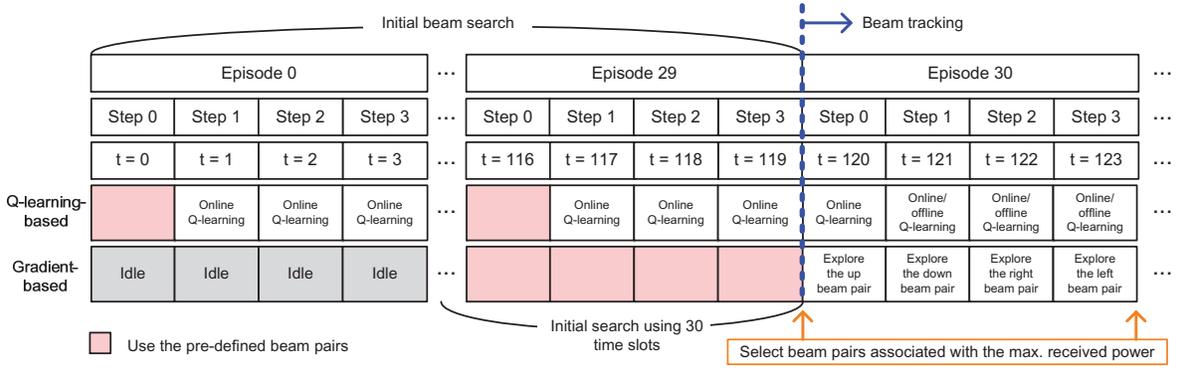


Fig. 6. Time frame defined by the episode and step. In the Q-learning-based beam tracking, we can use offline Q-learning in some steps for the purpose of overhead reduction (introduced in Subsections VI-B and VI-C).

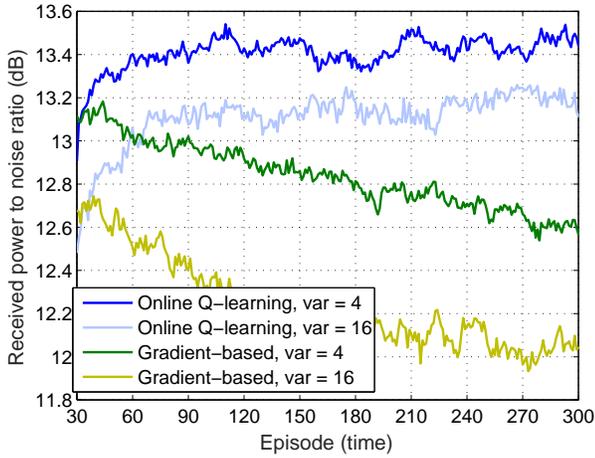


Fig. 7. Sum of the received power from $U = 3$ followers to noise ratio using the proposed Q-learning and reference gradient-based beam tracking methods with the variance of AoAs and AoDs $\sigma_\lambda^2 = 4, 16$. Compared with the gradient-based method, the Q-learning-based beam tracking can provide stable link quality over time.

and AoD with $\sigma_\lambda^2 = 16$ to explain the difference between these two methods.

At the beginning, both methods implement the initial beam search in the first 30 episodes or time slots⁷, where the time frame is sketched in Fig. 6. In the time frame, we assume that each episode includes 4 steps so that the gradient-based approach can evaluate the 4 neighboring beam pairs in an episode during the beam tracking. The gradient-based method uses 30 *time slots* to implement the initial beam search, while Q-learning method uses 30 *episodes* to implement the initial beam search and update the Q-values at each step. After the initial beam search, the beam tracking starts from Episode 30 with the state corresponding to the maximum received power obtained during the previous 30 episodes. The Q-learning method during the beam tracking may adopt online or offline Q-learning. To fairly compare with the reference method which gets the latest observations at each time slot,

⁷The 30 beam pairs for the initial beam search are uniformly chosen from a total of $N_F \times N_W = 12 \times 36 = 432$ potential beam pairs.

we use *online* Q-learning for all the steps in each episode to evaluate the proposed method in Fig. 5.

In Fig. 5, from Episode 30 to 80, the Q-learning-based beam tracking explored the range of AoA within $[60^\circ, 120^\circ]$ and the range of AoD within $[90^\circ, 170^\circ]$ using the beams steering to $\phi_{n_f} = 75^\circ, 105^\circ$ and $\phi_{n_w} = 105^\circ, 135^\circ, 165^\circ$, respectively.⁸ Q-learning records all the experience acquired during this time in the Q-table so that the agent uses current observations and the experience to predict the next beam pair. In such a way, it can stably track the appropriate beams. After Episode 80, there are probably not many data corresponding to the beam pairs with $\phi_{n_f} = 135^\circ, 165^\circ$ or $\phi_{n_w} = 45^\circ, 75^\circ$; therefore, it needs some time to update the Q-values as reference in the future. Next, let us look at the performance of the reference gradient-based scheme. The dilemma of gradient-based scheme is that it may get trapped into a local optimum and could only get out from it when AoA or AoD changes significantly. Q-learning method also finds the local optimal solution sometimes, but appropriate ϵ -greedy random actions can solve this problem. Moreover, compared with the gradient-based method, Q-learning has a *global* map (i.e., the Q-table), which provides useful information for beam tracking.

The performance comparison of the proposed and reference methods that support $U = 3$ followers simultaneously is shown in Fig. 7, where the received power is captured at the end of each episode as described in Fig. 6. Compared with the gradient-based method, the Q-learning-based scheme works stable over time, even when the variance σ_λ^2 of AoAs and AoDs is large. In terms of high-resilience demand for the multi-UAV system, the numerical results of the proposed method show that balancing the exploration and exploitation can outperform the one using exploration only. Although Q-learning needs some space and efforts to record the experience in the Q-table, it makes actions depending on not only current observations but also the experience and rewards so that the performance is not completely dominated by the current observations, while the gradient-based method totally relies on them.

⁸The beamwidth is around 30° ; ideally the beam switching occurs when AoA/AoD changes at $30^\circ, 60^\circ, \dots, 180^\circ$ in Fig. 5.

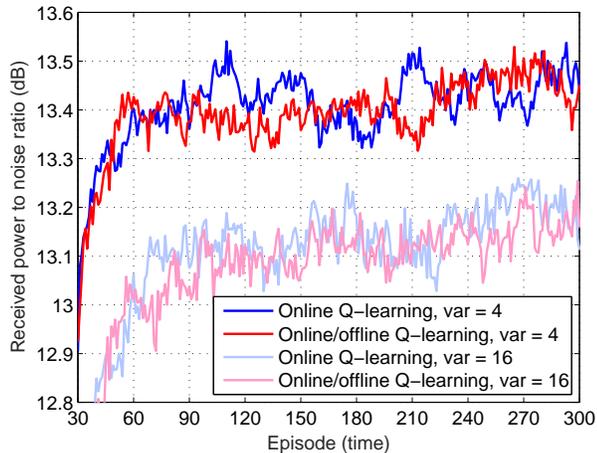


Fig. 8. Sum of the received power from $U = 3$ followers to noise ratio using the online and online/offline Q-learning-based beam tracking methods with $\sigma_\lambda^2 = 4, 16$. The curves of online Q-learning in this figure and Fig. 7 are identical. According to the results of overhead reduction in Fig. 9, using one or at most two steps per episode to track the beams is enough to maintain the link quality.

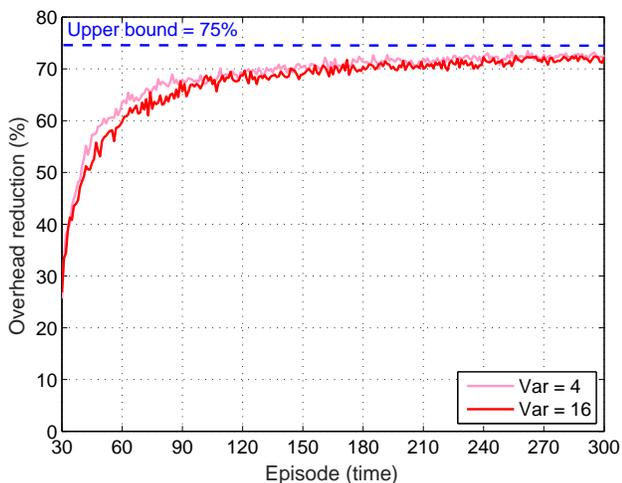


Fig. 9. Reduced overhead of pilot transmission using the online/offline Q-learning-based beam tracking with $\sigma_\lambda^2 = 4, 16$. In an episode, the first step is always used for pilot transmission so that the upper bound is 75%, while the other three steps may or may not be used for pilot transmission, depending on whether the corresponding states are explored or not.

B. Resilience and Efficiency of Offline Q-learning-Based Beam Tracking

In the previous subsection, we use *online* Q-learning to implement the beam tracking at all the steps in each episode in order to compare with the reference method. From the results shown in Fig. 7, we observe that online Q-learning provides stable link quality that can meet the high-resilience requirement for highly dynamic multi-UAV environments, but it also means that all the resources are used as pilot signals. From the perspective of system efficiency, it is inefficient design. Essentially, the trade-off between system efficiency and resilience has to be considered together. As a result, we introduce *offline* Q-learning that uses past observations to

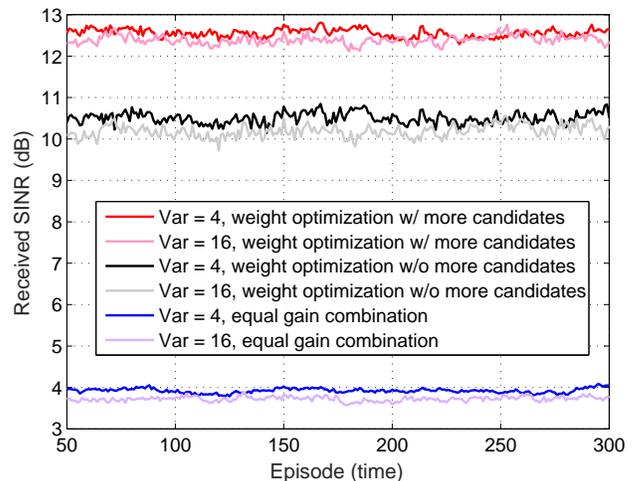


Fig. 10. Achievable SINR that uses equal gain or optimal weights to combine $U = 3$ analog beam pairs with $\sigma_\lambda^2 = 4, 16$. With and without the optimal digital weights, the difference in received SINR is 6.5 dB. With more than one candidate for the analog beam pairs, the received SINR can be further improved by at least 2 dB.

implement the beam tracking.

At the first step in each episode during the beam tracking, we let the followers transmit the pilot signals using the selected beam pairs, determined in the previous episodes, to update the observations. Therefore, the first step always adopts online Q-learning. In the rest three steps, the next state S_{t+1} (decided by Q-learning) may or may not be explored in the previous episodes. If the state was not explored, the agent still adopts online Q-learning in order to get the corresponding reward as well as Q-value. Instead, if the state was explored, Q-learning can use past observations to implement the beam tracking, which is offline Q-learning. However, we are not sure whether the next states in the rest three steps were explored. Therefore, the agent may adopt online or offline Q-learning, i.e., the case *online/offline* Q-learning in Figs. 6 and 8. In such a design, the upper bound of the reduced overhead of pilot transmission is 75%, see Fig. 9, since one of the four steps in an episode is dedicated to pilot transmission.

Figs. 8 and 9 show the comparison of the online and online/offline Q-learning-based beam tracking methods. The simulation results provide some interesting insights. After a certain time of exploration, the overhead of pilot transmission could be reduced up to 72% without any loss of performance. This implies that the experience stored in the Q-table provides enough information to solve the prediction problem, which is an advantage of machine learning.

C. SINR Maximization Using Digital Weights

We use SDMA to support multi-UAV communications at the price of spatial-domain interference. The goal of the digital beamforming is to minimize the interference plus noise given the selected analog beams. The performance is shown in Fig. 10, where the received SINR is obtained using $\hat{\mathbf{F}}_{P,t}$, $\hat{\mathbf{W}}_{P,t}$, and $\hat{\mathbf{F}}_{B,t}$.

First, the curves *equal gain combination* do not take into account more than one candidate for the analog beam pairs (i.e., both $\mathcal{I}_{\mathcal{F},t}$ and $\mathcal{I}_{\mathcal{W},t}$ have only one member respectively), and we let the digital beamforming equal to the $U \times U$ identity matrix (i.e., let $\tilde{\mathbf{F}}_{B,t} = \mathbf{I}_U$). Without trying to minimize the interference, the achievable SINR is only around 4 dB. When we design the digital weights to minimize the interference plus noise, see curves *w/o more candidates*, the SINR gain can be increased by around 6.5 dB. If we keep more candidates for the analog beams after the beam selection at the end of each episode (see Example 2), the SINR can be further improved. In curves *w/ more candidates*, reserving two candidates for each beam pair after the beam selection can have 2 dB gain in SINR. The reason is that the selected analog beams based on the received power do not ensure that the hybrid beamforming achieves the maximization SINR, even with the corresponding optimal digital weights.

Comparing the curves with low and high changes of the angle variables corresponding to low and high speeds of the UAVs in Fig. 10, we can see that the maximum difference in SINR is less than 0.5 dB. It shows that the proposed Q-learning-based hybrid beamforming is quite robust to the large variance of AoAs and AoDs.

VII. CONCLUSION

This paper solved a highly dynamic multi-UAV hybrid beamforming problem that is usually characterized by a high-resilience constraint. To meet the constraint, we apply Q-learning method to mmWave hybrid beamforming systems. Moreover, in a dynamic environment, how to efficiently obtain and use the observations matters to the beamforming performance. In the proposed analog beam tracking approach, we use current and past observations together with the designed rewards to solve the prediction problem. The numerical results show that the proposed method significantly increases the efficiency of data transmission and beam switching. To optimally combine the analog beams in a manner of SINR maximization, we present the solution of digital weights using the coupling coefficients given the selected beams. The solution can be simply extended to the case with more candidates for analog beams to further improve the received SINR.

APPENDIX A

DESIGN OF UPPER AND LOWER THRESHOLDS (c_u, c_l)

In the observation equation (13), given a channel matrix, the estimate of the power of signal $y_{u,t}(n_f, n_w)$ can be represented by

$$\begin{aligned} |y_{u,t}(n_f, n_w)|^2 &= \left| \tilde{\mathbf{f}}_{n_f}^H \mathbf{H}_{u,t} \tilde{\mathbf{w}}_{n_w} + z_t \right|^2 \\ &= \left| \tilde{\mathbf{f}}_{n_f}^H \mathbf{H}_{u,t} \tilde{\mathbf{w}}_{n_w} \right|^2 + \varepsilon_u(n_w, n_f) + \zeta, \end{aligned} \quad (33)$$

where $\left| \tilde{\mathbf{f}}_{n_f}^H \mathbf{H}_{u,t} \tilde{\mathbf{w}}_{n_w} \right|^2$ is a constant based on the given channel state, and the other two terms are given as follows. First,

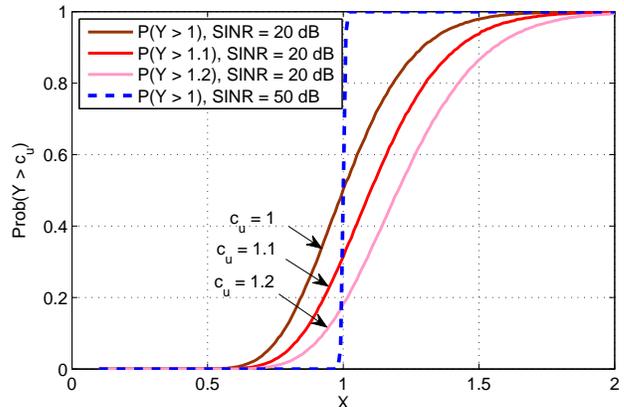


Fig. 11. The probabilities that the estimate of X is greater than the upper threshold c_u . Increasing the value of the upper threshold c_u can reduce the probability that the agent get a fail reward at SINR = 20 dB.

$\varepsilon_u(n_w, n_f)$ follows a normal distribution with mean zero and variance $2\sigma_z^2 \left| \tilde{\mathbf{f}}_{n_f}^H \mathbf{H}_{u,t} \tilde{\mathbf{w}}_{n_w} \right|^2$ given by

$$\begin{aligned} \varepsilon_u(n_w, n_f) &= 2 \Re \left(\tilde{\mathbf{f}}_{n_f}^H \mathbf{H}_{u,t} \tilde{\mathbf{w}}_{n_w} \right) \Re(z_t) + 2 \Im \left(\tilde{\mathbf{f}}_{n_f}^H \mathbf{H}_{u,t} \tilde{\mathbf{w}}_{n_w} \right) \Im(z_t) \\ &\sim \mathcal{N} \left(0, 2\sigma_z^2 \left| \tilde{\mathbf{f}}_{n_f}^H \mathbf{H}_{u,t} \tilde{\mathbf{w}}_{n_w} \right|^2 \right), \end{aligned} \quad (34)$$

and ζ follows a gamma distribution with shape parameter 1 and scale parameter σ_z^2 :

$$\zeta = \Re(z_t)^2 + \Im(z_t)^2 \sim \Gamma(1, \sigma_z^2). \quad (35)$$

Due to the fact that it is not able to obtain a closed-form expression for the density function of $\varepsilon_u(n_w, n_f) + \zeta$, we use a Monte Carlo method to find appropriate upper and lower thresholds. First, let us define the ratio of the power of coupling coefficients at time $t+1$ and t by

$$X = \frac{\left| \tilde{\mathbf{f}}_{n'_f}^H \mathbf{H}_{u,t+1} \tilde{\mathbf{w}}_{n'_w} \right|^2}{\left| \tilde{\mathbf{f}}_{n_f}^H \mathbf{H}_{u,t} \tilde{\mathbf{w}}_{n_w} \right|^2}, \quad (36)$$

where (n'_f, n'_w) is the beam index pair used at time $t+1$. In addition, the ratio of the received power at time $t+1$ and t is given by

$$\begin{aligned} Y &= \frac{|y_{u,t+1}(n'_f, n'_w)|^2}{|y_{u,t}(n_f, n_w)|^2} \\ &= \frac{\left| \tilde{\mathbf{f}}_{n'_f}^H \mathbf{H}_{u,t+1} \tilde{\mathbf{w}}_{n'_w} \right|^2 + \varepsilon_u(n'_w, n'_f) + \zeta_1}{\left| \tilde{\mathbf{f}}_{n_f}^H \mathbf{H}_{u,t} \tilde{\mathbf{w}}_{n_w} \right|^2 + \varepsilon_u(n_w, n_f) + \zeta_2}, \end{aligned} \quad (37)$$

where ζ_1 and ζ_2 follow the same Gamma distribution $\Gamma(1, \sigma_z^2)$.

When SINR = 50 dB, the noise variance σ_z^2 is pretty small so that we have $X \approx Y$, and the reward from the environment could be either positive (+1) or negative (-1). Therefore, it is fine to let $c_u = c_l = 1$. Ideally, $X < 1$ should lead to a negative reward, that is, $\text{Prob}(Y > c_u = 1) = 0$, as shown in the curve SINR = 50 dB in Fig. 11

In the case of SINR = 20 dB, the noise effect on the received power becomes serious. If we still assume that $c_u = 1$ at SINR = 20 dB, we can find that the probability that the agent get a positive reward when $X < 1$ is greater than 0. For example, when $X = 0.9$, the probability that the agent get a positive reward is $Prob(Y > c_u = 1) = 0.35$. The objective of the upper and lower thresholds are used to limit the reward from the environment when the values of the received power are unreliable. Increasing the value of c_u can effectively decrease this kind of error probability. However, it does not make sense to let c_u be very large because it will make the reward equal to zero even when $X > 1$, which is not beneficial for Q-learning. In the same manner, we can adjust the value of the other threshold c_l to reduce the probability of getting the fail reward.

APPENDIX B DERIVATION OF (31)

The objective function of the problem (31) is the generalized Rayleigh quotient [39]. To convert the problem of maximizing SINR to a simpler one of maximizing a normalized quadratic form, we define a vector $\tilde{\mathbf{x}}_u = \mathbf{B}_{u,t}^{-0.5} \mathbf{x}_u$, which is equivalent to $\mathbf{x}_u = \mathbf{B}_{u,t}^{-0.5} \tilde{\mathbf{x}}_u$. Replacing \mathbf{x}_u with $\mathbf{B}_{u,t}^{-0.5} \tilde{\mathbf{x}}_u$, the objective function of the problem becomes

$$\frac{\tilde{\mathbf{x}}_u^H \mathbf{B}_{u,t}^{-0.5} \mathbf{A}_{u,t} \mathbf{B}_{u,t}^{-0.5} \tilde{\mathbf{x}}_u}{\|\tilde{\mathbf{x}}_u\|_2^2}. \quad (38)$$

To maximize (38) is equivalent to maximize the numerator. Let $\tilde{\mathbf{x}}_u$ be the eigenvector of $\mathbf{B}_{u,t}^{-0.5} \mathbf{A}_{u,t} \mathbf{B}_{u,t}^{-0.5}$ corresponding to the maximum eigenvalue, the maximum value of (38) is therefore given by

$$\max_{\tilde{\mathbf{x}}_u} \frac{\tilde{\mathbf{x}}_u^H \mathbf{B}_{u,t}^{-0.5} \mathbf{A}_{u,t} \mathbf{B}_{u,t}^{-0.5} \tilde{\mathbf{x}}_u}{\|\tilde{\mathbf{x}}_u\|_2^2} = \lambda_{\max}(\mathbf{B}_{u,t}^{-0.5} \mathbf{A}_{u,t} \mathbf{B}_{u,t}^{-0.5}), \quad (39)$$

where $\lambda_{\max}(\mathbf{B}_{u,t}^{-0.5} \mathbf{A}_{u,t} \mathbf{B}_{u,t}^{-0.5})$ is the maximum eigenvalue of $\mathbf{B}_{u,t}^{-0.5} \mathbf{A}_{u,t} \mathbf{B}_{u,t}^{-0.5}$. As a result, we have the optimal solution of \mathbf{x}_u subject to the constraint $\|\mathbf{x}_u\|_2 = 1$:

$$\begin{aligned} \mathbf{x}_u^* &= \frac{\mathbf{B}_{u,t}^{-0.5} \tilde{\mathbf{x}}_u}{\|\mathbf{B}_{u,t}^{-0.5} \tilde{\mathbf{x}}_u\|_2} \\ &= \frac{\mathbf{B}_{u,t}^{-0.5} \mathbf{e}_{\max}(\mathbf{B}_{u,t}^{-0.5} \mathbf{A}_{u,t} \mathbf{B}_{u,t}^{-0.5})}{\|\mathbf{B}_{u,t}^{-0.5} \mathbf{e}_{\max}(\mathbf{B}_{u,t}^{-0.5} \mathbf{A}_{u,t} \mathbf{B}_{u,t}^{-0.5})\|_2}, \end{aligned} \quad (40)$$

where $\mathbf{e}_{\max}(\mathbf{B}_{u,t}^{-0.5} \mathbf{A}_{u,t} \mathbf{B}_{u,t}^{-0.5})$ is the dominant eigenvector of $\mathbf{B}_{u,t}^{-0.5} \mathbf{A}_{u,t} \mathbf{B}_{u,t}^{-0.5}$.

REFERENCES

- [1] Q. Nadeem, A. Kammoun, A. Chaaban, M. Debbah, and M. Alouini, "Asymptotic analysis of large intelligent surface assisted mimo communication," *Submitted to IEEE Trans. Wireless Commun.*, 2019. [Online]. Available: <https://arxiv.org/pdf/1903.08127.pdf>
- [2] N. Hossein Motlagh, T. Taleb, and O. Arouk, "Low-altitude unmanned aerial vehicles-based internet of things services: Comprehensive survey and future perspectives," *IEEE IoT J.*, vol. 3, no. 6, pp. 899–922, Dec 2016.
- [3] I. Uluturk, I. Uysal, and K. Chen, "Efficient 3d placement of access points in an aerial wireless network," in *IEEE Annu. Consumer Commun. Neww. Conf. (CCNC)*, Las Vegas, NV, USA, Jan. 2019.
- [4] S. Waharte and N. Trigoni, "Supporting search and rescue operations with uavs," in *Int. Conf. on Emerg. Security Technol.*, Canterbury, UK, Sep. 2010, pp. 142–147.
- [5] Y. Zeng, R. Zhang, and T. J. Lim, "Wireless communications with unmanned aerial vehicles: opportunities and challenges," *IEEE Commun. Mag.*, vol. 54, no. 5, pp. 36–42, May 2016.
- [6] C. Fiandrino, H. Assasa, P. Casari, and J. Widmer, "Scaling millimeter-wave networks to dense deployments and dynamic environments," *Proc. IEEE*, vol. 107, no. 4, pp. 732–745, Apr. 2019.
- [7] T. Rappaport, R. Heath, R. Daniels, and J. Murdock, *Millimeter Wave Wireless Communications*. Prentice Hall, 2014.
- [8] M. Xiao, S. Mumtaz, Y. Huang, L. Dai, Y. Li, M. Matthaiou, G. K. Karagiannidis, E. Björnson, K. Yang, C. I, and A. Ghosh, "Millimeter wave communications for future mobile networks," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 9, pp. 1909–1935, Sep. 2017.
- [9] M. K. Samimi and T. S. Rappaport, "3-D statistical channel model for millimeter-wave outdoor mobile broadband communications," in *IEEE Int. Conf. on Commun. (ICC)*, London, UK, Jun. 2015, pp. 2430–2436.
- [10] 3GPP TR 38.900 V14.3.1, "Study on channel model for frequency spectrum above 6 GHz (Release 14)," Tech. Rep., 2017.
- [11] J. Liberti and T. Rappaport, *Smart antennas for wireless communications: IS-95 and third generation CDMA applications*. Prentice Hall, 1999.
- [12] A. Hajimiri, H. Hashemi, A. Natarajan, X. Guan, and A. Komijani, "Integrated phased array systems in silicon," *Proc. IEEE*, vol. 93, no. 9, pp. 1637–1655, Sep. 2005.
- [13] E. Björnson, L. Van der Perre, S. Buzzi, and E. G. Larsson, "Massive mimo in sub-6 Ghz and mmwave: Physical, practical, and use-case differences," *IEEE Wireless Commun.*, vol. 26, no. 2, pp. 100–108, Apr. 2019.
- [14] X. Zhang, A. F. Molisch, and S.-Y. Kung, "Variable-phase-shift-based RF-baseband codesign for MIMO antenna selection," *IEEE Trans. Signal Process.*, vol. 53, no. 11, pp. 4091–4103, Nov. 2005.
- [15] O. E. Ayach, S. Rajagopal, S. Abu-Surra, Z. Pi, and R. W. Heath, "Spatially sparse precoding in millimeter wave MIMO systems," *IEEE Trans. Wireless Commun.*, vol. 13, no. 3, pp. 1499–1513, Mar. 2014.
- [16] H. Chiang, W. Rave, T. Kadur, and G. Fettweis, "A low-complexity beamforming method by orthogonal codebooks for millimeter wave links," in *IEEE Int. Conf. on Acoust., Speech and Signal Process. (ICASSP)*, New Orleans, LA, USA, Mar. 2017, pp. 3375 – 3379.
- [17] —, "Hybrid beamforming based on implicit channel state information for millimeter wave links," *IEEE J. Sel. Top. Signal Process.*, vol. 12, no. 2, pp. 326–339, May 2018.
- [18] A. Fotouhi, H. Qiang, M. Ding, M. Hassan, L. G. Giordano, A. Garcia-Rodriguez, and J. Yuan, "Survey on uav cellular communications: Practical aspects, standardization advancements, regulation, and security challenges," *IEEE Commun. Surveys Tutorials*, 2019.
- [19] P. Zhou, X. Fang, Y. Fang, R. He, Y. Long, and G. Huang, "Beam management and self-healing for mmwave uav mesh networks," *IEEE Trans. Veh. Technol.*, vol. 68, no. 2, pp. 1718–1732, Feb. 2019.
- [20] T. Kadur, W. Rave, H. Chiang, and G. Fettweis, "Experimental validation of a robust beam alignment algorithm in an indoor environment," in *IEEE Wireless Commun. and Networking Conf. (WCNC)*, Barcelona, Spain, Apr. 2018.
- [21] M. B. Booth, V. Suresh, N. Michelusi, and D. J. Love, "Multi-armed bandit beam alignment and tracking for mobile millimeter wave communications," *IEEE Commun. Lett.*, vol. 23, no. 7, pp. 1244–1248, July 2019.
- [22] V. Va, H. Vikalo, and R. W. Heath, "Beam tracking for mobile millimeter wave communication systems," in *IEEE Global Conf. on Signal and Inf. Process. (GlobalSIP)*, Washington, DC, USA, Dec 2016, pp. 743–747.
- [23] J. Supancic and D. Ramanan, "Tracking as online decision-making: Learning a policy from streaming videos with reinforcement learning," in *IEEE Int. Conf. on Computer Vision (ICCV)*, Venice, Italy, Oct. 2017, pp. 322–331.
- [24] A. Klautau, P. Batista, N. González-Prelcic, Y. Wang, and R. W. Heath, "5G mimo data for machine learning: Application to beam-selection using deep learning," in *Inform. Theory and Applicat. Workshop (ITA)*, San Diego, CA, USA, Feb. 2018.
- [25] Y. Chen, W. Cheng, and L. Wang, "Learning-assisted beam search for indoor mmwave networks," in *IEEE Wireless Commun. and Netw. Conf. Workshops (WCNCW)*, Barcelona, Spain, Apr. 2018, pp. 320–325.
- [26] Y. Ke, H. Gao, W. Xu, L. Li, L. Guo, and Z. Feng, "Position prediction based fast beam tracking scheme for multi-user uav-mmwave communications," in *IEEE Int. Conf. on Commun. (ICC)*, Shanghai, China, May 2019.

- [27] C. J. C. H. Watkins and P. Dayan, "Q-learning," *Machine Learning*, vol. 8, no. 3, pp. 279–292, May 1992. [Online]. Available: <https://doi.org/10.1007/BF00992698>
- [28] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. The MIT Press, 2018. [Online]. Available: <http://incompleteideas.net/book/the-book-2nd.html>
- [29] L. Li, H. Ren, Q. Cheng, K. Xue, W. Chen, M. Debbah, and Z. Han, "Millimeter-wave networking in sky: A machine learning and mean field game approach for joint beamforming and beam-steering," *Submitted to IEEE J. Sel. Areas Commun.*, 2019. [Online]. Available: <http://www.laneas.com/sites/default/files/publications/4223/JSAC-v9.2.pdf>
- [30] A. Alkhateeb, G. Leus, and R. W. Heath, "Limited feedback hybrid precoding for multi-user millimeter wave systems," *IEEE Trans. Wireless Commun.*, vol. 14, no. 11, pp. 6481–6494, Nov. 2015.
- [31] F. Sohrabi and W. Yu, "Hybrid digital and analog beamforming design for large-scale antenna arrays," *IEEE J. Sel. Topics Signal Process.*, vol. 10, no. 3, pp. 501–513, Apr. 2016.
- [32] H. Yin and H. Liu, "Performance of space-division multiple-access (SDMA) with scheduling," *IEEE Trans. Wireless Commun.*, vol. 1, no. 4, pp. 611–618, Oct. 2002.
- [33] A. Adhikary, E. A. Safadi, M. K. Samimi, R. Wang, G. Caire, T. S. Rappaport, and A. F. Molisch, "Joint spatial division and multiplexing for mm-wave channels," *IEEE J. Sel. Areas Commun.*, vol. 32, no. 6, pp. 1239–1255, Jun. 2014.
- [34] C. A. Balanis, *Antenna Theory: Analysis and Design*. Wiley-Interscience, 2005.
- [35] A. M. Sayeed, "Deconstructing multiantenna fading channels," *IEEE Trans. Signal Process.*, vol. 50, no. 10, pp. 2563–2579, Oct 2002.
- [36] L. Yang and W. Zhang, "Beam tracking and optimization for uav communications," *IEEE Trans. Wireless Commun. (Early Access)*, 2019.
- [37] M. Tan, "Multi-agent reinforcement learning: Independent vs. cooperative agents," in *Proc. 10th Intl. Conf. on Machine Learning*. Amherst, MA, USA: Morgan Kaufmann Publishers Inc., 1993, pp. 330–337.
- [38] K. Chen and H. Hung, "Wireless robotic communication for collaborative multi-agent systems," in *IEEE Int. Conf. on Commun. (ICC)*, Shanghai, China, May 2019.
- [39] C. D. Meyer, *Matrix Analysis and Applied Linear Algebra*, C. D. Meyer, Ed. Philadelphia, PA, USA: Society for Industrial and Applied Mathematics, 2000.