

Deep Reinforcement Learning Based on Location-Aware Imitation Environment for RIS-Aided mmWave MIMO Systems

Wangyang Xu, Jiancheng An, Chongwen Huang, Lu Gan, and Chau Yuen,
Fellow, IEEE

Abstract

Reconfigurable intelligent surface (RIS) has recently gained popularity as a promising solution for improving the signal transmission quality of wireless communications with less hardware cost and energy consumption. This letter offers a novel deep reinforcement learning (DRL) algorithm based on a location-aware imitation environment for the joint beamforming design in an RIS-aided mmWave multiple-input multiple-output system. Specifically, we design a neural network to imitate the transmission environment based on the geometric relationship between the user's location and the mmWave channel. Following this, a novel DRL-based method is developed that interacts with the imitation environment using the easily available location information. Finally, simulation results demonstrate that the proposed DRL-based algorithm provides more robust performance without excessive interaction overhead compared to the existing DRL-based approaches.

Index Terms

The work of Prof. Huang was supported by the China National Key R&D Program under Grant 2021YFA1000072, National Natural Science Foundation of China under Grant 62101492, Zhejiang Provincial Natural Science Foundation of China under Grant R22F0110230, Zhejiang University Education Foundation Qizhen Scholar Foundation, and Fundamental Research Funds for the Central Universities under Grant 2021FZZX001-21. The work of Lu Gan was supported by Yibin Science and Technology Program under Grant 2020FW007. (*Corresponding author: Lu Gan*).

W. Xu, J. An, and L. Gan are with the School of Information and Communication Engineering, University of Electronic Science and Technology of China (UESTC), Chengdu, Sichuan, 611731, China. Lu Gan is also with the Yibin Institute of UESTC, Yibin, Sichuan, 643000, China (E-mail: wangyangxu@std.uestc.edu.cn; jiancheng_an@163.com; ganlu@uestc.edu.cn).

C. Huang is with College of Information Science and Electronic Engineering, Zhejiang University, Hangzhou 310027, China, and with International Joint Innovation Center, Zhejiang University, Haining 314400, China, and also with Zhejiang-Singapore Innovation and AI Joint Research Lab and Zhejiang Provincial Key Laboratory of Info. Proc., Commun. & Netw. (IPCAN), Hangzhou 310027, China. (E-mail: [HYPERLINK "mailto:chongwenhuang@zju.edu.cn"](mailto:chongwenhuang@zju.edu.cn) chongwenhuang@zju.edu.cn).

C. Yuen is with Engineering Product Development (EPD) Pillar, Singapore University of Technology and Design, Singapore 487372, Singapore (E-mail: yuenchau@sutd.edu.sg).

Reconfigurable intelligent surface, deep reinforcement learning, imitation environment.

I. INTRODUCTION

Recently, reconfigurable intelligent surface (RIS) is emerged as a promising technology that significantly improves the system throughput, spectrum efficiency, and energy efficiency of wireless networks [1]. An RIS is equipped with a large number of hardware-efficient and nearly passive reflecting elements, which does not employ active radio frequency chains hence significantly reducing the energy consumption and hardware.

Nevertheless, the joint transmit beamforming and reflection coefficients design constitutes a challenge in RIS-aided mmWave multiple-input multiple-output (MIMO) systems. In [1]–[3], the joint beamforming has been investigated under the consideration of various optimization objectives and phase shift models, where several algorithms were applied such as semidefinite relaxation (SDR) [1], alternating optimization (AO) [2], and block coordinate descent (BCD) [3].

Additionally, deep learning (DL) techniques have been employed to gain various advantages, such as end-to-end, model-free, and data-driven optimizations [4]–[6]. Among these DL techniques, the deep reinforcement learning (DRL) enables efficient algorithm designs by observing the rewards from the environment and solving sophisticated optimization problems in the RIS-aided systems [5], [6]. Unlike supervised learning, DRL does not require any labels and is capable of adapting dynamic environments. Nonetheless, channel's variation remains a hurdle to the robust performance of the DRL. Furthermore, most DRL-based solutions are designed for the multiple-input single-output (MISO) systems and require the channel state information (CSI) as the actor network input, which is challenging to implement for RIS equipped with many passive elements.

In this letter, we study the problem of joint beamforming in RIS-aided mmWave MIMO wireless communications and propose a DRL algorithm based on the location-aware imitation environment network (IEN). In contrast to most previous works, the proposed algorithm employs the readily-available user's location information rather than the accurate CSI. In addition, to improve DRL's poor robustness facing diverse channels, a deep neural network (DNN) is built to imitate the actual environment for decreasing the excessive overhead imposed by DRL's interaction with the actual environment. Finally, simulation results are provided to verify our proposed algorithm.

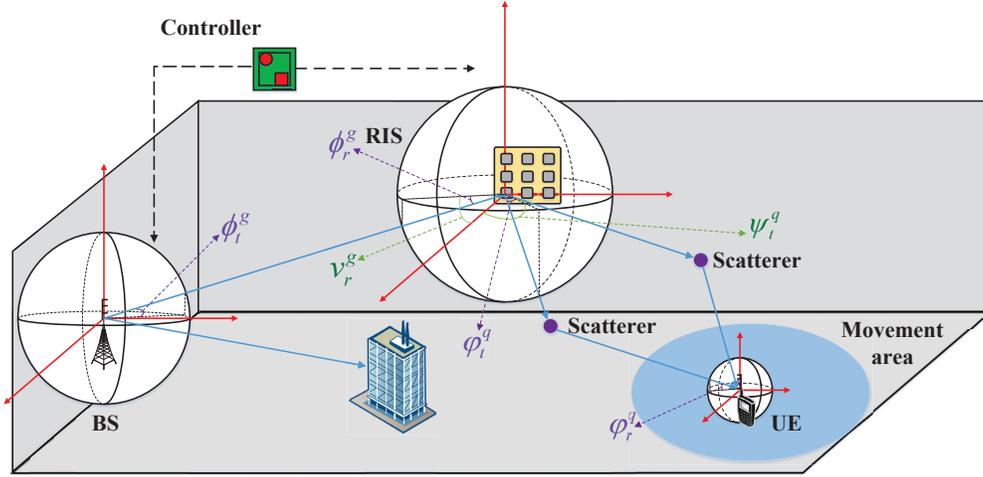


Fig. 1. An RIS-aided mmWave MIMO communication system.

II. SYSTEM MODEL AND PROBLEM FORMULATION

A. System Model

As illustrated in Fig. 1, we consider a downlink RIS-aided mmWave MIMO communication system, where an RIS equipped with N passive reflecting elements in a uniform planar array (UPA) is deployed for enhancing the transmission between the base station (BS) with M antennas and the user equipment (UE) with K antennas. The antennas of BS and UE are both arranged in the form of uniform linear array (ULA). Each RIS element is capable of rescattering the impinging signals with an individual phase shift, which can be dynamically adjusted by the RIS controller. Furthermore, we consider the narrowband communications over quasi-static block-fading channels. Let $\mathbf{G} \in \mathbb{C}^{N \times M}$ and $\mathbf{H}_r \in \mathbb{C}^{K \times N}$ denote the channel matrices from BS to RIS, and from RIS to UE, respectively. The direct link between the BS and the UE is assumed to be blocked by obstacles.

Specifically, the mmWave channel is characterized by the classic Saleh-Valenzuela model [7]. Hence, \mathbf{G} and \mathbf{H}_r are generated by

$$\mathbf{G} = \sqrt{\frac{MN}{L_G}} \sum_{g=1}^{L_G} \sqrt{PL_g} \mathbf{a}_P(\psi_{G,r}^g, \phi_{G,r}^g) \mathbf{a}_L^H(\phi_{G,t}^g), \quad (1)$$

$$\mathbf{H}_r = \sqrt{\frac{NK}{L_D}} \sum_{d=1}^{L_D} \sqrt{PL_d} \mathbf{a}_L(\phi_{h,r}^d) \mathbf{a}_P^H(\psi_{h,t}^d, \phi_{h,t}^d), \quad (2)$$

where L_G (L_D) and PL_g (PL_d) denote the multi-path number and the complex gain, respectively. $\psi_{G,r}^g$ ($\phi_{G,r}^g$) and $\phi_{G,t}^g$ ($g = 1, 2, \dots, L_G$) are the azimuth (elevation) angle of arrival (AoA), and

elevation angle of departure (AoD) of the g -th path of \mathbf{G} . Meanwhile, $\phi_{h,r}^d$ and $\psi_{h,t}^d$ ($\phi_{h,t}^d$) ($d = 1, 2, \dots, L_Q$) represent the elevation AoA and azimuth (elevation) AoD of \mathbf{H}_r , respectively. \mathbf{a}_L and \mathbf{a}_P denote the array response vectors of a half-wavelength spaced ULA and UPA, which are given by

$$\mathbf{a}_L(\phi) = \frac{1}{\sqrt{N_L}} [1, \dots, e^{j\pi n_l \sin \phi}, \dots, e^{j\pi(N_L-1) \sin \phi}]^T, \quad (3)$$

$$\mathbf{a}_P(\psi, \phi) = \frac{1}{\sqrt{N_x N_y}} [1, \dots, e^{j\pi(n_x \sin \psi \cos \phi + n_y \sin \phi)}, \dots, e^{j\pi((N_x-1) \sin \psi \cos \phi + (N_y-1) \sin \phi)}]^T, \quad (4)$$

where N_L is the antenna number of the ULA; N_x and N_y are the number of horizontal and vertical antennas of the UPA; n_l , n_x and n_y are the corresponding antenna indices; ψ and ϕ are the azimuth and elevation angles, respectively.

In the downlink transmission phase, the baseband signal $\mathbf{y} \in \mathbb{C}^{K \times 1}$ received at the UE can be expressed as

$$\mathbf{y} = \mathbf{H}_r \Theta \mathbf{G} \mathbf{x} + \mathbf{n}, \quad (5)$$

where $\mathbf{x} \in \mathbb{C}^{M \times 1}$ is the transmitted signal sequence; $\mathbf{n} \in \mathbb{C}^{K \times 1}$ is the additive white Gaussian noise satisfying $\mathbf{n} \sim \mathcal{CN}(0, \sigma^2 \mathbf{I}_K)$; $\Theta = \text{diag}(\boldsymbol{\theta})$ is the reflection coefficient matrix at the RIS, with $\boldsymbol{\theta} = \text{diag}(e^{j\varphi_1}, e^{j\varphi_2}, \dots, e^{j\varphi_N})$ and φ_l is the phase shift of the n -th RIS element. $\text{diag}(\cdot)$ denotes the diagonal operation. For simplicity, we assume that the phase shifts of each RIS element are continuously adjustable in the interval $[0, 2\pi)$.

B. Problem Formulation

In this letter, we aim to jointly design the transmit signal covariance matrix at the BS and reflection coefficient vector at the RIS for maximizing the achievable rate of RIS-aided mmWave MIMO systems. Accordingly, the optimization problem can be formulated as

$$\begin{aligned} \max_{\mathbf{Q}, \boldsymbol{\theta}} \quad & R = \log_2 \det \left(\mathbf{I}_K + \frac{\bar{\mathbf{H}} \mathbf{Q} \bar{\mathbf{H}}^H}{\sigma^2} \right) \\ \text{s.t.} \quad & \theta_n = e^{j\varphi_n}, \varphi_n \in [0, 2\pi), \quad \forall n = 1, 2, \dots, N, \\ & \text{tr}(\mathbf{Q}) \leq p, \mathbf{Q} \succeq 0, \end{aligned} \quad (6)$$

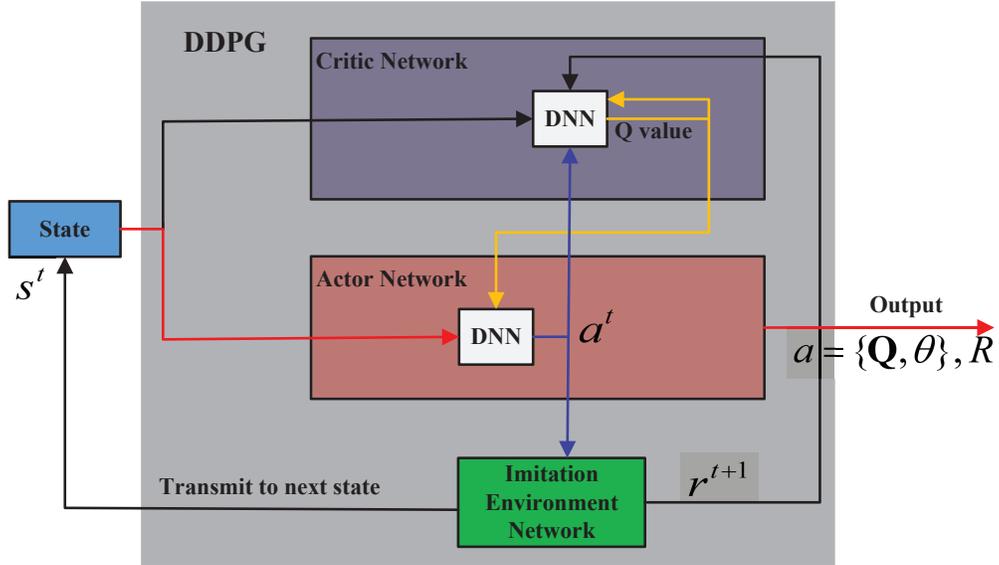


Fig. 2. The proposed DRL-based algorithm with IEN.

where $\bar{\mathbf{H}} = \mathbf{H}_r \Theta \mathbf{G}$ is the composite channel, $\mathbf{Q} \triangleq \mathbb{E}[\mathbf{x}\mathbf{x}^H]$ denotes the transmit signal covariance matrix, and we consider an average sum power constraint at the transmitter given by $\mathbb{E}[\|\mathbf{x}\|^2] \leq p$, which is equivalent to $\text{tr}(\mathbf{Q}) \leq p$.

We note that (6) is a non-convex problem. Although several novel approaches have been proposed for solving this complex problem [8], the BS generally requires prior knowledge of CSI, even in most data-driven DRL-based algorithms designed for the MISO systems. Furthermore, CSI acquisition remains a tremendous difficulty due to massive passive RIS elements. Besides, the robust performance of these DRL-based algorithms is also poor as UE moves.

III. DRL-BASED ALGORITHM WITH LOCATION-AWARE IEN

In this section, we propose a DRL-based algorithm with location-aware imitation environment for the joint beamforming design of RIS-aided mmWave MIMO systems. As shown in Fig. 2, the proposed DRL-based algorithm mainly uses the off-policy deep deterministic policy gradient (DDPG) network [9]. Besides, an IEN is proposed to replace the actual interaction environment.

A. The Imitation Environment Network

As shown in Fig. 1, BS and RIS are generally placed at fixed positions, leading the slowly time-varying BS-RIS channel \mathbf{G} . By contrast, the RIS-UE channel \mathbf{H}_r varies rapidly owing to the mobility of UE. However, due to the high attenuation in mmWave band, \mathbf{G} and \mathbf{H}_r rely

heavily on the locations of the BS, the RIS, and the UE. As a result, it is possible to recover the transmission channel from the BS to UE by employing only the UE's location information.

Since the RIS is generally installed at a high place for facilitating LoS propagation with BS/UE, it is reasonable to assume that the RIS serves a single UE with several unknown but fixed scatterers between it and BS/UE.

First, we build a DNN to imitate the actual transmission environment. Specifically, the IEN consists of two neural networks (NNs), termed the BS-RIS NN and RIS-UE NN, respectively, to provide a general paradigm for different cases of scatterers and NNs input-output mappings. The BS-RIS NN contains three dense layers with 128, 64, and $2MN$ neurons, respectively. The RIS-UE NN has the same structure as the BS-RIS NN, except that the neurons number of the third dense layer is $2KN$. The activation functions of the first two dense layers of these two NNs are tanh, while the third one is the linear function. The inputs of these two NNs are the tuples consisting of the 3D coordinates of the devices, $(x_{BS}, y_{BS}, z_{BS}, x_{RIS}, y_{RIS}, z_{RIS})$ and $(x_{RIS}, y_{RIS}, z_{RIS}, x_{UE}, y_{UE}, z_{UE})$. The outputs of the two NNs, $\text{vec}(\text{Re}\{\hat{\mathbf{G}}\}, \text{Im}\{\hat{\mathbf{G}}\})$ and $\text{vec}(\text{Re}\{\hat{\mathbf{H}}_r\}, \text{Im}\{\hat{\mathbf{H}}_r\})$, are the vectors of the real and imaginary parts of the predicted \mathbf{G} and \mathbf{H}_r because the NN can only handle real-valued data.

Next, we will introduce the details of the training set. We assume that the UE moves in a limited area served by the RIS. Therefore, it is possible to obtain the historical location data of potential UEs by some positioning technologies such as GPS. Specifically, we select U different historic locations of the UE's movement area, and then dynamically and randomly adjust the RIS reflection coefficient vector F times for each location. As a result, we can get a total of UF training inputs in the form of $\{loc_{BS}, loc_{RIS}, loc_{UE}^u, \boldsymbol{\theta}_f\}_{u=1,2,\dots,U; f=1,2,\dots,F}$. Note that loc_{BS} , loc_{RIS} , and loc_{UE}^u consisting of their 3D coordinates, which will be transformed to the form of the inputs of the BS-RIS NN and BS-RIS NN. Finally, we generate UF corresponding composite channels $\bar{\mathbf{H}}$ as training labels via the traditional channel estimation method [10]. We note that there may exist several inaccuracies in the data collection for all locations and composite channels, which will be left for our future research.

During the training phase, we first obtain the outputs, $\text{vec}(\text{Re}\{\hat{\mathbf{G}}\}, \text{Im}\{\hat{\mathbf{G}}\})$ and $\text{vec}(\text{Re}\{\hat{\mathbf{H}}_r\}, \text{Im}\{\hat{\mathbf{H}}_r\})$. By reconstructing the complex $\hat{\mathbf{G}}$ and $\hat{\mathbf{H}}_r$, the predicted composite channel is calculated by $\hat{\mathbf{H}} = \hat{\mathbf{H}}_r \Theta \hat{\mathbf{G}}$. The training of the IEN is based on stochastic gradient descent (SGD), and the

Algorithm 1 The Training of IEN

Input: $\{loc_{BS}, loc_{RIS}, loc_{UE}, \theta, \bar{\mathbf{H}}\}_{v,v=1,2,\dots,UF}$.

Output: The trained IEN.

- 1: **for** epoch $i = 1, 2, \dots, E$ **do**
 - 2: Construct $IN_{BR} = (x_{BS}, y_{BS}, z_{BS}, x_{RIS}, y_{RIS}, z_{RIS})$ and $IN_{RU} = (x_{RIS}, y_{RIS}, z_{RIS}, x_{UE}, y_{UE}, z_{UE})$ by $\{loc_{BS}, loc_{RIS}, loc_{UE}\}$;
 - 3: **for** $v = 1, 2, \dots, V$ **do**
 - 4: Input $IN_{BR,v}$ to the BS-RIS NN and output $\hat{\mathbf{G}}_v$;
 - 5: Input $IN_{RU,v}$ to the RIS-UE NN and output $\hat{\mathbf{H}}_{r,v}$;
 - 6: Obtain the predicted composite channel by $\hat{\mathbf{H}}_v = \hat{\mathbf{H}}_{r,v} \Theta_v \hat{\mathbf{G}}_v$;
 - 7: Calculate MSE by (7) and update parameters by SGD;
 - 8: **end for**
 - 9: **end for**
-

loss function in terms of the mean squared error (MSE) are defined as

$$\text{MSE} = \frac{1}{V} \sum_{v=1}^V \left\| \hat{\mathbf{H}}_v - \bar{\mathbf{H}}_v \right\|_F^2, \quad (7)$$

where V is the size of a training batch.

It is noted that the goal of the IEN is to learn the composite $\hat{\mathbf{H}}$ from location information, not the separated $\hat{\mathbf{G}}$ and $\hat{\mathbf{H}}_r$. The detailed training process is shown in Algorithm 1. After completing the training process, the predicted $\hat{\mathbf{H}}$ can thus be used to obtain the predicted achievable rate \hat{R} , denoted as

$$\hat{R} = \log_2 \det(\mathbf{I}_K + \frac{1}{\sigma^2} \hat{\mathbf{H}} \mathbf{Q} \hat{\mathbf{H}}^H). \quad (8)$$

B. The DRL-Based Algorithm

Inspired by the data-driven DRL approaches [5], [6], we proposed a novel DRL-based algorithm with the trained IEN, in which the DDPG network is invoked for solving the complex optimization problem (6).

Before proceeding, let's introduce some elements that the DDPG requires.

- *State*: a collection of observations that characterize the environment. The state $s^t \in S$ denotes the observation at the time step t , where S is the state space. In this letter, state s^t is

determined by the RIS reflection coefficient vector $\boldsymbol{\theta}^t$ at time step t , the transmit covariance matrix \mathbf{Q}^t at the time step t , the achievable rate R^t at time step t , and the location of the BS, the RIS, and the UE. To accommodate the NN, we convert the complex $\boldsymbol{\theta}^t$ and \mathbf{Q}^t into the vectors containing their real and imaginary parts. Therefore, s^t can be represented as $(\text{vec}(\text{Re}\{\mathbf{Q}^t\}, \text{Im}\{\mathbf{Q}^t\}), \text{Re}\{\boldsymbol{\theta}^t\}, \text{Im}\{\boldsymbol{\theta}^t\}, R^t, \text{loc}_{BS}, \text{loc}_{RIS}, \text{loc}_{UE})$.

- *Action*: a set of choices. The agent takes an action step by step during the learning process. Once the agent takes an action $a^t \in A$ at time step t , the state of the environment will transit from the current state s^t to the next state s^{t+1} . As a result, a reward r^t will be fed back to the agent. In this letter, the agent is the RIS controller, and action a^t is determined by \mathbf{Q}^{t+1} and $\boldsymbol{\theta}^{t+1}$. We also adopt the real-valued form of a^t as $(\text{vec}(\text{Re}\{\mathbf{Q}^{t+1}\}, \text{Im}\{\mathbf{Q}^{t+1}\}), \text{Re}\{\boldsymbol{\theta}^{t+1}\}, \text{Im}\{\boldsymbol{\theta}^{t+1}\})$, which needs to be scaled to satisfy the constraints in (6).
- *Q value and reward*: The reward r^t measures immediate return from action a^t given state s^t , whereas the Q value function measures potential future rewards which the agent may get from taking action a at the state s . In this letter, the reward r^t is the achievable rate R^t at time step t .
- *Experience*: defined as $(s^t, a^t, r^{t+1}, s^{t+1})$.

The purpose of DDPG is to maximize the output Q value. As shown in Fig. 2, the DDPG contains two basic parts, an actor network and a critic network. The actor network $\mu(s; \pi_\mu)$ takes the state as input and outputs the continuous action, which is in turn input to the critic network together with the state. The critics network $Q(s, a; \pi_Q)$ tries to approach the optimal Q value function. π_μ and π_Q are the parameters of the actor network and the critic network, respectively. In addition, two target networks copied from the actor network and the critic network, $\mu'(s; \pi_{\mu'})$ and $Q'(s, a; \pi_{Q'})$, are created for the better convergence of Q value. These two target networks have the same structure with the original two, but with different parameters. To measure the difference between the critic network's predicted value and the actual target value, a loss function is defined as

$$\text{Loss}(\pi_Q) = \frac{1}{V} \sum_{v=1}^V (y_v - Q(s^t, a^t; \pi_Q))^2, \quad (9)$$

where y_v is the actual target value of the v -th sample, defined as

$$y = r^{t+1} + \tau \max_{a'} Q(s^{t+1}, a'; \pi_{Q'}), \quad (10)$$

where $\tau \in (0, 1]$ is the discount rate, a' denotes the action output by the target actor network at the time step t .

Therefore, the critic network and the actor network can be updated by the SGD, expressed as

$$\pi_Q^{t+1} = \pi_Q^t - \lambda_Q \nabla_{\pi_Q} Loss(\pi_Q), \quad (11)$$

$$\pi_\mu^{t+1} = \pi_\mu^t - \lambda_\mu \nabla_{\pi_\mu} Q'(s^t, a; \pi_{Q'}) \nabla_{\pi_\mu} \mu(s^t; \pi_\mu), \quad (12)$$

where λ_Q and λ_μ are the corresponding learning rates, respectively. Moreover, the updates on the target actor network and the target critic network are given as

$$\pi_{\mu'} = \rho_\mu \pi_\mu + (1 - \rho_\mu) \pi_{\mu'}, \quad (13)$$

$$\pi_{Q'} = \rho_Q \pi_Q + (1 - \rho_Q) \pi_{Q'}, \quad (14)$$

respectively, where ρ_μ and ρ_Q are the corresponding learning rates for updating the target actor network and the target critic network.

However, the long convergence time of DDPG requires a huge number of transmission time slots and thus limits its application in practical communication systems. Unlike other DRL-based algorithms, the proposed algorithm employs the IEN to replace the interaction between the RIS and the actual environment. The detailed steps of the proposed location-aware DRL-based algorithm are shown in Algorithm 2.

IV. SIMULATION RESULTS

This section provides simulation results to verify the effectiveness of the proposed algorithm. We assume that BS is located at (20, 0, 10) m; RIS is deployed at (0, 30, 20) m; UE is distributed in a circular movement area with a radius of 5 m and the central location of (10, 50, 0) m. Both loc_{BS} and loc_{RIS} are assumed to be unchanged in the training processes of the IEN and the proposed DRL-based algorithm. Besides, loc_{UE} varies in the former and remains constant in the latter. For the sake of brevity, the BS-RIS channel is considered as the LoS propagation. We assume that there are two scatterers distributed between the RIS and the UE. The locations of the two scatterers are (5, 40, 10) m and (5, 45, 5) m, respectively. The distance-dependent

Algorithm 2 The Location-Aware DRL-based Algorithm

Input: $loc_{BS}, loc_{RIS}, loc_{UE}$.

Output: The optimal action $a = \{\mathbf{Q}, \boldsymbol{\theta}\}$ and the maximum achievable rate R of this algorithm.

- 1: **for** episode $j = 1, 2, \dots, J$ **do**
 - 2: Randomly generate the initial $(\mathbf{Q}^0, \boldsymbol{\theta}^0)$;
 - 3: Input $\{loc_{BS}, loc_{RIS}, loc_{UE}, \boldsymbol{\theta}^0\}$ to the trained IEN to obtain $\hat{\mathbf{H}}^0$, and then calculate \hat{R}^0 by (8) with \mathbf{Q}^0 ;
 - 4: Initialize state $s^0 = (\text{vec}(\text{Re}\{\mathbf{Q}^0\}, \text{Im}\{\mathbf{Q}^0\}), \text{Re}\{\boldsymbol{\theta}^0\}, \text{Im}\{\boldsymbol{\theta}^0\}, \hat{R}^0, loc_{BS}, loc_{RIS}, loc_{UE})$;
 - 5: Initialize a random process $\zeta \sim \mathcal{CN}$;
 - 6: **for** $t = 1, 2, \dots, T$ **do**
 - 7: Update action $a^t = \mu(s^t; \pi_\mu) + \zeta$;
 - 8: Obtain \hat{R}^t as step 2 by replacing the corresponding inputs of the IEN, and the next state s^{t+1} will be got accordingly.
 - 9: Store the experience $(s^t, a^t, r^{t+1}, s^{t+1})$ into the experience replay \mathcal{B} ;
 - 10: Sample V experiences $(s_v, a_v, r_{v+1}, s_{v+1})$ from \mathcal{B} ;
 - 11: Calculate the target Q value according to (10);
 - 12: Update the critic network $Q(s, a; \pi_Q)$ by (11);
 - 13: Update the actor network $\mu(s; \pi_\mu)$ by (12);
 - 14: Update the target actor network and the target critic network by (13) and (14).
 - 15: **end for**
 - 16: **end for**
-

path loss of each link is modeled by $PL = C_0 d^{-\alpha}$, where $C_0 = -20$ dB denotes the path loss at the reference distance of 1 m, while d and α denote the transmission distance and the path loss exponent, respectively. The path loss exponents of the BS-RIS and RIS-UE links are set to $\alpha_B = 2$ and $\alpha_R = 2.8$, respectively. Moreover, we set the transmit power at the BS as $p = 20$ dBm, while the average noise power at the UE is set to $\sigma^2 = -80$ dBm. Both actor and critic networks have four layers: one input layer, two hidden layers, and one output layer, respectively. The input layer of the actor network has $2(2M^2 + 2N) + 10$ neurons, which is

changed to $4(2M^2 + 2N) + 10$ in the critic network. The two hidden layers of both networks contain 500 and 300 neurons, respectively. The output layers of both networks have $2M^2 + 2N$ and 1 neurons, respectively. Furthermore, we exploit tanh as the activation function, which has a bounded range of outputs facilitating subsequent scaling on RIS phase shifts and the transmit covariance matrix. The learning rates, λ_Q , λ_μ , ρ_μ , and ρ_Q are set as 0.001. The buff size for experience replay \mathcal{B} and the total episodes J are 10000 and 1000, respectively. The batch size V is 16, while the discount rate $\tau = 0.99$.

Fig. 3(a) shows the MSE of the IEN's output versus the paths number of \mathbf{H}_r . It can be seen from Fig. 3(a) that the proposed IEN exhibits similar MSE performance for a given path number with the growing number of RIS elements, which means that the network can adapt to diverse RIS sizes. When the number of RIS elements remains fixed, the fitting performance of the network deteriorates as the number of paths increases, thus necessitating the network to learn more parameters and intrinsic features. Nevertheless, due to the limited number of scatterers in the mmWave band, the proposed network is capable of imitating the actual environment accurately.

Fig. 3(b) demonstrates the achievable rate performance versus the number of RIS elements, where we have $M = 4$ and $K = 4$. For the sake of illustration, we consider three comparison schemes, the AO algorithm (scheme 1) [8], the CSI-based DRL algorithm (scheme 2) [5], and the location-based DRL algorithm interacting with the actual environment (scheme 3). Observing from Fig. 3(b), all of the DRL-based i.e., scheme 2, 3, and the proposed algorithm can achieve comparable achievable rate to the traditional scheme 1. Besides, scheme 3 and the proposed algorithm employing location information perform the same but worse than the scheme 2 utilizing accurate CSI. The rate performance with the location error of UE is also investigated in Fig. 3(b). We defined the location error as follows, which is caused by the inaccurate positioning method and error during location information transmission.

$$\eta = \mathbb{E}[\|\mathbf{u} - \hat{\mathbf{u}}\|] / \mathbb{E}[\|\mathbf{u}\|], \quad (15)$$

where $\mathbf{u} = (x_{UE}, y_{UE}, z_{UE})$ is the position vector consisting of the 3D coordinates of the UE, and $\hat{\mathbf{u}}$ is the corresponding biased position vector. As we can see, the achievable rate of the proposed method decreases as η increases. As N increases, the achievable rate gap between the erroneous and perfect cases will gradually widen, which may be caused by the fact that the

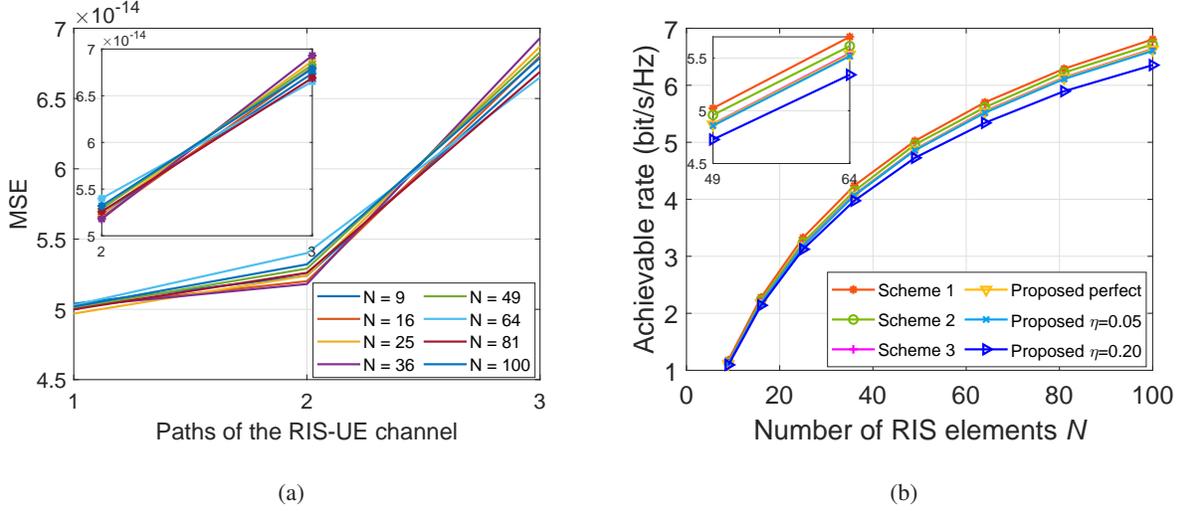


Fig. 3. (a) The MSE of the IEN's output versus the paths number of \mathbf{H}_r ; (b) The achievable rate versus the number of RIS elements, where we have $M = 4$, $K = 4$;

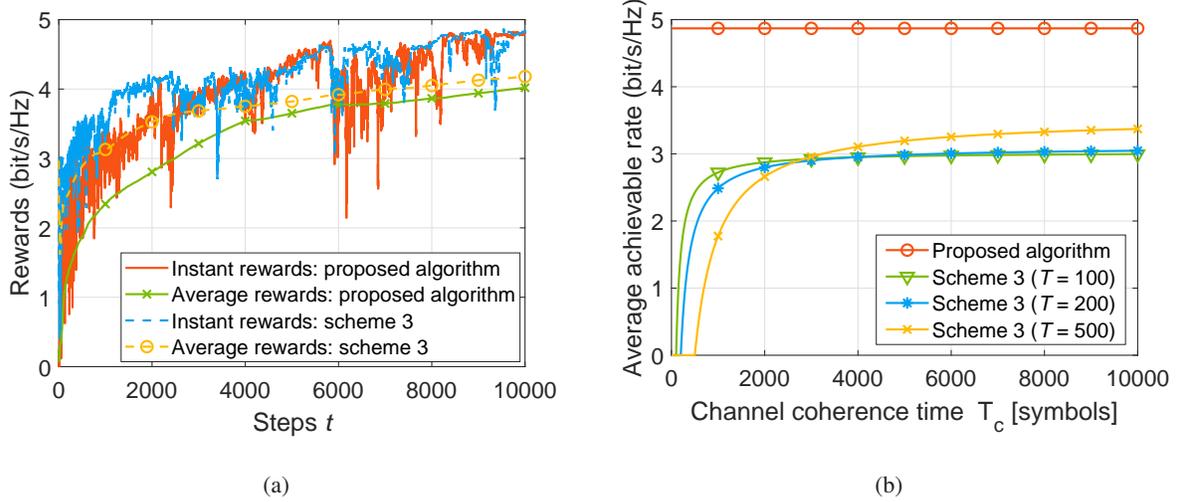


Fig. 4. (a) Rewards versus the steps, where we have $M = 4$, $N = 49$, $K = 4$. (b) The average achievable rate versus the channel coherence time T_c , where we have $M = 4$, $N = 49$, $K = 4$.

inaccurate location information can affect more links between the RIS elements and UE as N grows.

In Fig. 4(a), we evaluate the instant and average rewards between the proposed algorithm and scheme 3. The average reward $r_a(t)$ is defined as

$$r_a(t) = \frac{\sum_{i=1}^t r^i}{t}, \quad (16)$$

where r^t is the instant reward at the t -th time step. We note that these two algorithms converge with a similar trend as the training time increases, and scheme 3 earns a higher reward than

the proposed algorithm due to the mismatch of the imitation and actual environment. However, scheme 3 requires many time slots for interacting with the actual environment, whereas the proposed algorithm saves the overhead by interacting with the IEN.

Finally, Fig. 4(b) compares the average achievable rate of the proposed algorithm and scheme 3 versus the channel coherence time T_c . The average achievable rate is defined as

$$R_a = \max\left(0, \frac{T_c - T}{T_c}\right)R, \quad (17)$$

where T denotes the time slots of signal transmission for interaction, R denotes the achievable rate. Fig. 4(b) shows that the proposed algorithm results in a higher average achievable rate than scheme 3. Since the proposed algorithm interacts with the IEN and saves more time slots for data transmission, specifically, we have $T = 0$ for the proposed algorithm, implying that R_a remains constant with a given value of training time steps. On the contrary, since scheme 3 requires the signal transmission to interact with the actual environment, we have $T > 0$, which implies that R_a of scheme 3 will gradually increase for a growing value of coherence time.

V. CONCLUSION

In this letter, a novel DRL-based algorithm was proposed for RIS-aided mmWave MIMO wireless communication systems. In contrast to existing DRL-based algorithms that rely on the perfect CSI, the proposed algorithm employed the readily available UE's location information for circumventing the channel estimation. Furthermore, a network was designed to construct the map between the location information and the composite channel. This network was then employed as an imitation environment of the proposed DRL framework allowing the RIS controller to interact with the UE for the feedback rewards. The IEN decreased the interaction overhead, leading to an average achievable rate enhancement. Simulation results verified the effectiveness and advantages of the proposed algorithm over the existing DRL-based algorithms.

REFERENCES

- [1] Q. Wu and R. Zhang, "Intelligent reflecting surface enhanced wireless network via joint active and passive beamforming," *IEEE Trans. Wireless Commun.*, vol. 18, no. 11, pp. 5394–5409, Aug. 2019.
- [2] J. An and L. Gan, "The low-complexity design and optimal training overhead for IRS-assisted MISO systems," *IEEE Wireless Commun. Lett.*, vol. 10, no. 8, pp. 1820–1824, Aug. 2021.
- [3] H. Guo, Y.-C. Liang, J. Chen, and E. G. Larsson, "Weighted sum-rate maximization for reconfigurable intelligent surface aided wireless networks," *IEEE Trans. Wireless Commun.*, vol. 19, no. 5, pp. 3064–3076, Feb. 2020.

- [4] W. Xu, L. Gan, and C. Huang, "A robust deep learning-based beamforming design for RIS-assisted multiuser MISO communications with practical constraints," *IEEE Trans. Cogn. Commun. Netw.*, early access, 2021, doi:10.1109/TCCN.2021.3128605.
- [5] C. Huang, R. Mo, and C. Yuen, "Reconfigurable intelligent surface assisted multiuser MISO systems exploiting deep reinforcement learning," *IEEE J. Sel. Areas Commun.*, vol. 38, no. 8, pp. 1839–1850, Aug. 2020.
- [6] H. Yang, Z. Xiong, J. Zhao, D. Niyato, L. Xiao, and Q. Wu, "Deep reinforcement learning-based intelligent reflecting surface for secure wireless communications," *IEEE Trans. Wireless Commun.*, vol. 20, no. 1, pp. 375–388, Jan. 2021.
- [7] P. Wang, J. Fang, H. Duan, and H. Li, "Compressed channel estimation for intelligent reflecting surface-assisted millimeter wave systems," *IEEE Signal Process. Lett.*, vol. 27, pp. 905–909, May 2020.
- [8] S. Zhang and R. Zhang, "Capacity characterization for intelligent reflecting surface aided MIMO communication," *IEEE J. Sel. Areas Commun.*, vol. 38, no. 8, pp. 1823–1838, Aug. 2020.
- [9] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," *arXiv e-prints*, p. arXiv:1509.02971, Sep. 2015.
- [10] J. An, C. Xu, L. Gan, and L. Hanzo, "Low-complexity channel estimation and passive beamforming for RIS-assisted MIMO systems relying on discrete phase shifts," *IEEE Trans. Commun.*, vol. 70, no. 2, pp. 1245–1260, 2022.