

# Deep Reinforcement Learning for Intelligent Reflecting Surfaces: Towards Standalone Operation

Abdelrahman Taha\*, Yu Zhang\*, Faris B. Mismar†, and Ahmed Alkhateeb\*

\*Arizona State University, Tempe, AZ, USA, Emails: {a.taha, y.zhang, aalkhateeb}@asu.edu

†The University of Texas at Austin, Austin, TX, USA, Email: faris.mismar@utexas.edu

**Abstract**—The promising coverage and spectral efficiency gains of intelligent reflecting surfaces (IRSs) are attracting increasing interest. In order to realize these surfaces in practice, however, several challenges need to be addressed. One of these main challenges is how to configure the reflecting coefficients on these passive surfaces without requiring massive channel estimation or beam training overhead. Earlier work suggested leveraging supervised learning tools to design the IRS reflection matrices. While this approach has the potential of reducing the beam training overhead, it requires collecting large datasets for training the neural network models. In this paper, we propose a novel deep reinforcement learning framework for predicting the IRS reflection matrices with minimal training overhead. Simulation results show that the proposed online learning framework can converge to the optimal rate that assumes perfect channel knowledge. This represents an important step towards realizing a *standalone IRS operation*, where the surface configures itself without any control from the infrastructure.

**Index Terms**—reconfigurable intelligent surface, large intelligent surface, intelligent reflecting surface, smart reflect-array, beamforming, deep reinforcement learning

## I. INTRODUCTION

The increasing demand on data rates from the massive number of devices motivates the need to develop novel system architectures that are both energy and spectrally efficient. For the past years, state of the art research has focused on leveraging large-scale MIMO systems, such as massive and millimeter wave (mmWave) MIMO at the base stations (BSs) and mobile users. To further improve the coverage and the energy efficiency of these systems, intelligent reflecting surfaces (IRSs) have been recently proposed and attracted massive interest [1]–[5]. IRSs consist of a huge number of passive reflecting elements whose function is to reflect the incident signal *intelligently* into the desired directions, by means of software-controllable phase shifts. Since the IRS reflection beamforming design requires the perfect/imperfect channel knowledge, the channel estimation is a crucial aspect for the IRS interaction design problem. The massive number of passive IRS elements, however, impose a main challenge on acquiring the channel estimates; traditional channel estimation solutions will lead to either huge training overhead or prohibitive hardware complexity for the IRS architectures [5]. Given an end goal of achieving harmonic co-existence between all the heterogeneous wireless systems, setting an objective of developing *fully-standalone* IRS architectures seems as the next step forward for reaching that end goal.

Prior work focused on proposing solutions for both the channel estimation and the reflection beamforming design problems [5]–[8]. The authors in [5] proposed the first solution to the channel/beam training overhead challenge leveraging tools from both compressive sensing and supervised deep learning. The promising gains of these solutions motivated more research in these directions. For example, in [7], a supervised deep learning framework is used for channel estimation by mapping the received pilots to the direct and the cascaded channels. In [8], an IRS channel estimation scheme based on a minimum variance unbiased estimator is proposed. The solutions in [5], [7], [8], however, either considered supervised deep learning which requires large dataset collection phase before training, or assumed that the IRS is assisted/controlled by another base station/access point, not operating on its own.

This work presents a *novel* application of deep reinforcement learning in predicting the reflection coefficients of the IRS surfaces without requiring any prior training overhead. The main contributions of this paper can be summarized as follows.

- A novel deep reinforcement learning (DRL) based solution is proposed for the IRS interaction design, where the IRS learns how to reflect the incident signals in the best possible way by adjusting its reflection matrix. This solution eliminates the need for collecting large training dataset, hence requires almost no training overhead.
- The proposed framework is directed more towards *standalone IRS operation*, where the IRS architecture is not controlled/assisted by any base station, but rather operating on its own while interacting with the environment, and without any initial training phase requirement.

Simulation results based on accurate 3D ray-tracing datasets show that the achievable rates of the proposed DRL based solution can converge close to the upper bound with an added value of almost no training overhead, as opposed to supervised learning based solutions.

**Notation:**  $\mathbf{A}$  is a matrix,  $\mathbf{a}$  is a vector,  $a$  is a scalar, and  $\mathcal{A}$  is a set of vectors.  $\text{diag}(\mathbf{a})$  is a diagonal matrix with entries of  $\mathbf{a}$  on its diagonal.  $|\mathbf{A}|$  is the determinant of  $\mathbf{A}$ ,  $\mathbf{A}^T$  is its transpose,  $[\mathbf{A}]_{r,:}$  is the  $r^{\text{th}}$  row of  $\mathbf{A}$ , and  $\text{vec}(\mathbf{A})$  is a vector whose elements are the stacked columns of  $\mathbf{A}$ .  $\mathbf{I}$  is the identity matrix.  $\mathbf{A} \odot \mathbf{B}$  is the Hadamard product of  $\mathbf{A}$  and  $\mathbf{B}$ .  $\mathcal{N}(\mathbf{m}, \mathbf{R})$  is a complex Gaussian random vector with mean  $\mathbf{m}$  and covariance  $\mathbf{R}$ .  $\mathbb{E}[\cdot]$  is for expectation.

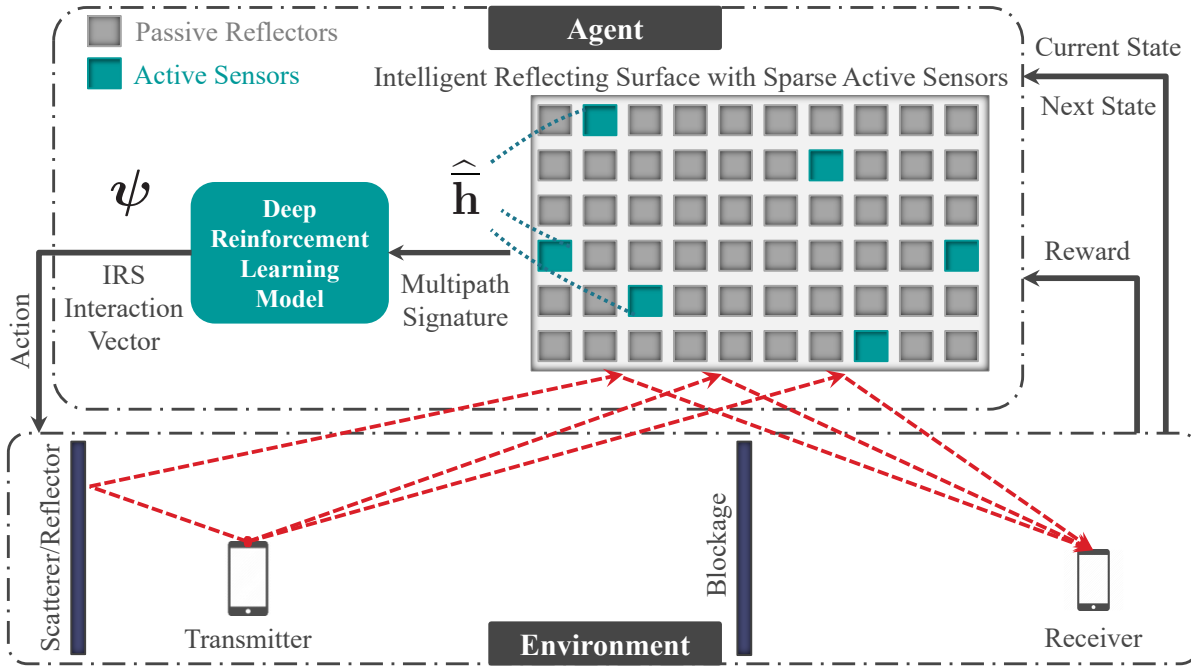


Fig. 1. The transmitter-receiver communication is assisted by an intelligent reflecting surface (IRS). The IRS is interacting with the incident signal through an interaction vector  $\psi$ . Active channel sensors are randomly distributed over the IRS. These active elements have two modes of operation (i) a channel sensing mode where it is connected to the baseband to estimate the channels and (ii) a reflection mode where it just reflects the incident signal by applying a phase shift. The rest of the IRS elements are passive reflectors. The environment is represented by the various scatterers, user locations, etc ... The IRS acts as a reinforcement learning agent by acquiring a state and a reward from the environment and exerting an action back on the environment.

## II. SYSTEM AND CHANNEL MODELS

### A. System Model

Consider an OFDM-based system of  $K$  subcarriers where a single-antenna transmitter is communicating with a single-antenna receiver due to the assistance of an  $M$ -elements intelligent reflecting surface (IRS), as in Fig. 1. Let  $\mathbf{h}_{T,k}, \mathbf{h}_{R,k} \in \mathbb{C}^{M \times 1}$  denote the channels from the transmitter/receiver to the IRS at the  $k^{\text{th}}$  subcarrier.  $s_k$  is the transmit signal, where  $\mathbb{E}[|s_k|^2] = \frac{P_T}{K}$ .  $P_T$  is the total transmit power.  $\Psi$  denotes the IRS interaction diagonal matrix.  $n_k \sim \mathcal{N}_{\mathbb{C}}(0, \sigma_n^2)$  is the receive noise. The receive signal at the receiver can be expressed as

$$y_k = \mathbf{h}_{R,k}^T \Psi \mathbf{h}_{T,k} s_k + n_k, \quad (1)$$

$$\stackrel{(a)}{=} (\mathbf{h}_{R,k} \odot \mathbf{h}_{T,k})^T \psi s_k + n_k, \quad (2)$$

where  $\psi$  is the IRS interaction vector, such that  $\Psi = \text{diag}(\psi)$ . Assume an IRS architecture of RF phase shifters, every interaction factor can be represented as  $[\psi]_m = e^{j\phi_m}$ , hence the choice of an interaction vector is constrained to a predefined codebook  $\mathcal{P}$ . Adopting the IRS architecture proposed in [5] and illustrated in Fig. 1, *active* elements are randomly distributed over the IRS. The *sampled* channel vector from the transmitter/receiver to the IRS active elements,  $\bar{\mathbf{h}}_{T,k}, \bar{\mathbf{h}}_{R,k} \in \mathbb{C}^{M \times 1}$ , can be expressed as  $\bar{\mathbf{h}}_{T,k} = \mathbf{G}_{\text{IRS}} \mathbf{h}_{T,k}$  and  $\bar{\mathbf{h}}_{R,k} = \mathbf{G}_{\text{IRS}} \mathbf{h}_{R,k}$ , where  $\mathbf{G}_{\text{IRS}}$  is an  $M \times M$  selection matrix that selects the entries corresponding to the active IRS elements. Finally, the overall IRS *sampled* channel vector can be expressed as  $\bar{\mathbf{h}}_k = \bar{\mathbf{h}}_{T,k} \odot \bar{\mathbf{h}}_{R,k}$ .

### B. Channel Model

A wideband gemoetric channel model is adopted [10]. Consider a transmitter-IRS channel,  $\mathbf{h}_{T,k}$ , (and similarly for the IRS-receiver channel) consisting of  $L$  clusters. Each cluster contributes with one ray from the transmitter to the IRS. The ray parameters are: azimuth/elevation angles of arrival,  $\theta_\ell, \phi_\ell \in [0, 2\pi)$ ; complex coefficient  $\alpha_\ell \in \mathbb{C}$ ; time delay  $\tau_\ell \in \mathbb{R}$ . The transmitter-IRS path loss is denoted by  $\rho_T$ . The pulse shaping function, with  $T_S$ -spaced signaling, is defined as  $p(\tau)$  at  $\tau$  seconds. The frequency domain channel vector,  $\mathbf{h}_{T,k}$ , can then be defined as

$$\mathbf{h}_{T,k} = \sqrt{\frac{M}{\rho_T}} \sum_{d=0}^{D-1} \sum_{\ell=1}^L \alpha_\ell \mathbf{a}(\theta_\ell, \phi_\ell) p(dT_S - \tau_\ell) e^{-j\frac{2\pi k}{K}d}, \quad (3)$$

where  $\mathbf{a}(\theta_\ell, \phi_\ell) \in \mathbb{C}^{M \times 1}$  is the IRS array response vector. Assume a block-fading channel model, where  $\mathbf{h}_{T,k}$  and  $\mathbf{h}_{R,k}$  are assumed to stay constant over the channel coherence time.

## III. PROBLEM FORMULATION

Given the objective of maximizing the achievable rate at the receiver, our problem is then to find the optimal interaction vector,  $\psi^*$ , that solves

$$\psi^* = \arg \max_{\psi \in \mathcal{P}} \sum_{k=1}^K \log_2 \left( 1 + \text{SNR} \left| (\mathbf{h}_{T,k} \odot \mathbf{h}_{R,k})^T \psi \right|^2 \right), \quad (4)$$

to achieve the optimal rate  $R^*$  defined as

$$R^* = \frac{1}{K} \sum_{k=1}^K \log_2 \left( 1 + \text{SNR} \left| (\mathbf{h}_{T,k} \odot \mathbf{h}_{R,k})^T \boldsymbol{\psi}^* \right|^2 \right). \quad (5)$$

Unfortunately, there is no closed form solution for the optimization problem in (4) due to the quantized codebook constraint and the use of one interaction vector  $\boldsymbol{\psi}$  fixed over all subcarriers. Accordingly, finding the optimal interaction vector for the IRS,  $\boldsymbol{\psi}^*$ , requires an exhaustive search over the codebook  $\mathcal{P}$ . This search, however, leads either to prohibitive training overhead, hardware complexity, or power consumption, as detailed in [5]. Our objective is then to find an efficient solution for the IRS systems that approaches the optimal rate in (5) with **almost no training overhead** and with an **energy-efficient hardware**. In the next section, we propose a *novel* application of deep reinforcement learning in the interaction design problem of intelligent reflecting surfaces. This solution actually eliminates the need for collecting large training datasets as opposed to the supervised learning solution proposed in [5]. The supervised learning solution, however, approaches the optimal rate with fewer iterations.

#### IV. DEEP REINFORCEMENT LEARNING BASED IRS INTERACTION DESIGN

##### A. Key Idea

From (4), the optimal interaction vector is a function of the channels between the two communication ends and the IRS. To avoid the prohibitive overhead of estimating the full IRS channels, the optimal interaction vector choice can be mapped to the surrounding *environment*, which the full IRS channels inherently describe. Modeling the various elements of the environment, mathematically, is notoriously complicated. In contrast, leveraging an awareness of the environment using a multipath signature [10] can be sufficient. In such case, deep reinforcement learning models can be adopted to learn the mapping function from multipath signatures to the optimal interaction vectors as illustrated in Fig. 1. The IRS active elements play a crucial role in capturing one form of multipath signatures: the *sampled* channels,  $\bar{\mathbf{h}}_{T,k}, \bar{\mathbf{h}}_{R,k}$ . Fortunately, estimating the sampled channel vectors can be accomplished with a few pilot signals; i.e., negligible training overhead. This solution also involves energy-efficient low-complexity hardware architectures (few sparse active IRS elements) [5].

##### B. Proposed Solution

The proposed deep reinforcement learning (DRL) based IRS interaction design approach operates in two parts: (I) the agent interaction and (II) the agent learning, as in Algorithm 1. The IRS interchanges between these two parts continuously.

###### PART I: Agent Interaction

The IRS interaction with the *environment* can be outlined as follows: the IRS observes the current *state*,  $s$ , of the environment and takes an *action*,  $a$ , predicated upon the observed state. The IRS then receives a *reward*,  $r$ , for the action taken and a new *state* observation,  $s'$ , from the environment. Once

the experience is acquired,  $\langle s, a, r, s' \rangle$ , the IRS trains the DRL model using current and past experiences, in the second part.

Let the term “experience” indicates the information captured in one learning episode, and define the concatenated *sampled* channel vector as

$$\bar{\mathbf{h}} = \text{vec} \left( [\bar{\mathbf{h}}_1, \bar{\mathbf{h}}_2, \dots, \bar{\mathbf{h}}_K] \right). \quad (6)$$

Assume that the one learning episode occurs every coherence block and let  $T$  be the maximum number of episodes,  $\bar{\mathbf{h}}(t)$  denotes the concatenated *sampled* channel vector at the  $t^{\text{th}}$  episode, where  $t = 1, \dots, T$ . Part I steps are summarized as follows.

**1. Sampled channel estimation (lines 3,13):** The transmitter and receiver transmits two orthogonal uplink pilots. The IRS active elements will receive these pilots and estimate the *sampled* channel vectors to construct the multipath signature.

$$\hat{\mathbf{h}}_{T,k}(t) = \bar{\mathbf{h}}_{T,k}(t) + \mathbf{v}_k, \hat{\mathbf{h}}_{R,k}(t) = \bar{\mathbf{h}}_{R,k}(t) + \mathbf{w}_k, \quad (7)$$

$$\hat{\mathbf{h}}_k(t) = \hat{\mathbf{h}}_{T,k}(t) \odot \hat{\mathbf{h}}_{R,k}(t), \quad (8)$$

$$\hat{\mathbf{h}}(t) = \text{vec} \left( [\hat{\mathbf{h}}_1(t), \hat{\mathbf{h}}_2(t), \dots, \hat{\mathbf{h}}_K(t)] \right). \quad (9)$$

where  $\mathbf{v}_k, \mathbf{w}_k \sim \mathcal{N}_{\mathbb{C}}(\mathbf{0}, \sigma_n^2 \mathbf{I})$  are the receive noise vectors.

**2. Data transmission (lines 5-10):** The multipath signature is used to predict the interaction vector. To account for exploration (i.e., randomly sampling from the action space) besides exploitation (i.e., using prior learning experience), the factor  $\epsilon$  is introduced such that an interaction vector can be randomly chosen out of the codebook  $\mathcal{P}$  with  $\epsilon$  probability. Otherwise, the interaction vector is predicted from the current network. After that, the interaction vector chosen, reflects the transmitted data from the transmitter.

**3. Feedback reception (lines 11,12):** The IRS receives a feedback from the receiver indicating the achievable rate,  $R(t)$ , attained by using the interaction vector, which is defined as

$$R(t) = \frac{1}{K} \sum_{k=1}^K \log_2 \left( 1 + \text{SNR} \left| (\mathbf{h}_{T,k}(t) \odot \mathbf{h}_{R,k}(t))^T \boldsymbol{\psi}_a \right|^2 \right). \quad (10)$$

After that, the rate is quantized based on a threshold level, such that  $R_Q(t) = 1$  if  $R(t) > R^{\text{TH}}$ , otherwise,  $R_Q(t) = -1$ . Reward clipping is substantial for learning convergence [11].

###### PART II: Agent Learning

The IRS leverages the acquired experiences to train the DRL model. Part II steps are summarized as follows.

**1. Constructing a new experience (lines 14,15):** The new experience acquired is now stored in the experience replay buffer  $\mathcal{D}$  for training of the deep Q-network [12].

**2. Model training (lines 16-23):** The deep Q-network is now trained to minimize the prediction loss. To do so, we use the stochastic gradient descent algorithm (SGD). The training operates sequentially using minibatches from the replay buffer  $\mathcal{D}$ . It learns how to map an input state (*sampled* channel vector) to an output action (interaction vector).

---

**Algorithm 1** Deep Reinforcement Learning Based IRS Interaction Design
 

---

**Input:** Reflection beamforming codebook  $\mathcal{P}$ .

**Output:** Trained network  $Q(s, a|\theta)$ .

- 1: **Initialization:** Network  $Q(s, a|\theta)$ , replay buffer  $\mathcal{D}$ .
  - 2: **repeat**
  - 3:   IRS receives two pilots to estimate  $\widehat{\mathbf{h}}(1)$ .  $\triangleright$  **Current state**
  - 4:   **for** episode  $t = 1$  **to**  $T$  **do**  $\triangleright$  **For every episode**
  - 5:     **PART I: Agent Interaction**
  - 6:     Sample  $\xi \sim \text{Uniform}(0, 1)$
  - 7:     **if**  $\xi \leq \epsilon$  **then**  $\triangleright$  **Select action**
  - 8:       Select interaction vector,  $\psi(t) \in \mathcal{P}$  at random.
  - 9:     **else**
  - 10:       Select interaction vector,  $\psi(t) = \arg \max_{a'} Q(s, a'|\theta)$ .
  - 11:     IRS reflects using  $\psi(t)$  beam.  $\triangleright$  **Carry out action**
  - 12:     IRS receives the feedback  $R(t)$ .  $\triangleright$  **Observe reward**
  - 13:     IRS quantizes the reward,  $R_Q(t) \in \{\pm 1\}$ .
  - 14:     IRS receives two pilots to estimate  $\widehat{\mathbf{h}}(t+1)$ .  $\triangleright$  **Next state**
  - 15:     **PART II: Agent Learning**
  - 16:      $\langle s, a, r, s' \rangle \leftarrow \langle \widehat{\mathbf{h}}(t), \psi(t), R_Q(t), \widehat{\mathbf{h}}(t+1) \rangle$ .
  - 17:     Store the experience  $\langle s, a, r, s' \rangle$  in  $\mathcal{D}$ .
  - 18:     Minibatch experiences from  $\mathcal{D}$  for training.
  - 19:     feedforward  $s$  to calculate  $\widehat{\mathbf{R}}(t) \leftarrow Q(s, a|\theta) \forall a$ .
  - 20:     feedforward  $s'$  to calculate  $\Gamma \leftarrow \max_{a'} Q(s', a'|\theta)$
  - 21:     and calculate  $a^* \leftarrow \arg \max_{a'} Q(s', a'|\theta)$ .
  - 22:     Construct the target vector,  $\overline{\mathbf{R}}(t)$ :
  - 23:      $[\overline{\mathbf{R}}(t)]_{a^*} \leftarrow R_Q(t) + \gamma \Gamma$ ,
  - 24:      $[\overline{\mathbf{R}}(t)]_{a' \neq a^*} \leftarrow [\widehat{\mathbf{R}}(t)]_{a' \neq a^*}, a' \in \{1, \dots, |\mathcal{P}|\}$ .
  - 25:     Perform SGD on  $\text{MSE}(\overline{\mathbf{R}}(t), \widehat{\mathbf{R}}(t))$  to find  $\theta^*$ .
  - 26:     Update network weights  $\theta(t) \leftarrow \theta^*$ .
  - 27:     Decrease  $\epsilon$  gradually.
  - 28:      $s \leftarrow s'$ .  $\triangleright$  **Assign next state to current state**
  - 29: **until** reaching a terminal goal
- 

### C. Machine Learning Design

- **Input Representation:** the concatenated *sampled* channel vector,  $\widehat{\mathbf{h}}$ , is the input to the deep Q-network. The normalization method used is a simple per-dataset scaling [13], [14]; all samples are normalized by the maximum absolute value over the whole input data. This method preserves distance information encoded in the multipath signatures. Each complex entry of the input data is split into real and imaginary values, doubling the dimensionality of each input vector to  $2K\overline{M}$ .

- **Q-Network Architecture:** The Q-network is designed as a Multi-Layer Perceptron network of  $U$  layers. The first  $U - 1$  of them alternate between fully-connected and rectified linear unit layers and the last one (output layer) is a fully-connected layer. The  $u^{\text{th}}$  layer in the network has a stack of  $A_u$  neurons. Two deep Q-networks are used for training stability [15].

- **Training Loss Function:** Given the objective of predicting the best interaction vector, having the highest achievable rate estimate, the model is trained using a regression loss function. at the  $t^{\text{th}}$  episode, the training is guided through minimizing the loss function,  $\text{MSE}(\overline{\mathbf{R}}(t), \widehat{\mathbf{R}}(t))$ , which is the mean-squared-error between the desired and the predicted output,  $\overline{\mathbf{R}}(t)$  and  $\widehat{\mathbf{R}}(t)$ .

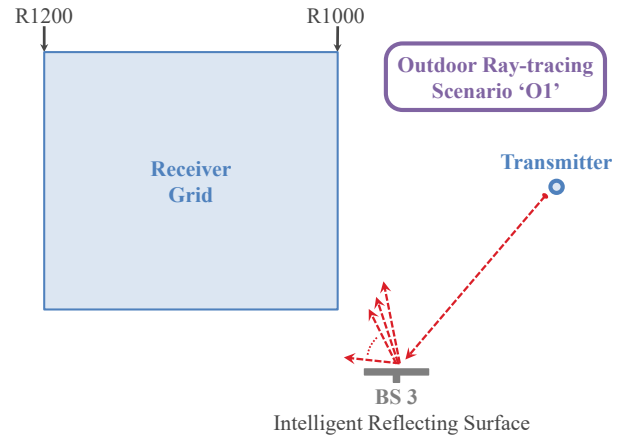


Fig. 2. The adopted ray-tracing scenario where an IRS is reflecting the signal received from one fixed transmitter to a receiver. The receiver is selected from a grid of candidate locations. This scenario is generated using Remcom Wireless InSite [16], and is available on the DeepMIMO dataset [17].

## V. SIMULATION RESULTS

In this section, we evaluate the performance of the proposed deep reinforcement learning solution.

### A. Simulation Setup

The DeepMIMO dataset in [17] is adopted to generate the channels based on the outdoor ray-tracing scenario ‘O1’. The dataset parameters are summarized in Table I. The transmitter’s position is fixed while the receiver can take any random position in a specified x-y grid, as illustrated in Fig. 2. We select BS 3 to be the IRS. For a detailed description of the simulation setup, please refer to the simulation setup in [5].

TABLE I  
THE ADOPTED DEEPMIMO DATASET PARAMETERS

DeepMIMO Dataset Parameter	Value
Frequency band	3.5 GHz
Active BSs	3
Active users (receivers)	From row R1000 to row R1200
Active user (transmitter)	row R850 column 90
Number of BS Antennas	$(M_x, M_y, M_z) = (1, 40, 10)$
Antenna spacing	$0.5\lambda$
System bandwidth	100 MHz
Number of OFDM subcarriers	512
OFDM sampling factor	1
OFDM limit	64
Number of paths	1, 15

**Deep reinforcement learning parameters:** We adopt the DRL model described in Section IV-C. States are represented by the normalized concatenated *sampled* channel of each user pair, and actions are represented by each candidate interaction vector,  $\psi \in \mathcal{P}$ . To reduce the Q-network complexity, we input the normalized *sampled* channels only at the first 64 subcarriers. The neural network architecture consists of four fully-connected layers of 4096, 16384, 16384, 4096 nodes, respectively. Given the size of the receiver x-y grid, the DRL dataset has 36200 data points. We split this dataset into two

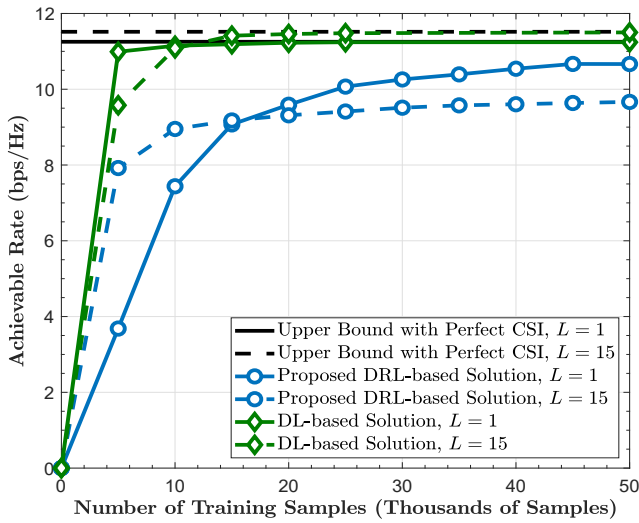


Fig. 3. The achievable rates of both the proposed deep reinforcement learning (DRL) solution and the supervised deep learning (DL) solution in [5], are compared to the upper bound, using  $\bar{M} = 4$  active elements for a 3.5GHz scenario with  $L \in \{1, 15\}$  channel path(s). The upper bound,  $R^*$  in (5), assumes perfect channel knowledge. The figure shows the potential of the proposed DRL solution in approaching the optimal rate with almost no training overhead and a small fraction of the IRS elements to be active.

sets: training and testing sets, with 70% and 30% of the points, respectively. We consider a replay buffer of 8192 samples and a batch size of 512 samples.  $\epsilon$  starts from 0.99 and decrease gradually by a factor of 0.5% every 40 training iterations till it reaches 0.1.  $\gamma = 0$ .  $R^{\text{TH}} = 8.9$  bps/Hz is set to the min-max rate of the dataset.

### B. Achievable Rates with Deep Reinforcement Learning

Fig. 3 illustrates the achievable rate of both the proposed DRL based solution and the supervised deep learning (DL) based solution in [5], using 4 active elements with  $L \in \{1, 15\}$  channel paths. Their performances are compared to the upper bound with perfect full channel knowledge, calculated according to (5). As shown, the proposed DRL solution is capable of approaching the optimal rate with more training samples than the one needed by the DL solution. **In contrast, the proposed DRL solution uses only one beam for each training episode, which constitute almost 0.3% of the beams used by the DL solution in the training phase (400 beams).** This emphasizes the efficiency of the DRL solution in operating with almost no training overhead.

Another candidate approach for refining the DRL prediction is to use the trained DRL model in predicting the most promising  $k_B$  beams. Then, these beams are used for beam training to identify the best beam that will be utilized for the rest of the coherence block. Fig. 4 illustrates the achievable rate of the proposed DRL based solution compared to the upper bound, at different values of  $k_B$  (1, 3), using 4 active elements. As demonstrated, the beam training of the promising  $k_B$  beams achieves better performance than just relying on the best network-predicted beam to reflect the incident signals. To test the effectiveness of the proposed framework, we examined

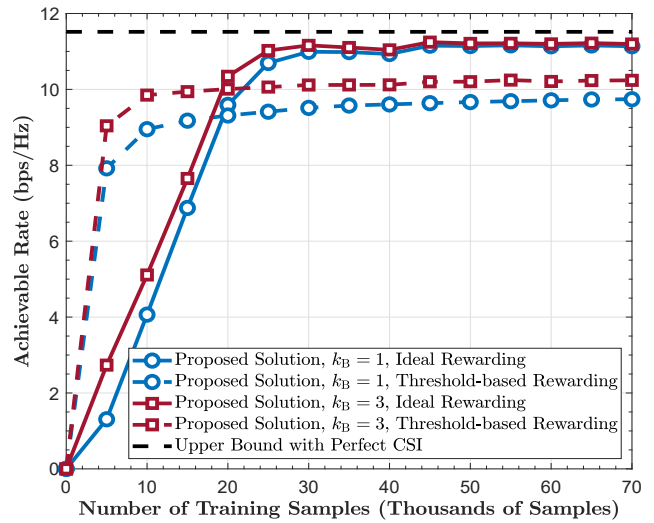


Fig. 4. The achievable rate of the proposed DRL based approach is compared to the upper bound  $R^*$ , using  $\bar{M} = 4$  active elements with  $L = 15$  channel paths. The figure illustrates the achievable rate gain when the beams selected by the deep reinforcement learning model are further refined through beam training over  $k_B$  beams.

another variant of the algorithm by updating its reward policy such that  $R_Q(t) = 1$  if  $R(t) = R^*(t)$ ; otherwise,  $R_Q(t) = -1$ , as illustrated in Fig. 4. The proposed DRL solution under this ideal rewarding assumption can converge to the optimal rate. This indicates that the small gap between the performance of the proposed solution and the upper bound can be explained by the practical assumptions of using threshold-based rewarding and operating in an environment with 15 channel paths. These results shows the gains from exploring deep reinforcement learning frameworks to develop *standalone* IRS architectures.

## VI. CONCLUSION

For an IRS-assisted wireless communication systems, we developed an efficient solution for designing the IRS interaction matrices. Given an objective of designing standalone IRS architectures, the proposed solution exploits deep reinforcement learning frameworks for the IRS to learn how to predict, on its own, the optimal interaction matrices directly from the sampled channel knowledge. This solution does not require an initial dataset collection phase as opposed to the supervised learning based solutions. Simulation results based on accurate ray-tracing channels showed that the proposed solution can converge near the optimal data rates with almost no training overhead and with few active elements.

## REFERENCES

- [1] C. Huang, S. Hu, G. C. Alexandropoulos, A. Zappone, C. Yuen, R. Zhang, M. Di Renzo, and M. Debbah, "Holographic MIMO Surfaces for 6G Wireless Networks: Opportunities, Challenges, and Trends," *arXiv preprint arXiv:1911.12296*, 2019.
- [2] E. Basar, M. Di Renzo, J. De Rosny, M. Debbah, M.-S. Alouini, and R. Zhang, "Wireless Communications Through Reconfigurable Intelligent Surfaces," *IEEE Access*, vol. 7, pp. 116 753–116 773, 2019.
- [3] Y.-C. Liang, R. Long, Q. Zhang, J. Chen, H. V. Cheng, and H. Guo, "Large Intelligent Surface/Antennas (LISA): Making Reflective Radios Smart," *arXiv preprint arXiv:1906.06578*, 2019.

- [4] X. Yuan, Y.-J. Zhang, Y. Shi, W. Yan, and H. Liu, "Reconfigurable-Intelligent-Surface Empowered 6G Wireless Communications: Challenges and Opportunities," *arXiv preprint arXiv:2001.00364*, 2020.
- [5] A. Taha, M. Alrabeiah, and A. Alkhateeb, "Enabling Large Intelligent Surfaces with Compressive Sensing and Deep Learning," *arXiv preprint arXiv:1904.10136*, Apr 2019.
- [6] C. Huang, G. C. Alexandropoulos, C. Yuen, and M. Debbah, "Indoor Signal Focusing with Deep Learning Designed Reconfigurable Intelligent Surfaces," in *2019 IEEE 20th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*. IEEE, 2019, pp. 1–5.
- [7] A. M. Elbir, A. Papazafeiropoulos, P. Kourtessis, and S. Chatzinotas, "Deep Channel Learning For Large Intelligent Surfaces Aided mm-Wave Massive MIMO Systems," *arXiv preprint arXiv:2001.11085*, 2020.
- [8] T. L. Jensen and E. De Carvalho, "An Optimal Channel Estimation Scheme for Intelligent Reflecting Surfaces Based on a Minimum Variance Unbiased Estimator," *arXiv preprint arXiv:1909.09440*, 2019.
- [9] A. Taha, M. Alrabeiah, and A. Alkhateeb, "Deep Learning for Large Intelligent Surfaces in Millimeter Wave and Massive MIMO Systems," in *Proc. of IEEE Global Communications Conference (GLOBECOM)*. IEEE, 2019, pp. 1–6.
- [10] A. Alkhateeb, S. Alex, P. Varkey, Y. Li, Q. Qu, and D. Tujkovic, "Deep Learning Coordinated Beamforming for Highly-Mobile Millimeter Wave Systems," *IEEE Access*, vol. 6, pp. 37 328–37 348, 2018.
- [11] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, "Human-level Control through Deep Reinforcement Learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [12] F. B. Mismar, B. L. Evans, and A. Alkhateeb, "Deep Reinforcement Learning for 5G Networks: Joint Beamforming, Power Control, and Interference Coordination," *IEEE Transactions on Communications*, 2019.
- [13] Y. Zhang, M. Alrabeiah, and A. Alkhateeb, "Deep Learning for Massive MIMO with 1-Bit ADCs: When More Antennas Need Fewer Pilots," *arXiv preprint arXiv:1910.06960*, 2019.
- [14] X. Li and A. Alkhateeb, "Deep Learning for Direct Hybrid Precoding in Millimeter Wave Massive MIMO Systems," in *Asilomar Conference on Signals, Systems, and Computers*, *arXiv preprint arXiv:1905.13212*, Nov. 2019.
- [15] H. Van Hasselt, A. Guez, and D. Silver, "Deep Reinforcement Learning with Double Q-learning," in *Proc. of AAAI conference on artificial intelligence*, 2016.
- [16] Remcom, "Wireless InSite," <http://www.remcom.com/wireless-insite>.
- [17] A. Alkhateeb, "DeepMIMO: A Generic Deep Learning Dataset for Millimeter Wave and Massive MIMO Applications," in *Proc. of Information Theory and Applications Workshop (ITA)*, San Diego, CA, Feb 2019, pp. 1–8. [Online]. Available: <https://www.deepmimo.net/>