

# Federated Deep Reinforcement Learning for THz-Beam Search with Limited CSI

Po-Chun Hsu, Li-Hsiang Shen, Chun-Hung Liu\*, and Kai-Ten Feng

Department of Electrical and Computer Engineering, National Yang Ming Chiao Tung University, Hsinchu, Taiwan

Department of Electrical and Computer Engineering, Mississippi State University, MS, USA\*

e-mail: {g309513013.c, ktfeng, gp3xu4vu6.cm04g}@nycu.edu.tw and chliu@ece.msstate.edu\*

**Abstract**—Terahertz (THz) communication with ultra-wide available spectrum is a promising technique that can achieve the stringent requirement of high data rate in the next-generation wireless networks, yet its severe propagation attenuation significantly hinders its implementation in practice. Finding beam directions for a large-scale antenna array to effectively overcome severe propagation attenuation of THz signals is a pressing need. This paper proposes a novel approach of federated deep reinforcement learning (FDRL) to swiftly perform THz-beam search for multiple base stations (BSs) coordinated by an edge server in a cellular network. All the BSs conduct deep deterministic policy gradient (DDPG)-based DRL to obtain THz beamforming policy with limited channel state information (CSI). They update their DDPG models with hidden information in order to mitigate inter-cell interference. We demonstrate that the cell network can achieve higher throughput as more THz CSI and hidden neurons of DDPG are adopted. We also show that FDRL with partial model update is able to nearly achieve the same performance of FDRL with full model update, which indicates an effective means to reduce communication load between the edge server and the BSs by partial model uploading. Moreover, the proposed FDRL outperforms conventional non-learning-based and existing non-FDRL benchmark optimization methods.

**Index Terms**—Terahertz, federated learning, deep reinforcement learning, beamforming, edge computing.

## I. INTRODUCTION

For the purpose of meeting the increasing ultra-high data requirement, such as virtual reality (VR), augmented reality (AR) and hologram technologies in the future sixth generation (6G) communication network, the prospect of terahertz (THz) communication is considered to be able to provide higher frequency bands ranging from 0.1 to 10 THz. However, compared to conventional millimeter wave (mmWave) at gigahertz (GHz) bands, the most different and important challenges for THz are severe power attenuation, blockages and additional molecular absorption that leads to a much shorter propagation distance. Therefore, beamforming techniques are utilized to enhance the transmit direction towards the desired receiving user equipment (UE) rather than omni-directional ones which

do not require a large-scale antenna array deployed to obtain high beam gain and spatial diversity. Although there exist potentially rich channel paths in a multi-antenna system, only a single line-of-sight (LoS) link can be utilized in the THz band in most cases. The related beamforming designs in [1], [2] were proposed for THz systems. However, when associating with enormous UEs in different cells, serious interfered beamforming issues would probably exist in short transmission distances of a THz network, and thereby it would be necessary to appropriately coordinate multiple BSs in the THz network so as to improve the overall transmission performance of the THz network. In addition, references in [3]–[5] employed full CSI of beamforming that requires highly complicate and accurate channel estimation, which is fairly difficult to be implemented for dynamic large-scale antenna arrays. Furthermore, estimating the full CSI of a large-scale antenna array is always a difficult task because wireless channels are prone to be affected by dynamic environmental variances. Accordingly, traditional optimization methods are hardly to perfectly tackle the beam search problem in wireless network deployed with large-scale antenna arrays of THz.

Recently, deep learning techniques are widely applied in the different fields in wireless communication systems. As a prospect, deep reinforcement learning (DRL) enables the agent, which may be BS or UE to adjust its wireless state and action, i.e., policy output according to the changed environment. Different from model-free reinforcement learning, the deep Q network (DQN) architecture is implemented via deep neural network (DNN) to decide the Q-value instead of the Q-table, which benefits problems with non-countable or near-continuous variables with infinite solution sets. However, when the action space is high-dimensional and continuous, it is very inefficient to apply a regular DQN to obtain the decision space. Deep deterministic policy gradient (DDPG) using a two-layered DNN as actor-critic (AC) network is known to deal with this problem. Federated edge learning (FEL) has been considerably studied for improving the training progress via learning model exchange with less information uploaded to edge server [6]. The main concept of FEL is to integrate local training models from different (mobile) clients in order to acquire a more complete global model, which can include certain hidden information in different clients. References in [7]–[9] studied how to improve the learning speed when

<sup>1</sup>The work of P.-C. Hsu, L.-H. Shen, and K.-T. Feng was supported in part funded by Ministry of Science and Technology (MoST) Grants 110-2221-E-A49-041-MY3, 111-2218-E-A49-024, STEM Project, the National Defense Science and Technology Academic Collaborative Research Project in 2022, and Higher Education Sprout Project of the National Yang Ming Chiao Tung University and Ministry of Education (MoE), Taiwan. The work of C.-H. Liu was supported in part by the U.S. National Science Foundation (NSF) under Award CNS-2006453.

minimizing computing delay and energy consumption through FEL. In [10], how to enhance the privacy of every participating client by using FEL. Although these aforementioned prior works advanced the study of FEL over wireless communication, they overlooked two main practical issues that were not considered in their system models, i.e., co-channel interference and limited CSI for large-scale antenna array. Therefore, there exists a pressing need to design interference mitigation scheme in THz beamforming by combining FEL with deep learning techniques.

The contributions of this paper are summarized as follows.

- We have conceived a federated deep reinforcement learning (FDRL) leveraging the benefits of both FEL and DRL architectures. FEL aims at model exchange from neural networks extracting hidden information of partially estimated CSI, which potentially alleviates interference from other BSs. While, AC-based DDPG is designed to search candidate THz beams to maximize the total throughput performance.
- We characterize the performance in terms of computational complexity and network throughput. We can infer that higher network throughput can be achieved with more antennas, exchanged data, and more neurons of FDRL under a compromised computational complexity of deep learning. The proposed FDRL scheme outperforms the baseline using pure deep Q-learning and conventional non-deep learning based beamforming methods.

The rest of the paper is organized as follows: Section II describes the system model and formulates the THz beamforming problem. Section III elaborates on our proposed FDRL algorithm for coordinating THz beamforming under a multi-BS network. Section IV shows simulation results, whilst conclusions are drawn in Section V.

## II. SYSTEM MODEL

In this paper, we consider a cellular network in which each UE is equipped with a single antenna and there are  $K$  base stations (BSs) operating in the THz (frequency) band, each of which equipped with  $N$  antennas. In the downlink, each BS is assumed to adopt different resource blocks to serve different UEs in its cell. Namely, no UEs in the same cell share the same resource blocks and thereby there is no intra-cell interference in the network. We also assume that the frequency reuse factor is one in the network so that the  $K$  BSs interferes each other in the downlink and UE  $k$  thus receives interference from the other  $K - 1$  BSs when it is served by BS  $k$ . All the BSs are connected with high-speed optical links to an edge server where edge computing can be conducted. A schematic diagram of the cellular network with edge computing considered in this paper is shown in Fig. 1. In the following, we will first introduce the channel model in the THz band and then specify the signal model transmitted over a THz channel.

### A. THz Channel Model

Due to THz signals' nature of extremely high frequency, transmitting them significantly suffers from two serious envi-

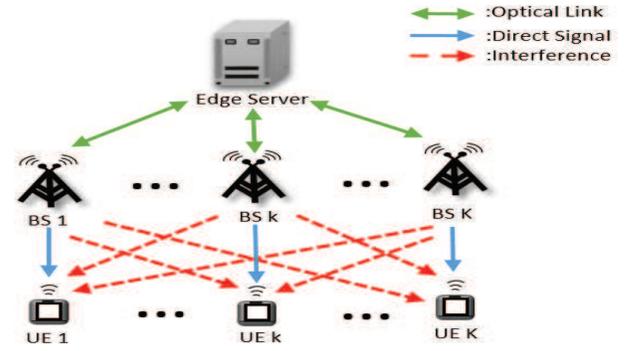


Fig. 1. A schematic diagram of the cellular network with edge computing considered in this paper. In the network, there are  $K$  base stations, each of them equipped with  $N$  transmit antennas. A downlink communication scenario is considered and BS  $k$  serves UE  $k$  equipped with a single antenna in the THz frequency band. Each of the BSs is connected to an edge server through a high-speed optical link.

ronmental impairments, i.e., severe attenuation and molecular absorption [11]–[14]. As such, THz signals undergo much higher path loss than mmWave as well as UHF signals. For a THz channel with frequency  $f$ , its channel response for transmitting a signal over distance  $d$ , denoted by complex vector  $\mathbf{h} \in \mathbb{C}^N$ , can be modeled as

$$\mathbf{h} = G \left[ 1 + \sum_{l=1}^L \Lambda_l(f) \right] a_L(f, d) \mathbf{a}_t(\theta_t), \quad (1)$$

where  $G$  is called the integrated antenna gain consisting of the transmitted and received antenna gains of the antenna array,  $L$  denotes the number of non-line-of-sight (NLoS) paths,  $\Lambda_l(f)$  is a frequency-dependent constant consisting of the reflection factor and roughness coefficient of NLoS paths affected by the reflective interfaces and material impedance. Moreover,  $a_L(f, d)$  is defined as

$$a_L(f, d) \triangleq \frac{c}{4\pi f d} e^{-\frac{1}{2}\rho(f)d}$$

in which  $c$  is the speed of light and  $\rho(f)$  is the medium absorption factor of frequency  $f$ . Furthermore, we consider a uniform linear array at each BS so that  $\mathbf{a}_t(\theta_t)$  can be defined as

$$\mathbf{a}_t(\theta_t) \triangleq \frac{1}{N} \left[ 1, e^{j\frac{2\pi}{\lambda}d_a \sin(\theta_t)}, \dots, e^{j\frac{2\pi}{\lambda}d_a(N-1) \sin(\theta_t)} \right]^T, \quad (2)$$

where  $\theta_t \in [-\frac{\pi}{2}, \frac{\pi}{2}]$  is the angle of departure,  $d_a$  is the distance between two antennas, and  $T$  denotes the transpose operation of a vector. Note that  $\mathbf{h}$  in (1) consists of LoS and NLoS components, that is,  $G a_L(f, d) \mathbf{a}_t(\theta_t)$  is the LoS component whereas the other term is the NLoS component.

### B. THz Signal Model

Let  $\mathbf{w}_k \in \mathbb{C}^N$  be the beamforming vector for BS  $k$  and  $x_k \in \mathbb{C}$  be the signal with unit power transmitted to the  $k$ th UE by BS  $k$ . Since there are  $K$  BSs in the network, we consider the worst scenario that all the BSs interferes each other when they serve their UE. As a result, the signal received by UE  $k$

can be specifically written as

$$y_k = \underbrace{\sqrt{P} \mathbf{h}_{kk}^H \mathbf{w}_k x_k}_{\text{desired signal}} + \underbrace{\sum_{j=1, j \neq i}^K \sqrt{P} \mathbf{h}_{jk}^H \mathbf{w}_j x_j}_{\text{interference signal}} + \underbrace{n_k}_{\text{noise}}, \quad (3)$$

where  $k \in \{1, \dots, K\}$ ,  $P$  is the transmit power of each BS, and superscript  $\mathcal{H}$  stands for the Hermitian operation of a complex vector,  $n_k \in \mathbb{C}$  denotes the Gaussian noise, and  $\mathbf{h}_{kk} \in \mathbb{C}^N$  and  $\mathbf{h}_{jk} \in \mathbb{C}^N$  are the channel vectors from BS  $k$  to UE  $k$  and from (interfering) BS  $j$  to UE  $k$ , respectively. Note that  $\mathbf{h}_{kk}$  and  $\mathbf{h}_{jk}$  adopt the channel model defined in (1). As such, the signal-to-noise-plus-interference ratio (SINR) received at the  $k$ th UE can be defined as

$$\Gamma_k = \frac{P |\mathbf{h}_{kk}^H \mathbf{w}_k|^2}{\sum_{j=1, j \neq i}^K P |\mathbf{h}_{jk}^H \mathbf{w}_j|^2 + \sigma_n^2}, \quad (4)$$

where  $|\cdot|$  represents the operator of absolute value and  $\sigma_n^2$  is the power of the Gaussian noise  $n_k$  for all  $k \in \{1, \dots, K\}$ . According to (4), the downlink achievable rate (spectral efficiency) of BS  $k$  can be written as

$$C_k = \log_2(1 + \Gamma_k), \quad (\text{bits/sec/Hz}) \quad (5)$$

for all  $k \in \{1, \dots, K\}$ . In the following, we will use  $C_k$  to formulate an optimization problem of beam search that is able to maximize the sum rate of all the BSs in the scenario that only limited CSI is available at each BS.

### C. Sum-Rate Optimization with Limited CSI

In this work, we aim to maximize the sum rate of the network by optimizing each beamforming vector  $\mathbf{w}_k$  with limited CSI  $g(\mathbf{h}_{kk}) = \mathbf{h}_{kk}^{lim} \in \mathbb{C}^N$  that is estimated at each BS  $k$ . The limited CSI means that only part of the total CSI at each BS can be effectively estimated due to the large size of each antenna array, and thereby the achievable SINR  $\Gamma_k$  in (4) reduces to the limited SINR  $\Gamma_k^{lim}$ , which equals  $\Gamma_k$  that adopts  $\mathbf{h}_{kk}^{lim}$  in place of  $\mathbf{h}_{kk}$ . As a result, we can formulate a optimization problem of the (mean) network throughput (sum rate) with limited CSI as follows:

$$\begin{aligned} \max_{\mathbf{w}_k} \quad & \sum_{k=1}^K C_k^{lim} \\ \text{s.t.} \quad & \text{tr}(\mathbf{w}_k \mathbf{w}_k^H) \leq 1, \quad g(\mathbf{h}_{kk}) = \mathbf{h}_{kk}^{lim}, \quad k \in \{1, \dots, K\}, \end{aligned} \quad (6)$$

where  $C_k^{lim} \triangleq \log_2(1 + \Gamma_k^{lim})$  and  $\text{tr}(\cdot)$  denotes the trace operator of a matrix. However, the optimization problem in (6) cannot be readily solved via conventional optimization methods with respect to the digital beamforming and partially-attainable CSI. Furthermore, highly computational complexity of global optimum and huge overhead of CSI exchange caused by many antenna arrays in terahertz networks make the conventional methods fairly difficult to analyze the sophisticated and unpredictable communication network. As a consequence, we design a deep-learning-based scheme by jointly leveraging DRL and federated learning architectures, which helps to reduce the amount of CSI exchange and tries to relieve the

impact of interference, to resolve the complex optimization problem.

## III. FEDERATED DEEP REINFORCEMENT LEARNING (FDRL) FOR THZ BEAM SEARCH

Since the optimization problem of the network throughput with limited CSI in (6) is not analytically tractable, we propose an FEL-based DDPG learning scheme that can iteratively coordinate to attain an optimal policy of THz beamforming. We consider that each BS conducts a DDPG to obtain a policy of THz beamforming with limited CSI when all the BSs are controlled by a single FEL server to exchange training model with hidden information so as to mitigate interference.

### A. DRL-based DDPG Network

We consider that the DRL framework contains a state set  $\mathcal{S}$ , an action set  $\mathcal{A}$ , a reward set  $\mathcal{R}$ , and an agent (BS or UE) that applies a certain action to obtain the corresponding reward while updating the current status. The action will be reinforced iteratively in order to receive better rewards in a varying environment. On account of the large state-action set existing in our THz beamforming problem, using conventional DRL approaches to tackle the problem is certainly inappropriate because it definitely suffers from huge storage of table-mapping and slow convergence. Moreover, since a THz beamforming vector is deemed to be continuous variables with substantially-high quantization levels, this thus motivates us to adopt DDPG to establish a two-layered actor-critic network to resolve the problem with continuous solutions. For the DRL-based DDPG in the THz network, we define the state, action and rewards as follows.

1) *State set*  $\mathcal{S}$ : In THz, the state set collects the statuses of each BS under current THz channels, denoted by  $\mathcal{S} = \{s_k | k = 1, 2, \dots, K\}$ , which consists of the serving CSI  $\mathbf{h}_{kk}$  linked to the  $i$ -th UE, and SINR  $\Gamma_k$  feedback from the UE  $i$ . Note that  $\mathbf{h}_{kk}$  may be partially attainable due to limited measured CSI under a large-scale THz antenna array. Therefore, state of each BS should be  $s_k = \{\mathbf{h}_{kk}^{lim}, \Gamma_k | k = 1, 2, \dots, K\}$ .

2) *Action set*  $\mathcal{A}$ : The action set represents the decision-making of THz beamforming vector defined as  $\mathcal{A} = \{\mathbf{a}_k = \mathbf{w}_k | k = 1, 2, \dots, K\}$ . Note that each BS will only determine its own action, i.e.,  $\mathbf{w}_k$  according to current input state and reward.

3) *Reward set*  $\mathcal{R}$ : We define the overall reward as  $\mathcal{R} = \{r_k | k = 1, 2, \dots, K\}$ . Since we aim at maximizing the sum throughput in (6), we consider the reward function as individual throughput of each BS, i.e.,  $r_k = C_k^{lim}$ .

As shown in Fig. 2, the DDPG architecture contains the main and target networks that individually consist of actor and critic sub-networks wherein  $\theta_k^\mu$  and  $\theta_k^Q$  denote DNN-enabled actor/critic weights in the main network, and  $\theta_k^{\mu'}$  and  $\theta_k^{Q'}$  denote the actor/critic weights in the target network. The main network determines the beamforming action of the  $i$ -th THz BS as  $a_{k,t} = \mu(s_{k,t} | \theta_k^\mu) + N_G$ , where  $\mu(s_{k,t} | \theta_k^\mu)$  is the output layer of the DNN-based actor network. To perform the exploration of the environment, the deterministic policy will

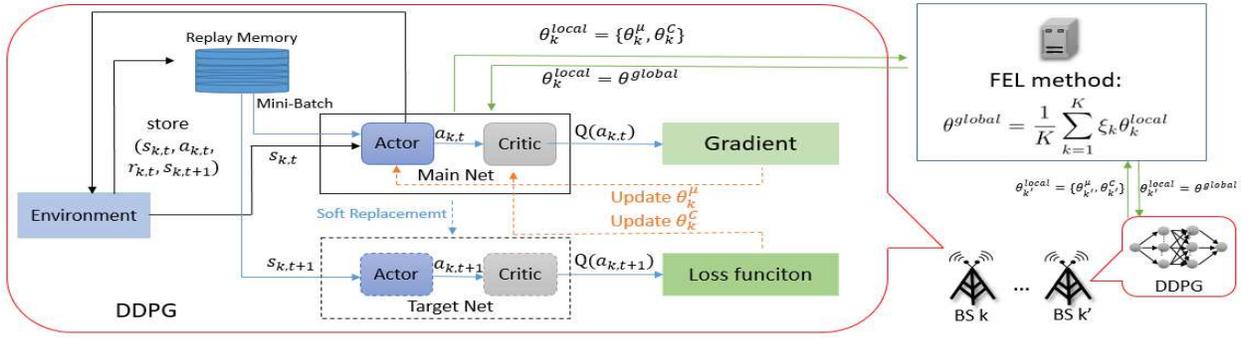


Fig. 2. The DRL-based DDPG algorithm. The main/target networks contain their actor-critic sub-networks established by DNN.

obtain the probabilistic action by adding the perturbation  $N_G$  as Gaussian noise. On the other hand, the target network input is fed by the action of actor network, which provides the Q-value outcome  $Q(s_{k,t}, a_{k,t} | \theta_k^Q)$  via hidden DNN layers at  $t$ -th epoch to evaluate the selected action, which is written as

$$Q(s_{k,t}, a_{k,t} | \theta_k^Q) = \mathbb{E} \left[ r_k + \gamma Q(s_{k,t+1}, a_{k,t+1} | \theta_k^Q) \right], \quad (7)$$

where  $\gamma$  is a discount factor and  $\mathbb{E}[\cdot]$  is the expectation for the trajectory since it is difficult to sample thousands of situations. In order to update the DDPG network, the gradient of actor network is acquired as

$$\begin{aligned} \nabla_{\theta_k^\mu} J_i &= \nabla_{\theta_k^\mu} \mathbb{E} \left[ Q(s_{k,t}, a_{k,t} | \theta_k^Q) \right] \\ &= \mathbb{E} \left[ \nabla_{a_{k,t}} Q(s_{k,t}, a_{k,t} | \theta_k^Q) \cdot \nabla_{\theta_k^\mu} \mu(s_{k,t} | \theta_k^\mu) \right], \end{aligned} \quad (8)$$

where critic loss function can be given by

$$L_k = \mathbb{E} \left[ y_k - Q(s_k, a_k | \theta_k^Q) \right]^2 \quad (9)$$

with  $y_k = r_k + \gamma Q(s_{k,t+1}, a_{k,t+1} | \theta_k^Q)$ . The target network will periodically update the network weights from the main network based on the soft update [15] for both actor-critic sub-networks which is represented by

$$\theta_k^{Q'} = \tau_a \theta_k^Q + (1 - \tau_a) \theta_k^{Q'}, \quad (10)$$

$$\theta_k^{\mu'} = \tau_c \theta_k^\mu + (1 - \tau_c) \theta_k^{\mu'}, \quad (11)$$

where  $\tau_a$  and  $\tau_c$  are constants indicating the significance of parameters in the target and main networks.

### B. Federated Edge Learning for Interference Mitigation

After DDPG learning for THz beamforming adjustment, the local DDPG training models of NN's weights will be sent from BSs to the edge server per  $T$  iteration. The edge server will then aggregates local training model in order to exchange the hidden information of interference information rather than directly upload high-overhead full CSI data. Inspired by the FEL method in [15], the model aggregation can be presented as

$$\theta^{global} = \frac{1}{K} \sum_{k=1}^K \xi_k \theta_k^{local}, \quad (12)$$

### Algorithm 1: Proposed FDRL Algorithm

- 1: Input:  $\mathbf{h}_{kk}^{lim}, \Gamma_k, \forall k$
- 2: Output:  $a_k = \mathbf{w}_k, \forall i$
- 3: Initialize:  $\theta_k^\mu, \theta_k^{\mu'}, \theta_k^Q, \theta_k^{Q'} \forall k, \theta^{global}$ , replay memory  $M_r$
- 4: **for**  $t = 1, 2, \dots, E$  **do**
- 5:   **for each** BS  $k$  **do**
- 6:     Decide the action  $a_{k,t} = \mu(s_{k,t} | \theta_k^\mu) + N_G$
- 7:     Interact with the environment and save result of  $(s_{k,t}, a_{k,t}, r_k, s_{k,t+1})$  to replay memory  $M_r$
- 8:     Off-line actor/critic model training by mini-batching data with a size of  $B$
- 9:     Soft update  $\theta_k^{\mu'}, \theta_k^{Q'}$
- 10:   **end for**
- 11:   **if**  $\text{mod}(t, T) = 0$  **then**
- 12:     FEL model aggregation:  $\theta^{global} = \frac{1}{K} \sum_{i=1}^K \xi_k \theta_k^{local}$
- 13:     Model update after aggregation:  $\theta_k^{local} = \theta^{global}$
- 14:   **end if**
- 15: **end for**

where  $\xi_k \in [0, 1]$  is a ratio indicating the importance of each training model depending on certain property of dataset in each BS and  $\sum_{k=1}^K \xi_k = 1$ . In (12),  $\theta_k^{local} = \{\theta_k^\mu, \theta_k^Q\}$  consists of neural weights of main actor/critic network. Note that in this case, we consider  $\xi_k = 1$  as equivalent importance of each beamforming model since the THz BS could provide potentially useful information of limited estimated THz CSI. After finishing the model aggregation, the edge server returns the global parameters to each BS, which is repeatedly performed until convergence. Thus, the candidate beamforming of each THz BS  $\mathbf{w}_k$  will converge to near optimum by searching for the higher DDPG reward through the iterative training. Additionally, to further address full model upload problem, we design a partial FEL mechanism by uploading partial weight elements  $\theta_{k,P}^{local}$  of original training model of  $\theta_k^{local}$ . That is, we upload the weight elements with a ratio of  $\mathcal{N}(\theta_{i,P}^{local}) / \mathcal{N}(\theta_i^{local})$ , where  $\mathcal{N}(\cdot)$  is defined as a function calculating the number of elements in the weight vectors. The concrete algorithm of proposed FDRL is proposed in Algorithm 1.

TABLE I  
PARAMETERS OF FDRL-ENABLED THZ NETWORK

Definition	Symbol	Value
Discount factor	$\gamma$	0.9
Significance of actor network	$\tau_a$	0.01
Significance of critic network	$\tau_c$	0.01
Number of epoch	$E$	300
Buffer of the memory size	$M_r$	10
Batch size of DDPG training	$B$	5
The cycle of FEL	$T$	20
Number of antennas	$N$	{8, 16, 32, 64, 128, 256}
Number of BSs	$K$	{2, 3, 6}
Operating frequency	$f$	0.3 THz
Distance between BS and UE	$d$	[10, 100] m
Medium absorption factor	$\rho(f)$	0.1
Number of NLoS paths	$L$	5
Integrated antenna gains	$G$	10 dB
Noise power	$\sigma_n^2$	-74 dBm
Transmit power	$P$	10 dBm

#### IV. SIMULATION RESULTS

In this section, we have performed simulations of proposed FDRL-enabled THz beamforming with maximization of system throughput while mitigating network interference. The THz BS and serving UEs are uniformly-randomly distributed in the radius from 10 to 100 meters. We consider  $N_{NL} = 5$  NLoS paths and THz frequency is set to be 0.3 THz. The remaining parameter setting of the THz network is listed in Table I.

In Fig. 3, the convergence of the throughput with  $K = 3$ , clients served by three BS equipped with  $N_t = 8$  antennas. Each BS have a two-layer actor-critic neural network with {100, 70} neurons. The actor network decides the beamforming vector  $\mathbf{w}_k$ , while the critic network using the Q-learning network evaluates the decision of actor network. Initially, the beamforming vector is randomly selected with the perturbation  $N_G$  with variance equal to 3 leveraging exploitation and exploration. However, it will gradually decay to 0.99 as deterministic decision. At around 100-th epoch, the client 1 tends to be stable, but the others are still searching for the potential solution from DDPG network. At around 150-th epoch, the performance of throughput is almost converged. Note that the result only shows a single run of certain channel condition in Fig. 3; however, we will conduct 100 Monte Carlo runs leveraging different channel conditions in the following comparisons.

In Fig. 4, we illustrate the throughput that is affected by the number of the actor and critic neurons. Each of the BSs is equipped with  $N_t = 128$  antennas and  $K = 4$  is considered. As the actor-critic neurons is set as (20, 20), the performance of full FEL upload achieves higher throughput than that of partial upload with around 2 and 0.5 Gbps for 10% and 50% upload of FEL parameters. However, the operational overhead, i.e., operation for training in local network and FEL server is quite higher using full upload than that of 10% upload. We can observe that the 10%-upload overhead possesses half the overhead compared to that of full FEL upload while sustaining adequately high throughput performance, i.e., the 10% exchanged hidden information is enough to alleviate induced interference. Moreover, when the number of actor-

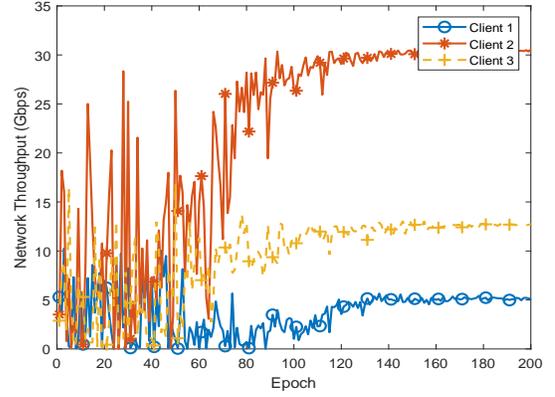
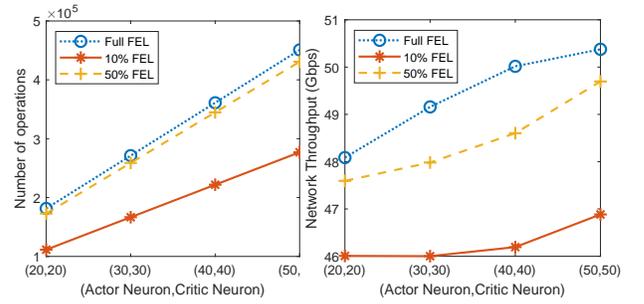


Fig. 3. The network throughput achieved by the proposed FDRL with 3 BS (i.e.,  $K = 3$ ) along the training epochs.



(a) Complexity of FDRL (b) Throughput of FDRL  
Fig. 4. The proposed FDRL compared to different numbers of actor and critic neurons with {(20, 20), (30, 30), (40, 40), (50, 50)}.

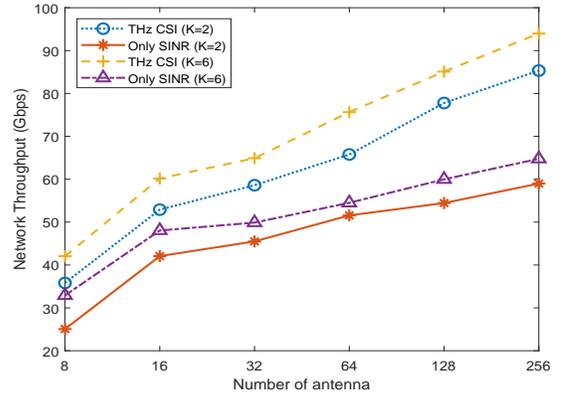


Fig. 5. Network throughput of the proposed FDRL with different number of antennas considering limited CSI and only SINR feedback with  $K = \{2, 6\}$  BS.

critic neurons is increased to (30, 30), it reaches high throughput performance but provokes higher computational overhead, which strikes a compelling tradeoff between complexity and throughput performance.

Fig. 5 demonstrates the throughput considering input states of CSI and only SINR feedback with different number of antennas and THz BS deployment with  $K = \{2, 6\}$  and  $N_t = \{8, 16, 32, 64, 128, 256\}$  antennas. We can observe the result of  $K = 2$  and  $K = 6$  that higher performance can be

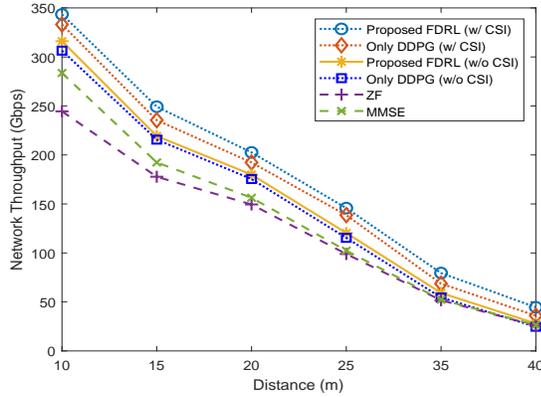


Fig. 6. The proposed FDRL algorithm compared to conventional and non-FEL benchmark optimization methods.

obtained due to advantageous FDRL of interference mitigation. Furthermore, the throughput performance is proportional to the number of antennas due to higher spatial diversity. In addition, throughput difference becomes increasingly larger when  $N_t = 256$  antennas are equipped by comparing the mechanism of THz CSI and of only SINR feedback. This is because that more hidden information from the estimated CSI is extracted and exchanged by FEL server, which is advantageous to interference cancellation.

Fig. 6 demonstrates the throughput considering distances between THz BS-UE. The throughput of all algorithms decreases to near zero due to limitation of severe intrinsic path loss from the THz channel. The proposed FDRL algorithm with estimated CSI can exchange more hidden information than that without CSI, which achieves higher throughput performance. In addition, DDPG lacks training model exchange from FEL because each BS conducts local training without model exchange, which results in lower throughput than the proposed FDRL algorithm. Moreover, the proposed FDRL is capable of exchanging sufficient hidden training models from powerful deep learning based DDPG, which outperforms the conventional beamforming methods of zero forcing (ZF) and minimum mean square error (MMSE).

## V. CONCLUSIONS

In this paper, an FDRL scheme was proposed to improve the network throughput by optimally searching the beamformers of the BSs under the situation of limited THz CSI. We have numerically demonstrated that FDRL is capable of exchanging more representative features among THz BSs to alleviate interference as well as to achieve higher throughput even when the BSs have limited CSI. With more deployed antenna arrays, the network is able to achieve a higher throughput due to higher spatial diversity. Moreover, we can observe a compelling tradeoff between overhead of exchange information and network throughput. Namely, the network can achieve high throughput at a low cost of uploading partial FDRL models. In a nutshell, the proposed FDRL scheme using DDPG

and FEL architectures outperforms the baseline methods with pure deep Q-learning because it can take the advantage of interference mitigation from information exchange. Also, we show that the proposed FDRL is superior to the existing beamforming techniques using non-learning methods, such as ZF and MMSE.

## REFERENCES

- [1] B. Ning, Z. Chen, W. Chen, Y. Du, and J. Fang, "Terahertz multi-user massive MIMO with intelligent reflecting surface: Beam training and hybrid beamforming," *IEEE Trans. Veh. Technol.*, vol. 70, no. 2, pp. 1376–1393, Feb. 2021.
- [2] C. Huang, Z. Yang, G. C. Alexandropoulos, K. Xiong, L. Wei, C. Yuen, and Z. Zhang, "Hybrid beamforming for RIS-empowered multi-hop terahertz communications: A DRL-based method," in *Proc. IEEE Globecom Workshops (GC Wkshps)*, Dec. 2020, pp. 1–6.
- [3] P. V. Tuan, T. Trung Duy, and I. Koo, "Multiuser MISO beamforming design for balancing the received powers in secure cognitive radio networks," in *Proc. IEEE Seventh International Conference on Communications and Electronics (ICCE)*, 2018, pp. 39–43.
- [4] G. Taricco, "On the beamforming capacity of MISO channels," vol. 1, no. 2, pp. 141–144, 2012.
- [5] F. Sotirani and W. Yu, "Hybrid analog and digital beamforming for mmwave OFDM large-scale antenna arrays," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 7, pp. 1432–1443, Jul. 2017.
- [6] S. Yang and Y. Liu, "Training efficiency of federated learning: A wireless communication perspective," in *Proc. International Conference on Wireless Communications and Signal Processing (WCSP)*, 2020, pp. 922–926.
- [7] X. Mo and J. Xu, "Energy-efficient federated edge learning with joint communication and computation design," *Journal of Communications and Information Networks*, vol. 6, no. 2, pp. 110–124, 2021.
- [8] K.-H. Liu, Y.-H. Hsu, W.-N. Lin, and W. Liao, "Fine-grained offloading for multi-access edge computing with actor-critic federated learning," in *Proc. IEEE Wireless Communications and Networking Conference (WCNC)*, 2021, pp. 1–6.
- [9] S. Zarandi and H. Tabassum, "Federated double deep q-learning for joint delay and energy minimization in IoT networks," in *Proc. IEEE International Conference on Communications Workshops (ICC Workshops)*, 2021, pp. 1–6.
- [10] W. Y. B. Lim, N. C. Luong, D. T. Hoang, Y. Jiao, Y.-C. Liang, Q. Yang, D. Niyato, and C. Miao, "Federated learning in mobile edge networks: A comprehensive survey," *IEEE Commun. Surveys Tuts.*, vol. 22, no. 3, pp. 2031–2063, 2020.
- [11] C. Lin and G. Y. Li, "Adaptive beamforming with resource allocation for distance-aware multi-user indoor terahertz communications," *IEEE Trans. Commun.*, vol. 63, no. 8, pp. 2985–2995, Aug. 2015.
- [12] S. Priebe and T. Kurner, "Stochastic modeling of THz indoor radio channels," *IEEE Trans. Wireless Commun.*, vol. 12, no. 9, pp. 4445–4455, Sep. 2013.
- [13] C. Han, A. O. Bicen, and I. F. Akyildiz, "Multi-ray channel modeling and wideband characterization for wireless communications in the terahertz band," *IEEE Trans. Wireless Commun.*, vol. 14, no. 5, pp. 2402–2412, May 2015.
- [14] X. Gao, L. Dai, Y. Zhang, T. Xie, X. Dai, and Z. Wang, "Fast channel tracking for terahertz beamspace massive MIMO systems," *IEEE Trans. Veh. Technol.*, vol. 66, no. 7, pp. 5689–5696, Jul. 2017.
- [15] 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," *Computer Science*, vol. 8, no. 6, 2015.