

Deep Reinforcement Learning-Based Intelligent Reflecting Surface Optimization for TDD Multi-User MIMO Systems

Fengyu Zhao^{ID}, Wen Chen^{ID}, Senior Member, IEEE, Ziwei Liu^{ID}, Jun Li^{ID}, Senior Member, IEEE, and Qingqing Wu^{ID}, Senior Member, IEEE

Abstract—In this letter, we investigate the discrete phase shift design of the intelligent reflecting surface (IRS) in a time-division duplexing (TDD) multi-user multiple-input-multiple-output (MIMO) system. We modify the design of deep reinforcement learning (DRL) scheme so that we can maximizing the average downlink data transmission rate free from the sub-channel channel state information (CSI). Based on the characteristics of the model, we modify the “proximal policy optimization (PPO)” algorithm and integrate gated recurrent unit (GRU) to tackle the non-convex optimization problem. Simulation results show that the performance of the proposed PPO-GRU surpasses the benchmarks in terms of performance, convergence speed, and training stability.

Index Terms—Intelligent reflecting surface (IRS), time-division duplexing (TDD), multi-user multiple-input-multiple-output (MU-MIMO), deep reinforcement learning (DRL).

I. INTRODUCTION

INTELLIGENT reflecting surface (IRS) is a low power technology that smartly tunes the radio signal propagation in wireless networks via a plurality of low-cost passive reflecting elements. Numerous influential works have been done on the configuration of continuous phase shifts of IRS with different design objectives. The authors in [1] studied an IRS-aided radar-communication (Radcom) scenario considering the cross-correlation design and the interference introduced by the IRS on the Radcom base station (BS). In [2], IRS-assisted simultaneous wireless information and power transfer (SWIPT) non-orthogonal multiple access (NOMA) networks are investigated to minimize BS transmit power. In [3], multiple access schemes are investigated in IRS-aided wireless-powered mobile edge computing (WP-MEC). However, all the aforementioned papers are based on the instantaneous/perfect channel state information (CSI) assumption. It is a practically difficult task to acquire the CSI of

Manuscript received 3 July 2023; revised 28 July 2023; accepted 31 July 2023. Date of publication 3 August 2023; date of current version 9 November 2023. This work was supported in part by the National Key Project under Grant 2020YFB1807700; in part by NSFC under Grant 62071296; and in part by Shanghai under Grant 22JC1404000, Grant 20JC1416502, and Grant PKX2021-D02. The associate editor coordinating the review of this article and approving it for publication was D. Mishra. (Corresponding author: Wen Chen.)

Fengyu Zhao, Wen Chen, Ziwei Liu, and Qingqing Wu are with the Department of Electronic Engineering, Shanghai Jiao Tong University, Shanghai 200240, China (e-mail: Aaronfy@sjtu.edu.cn; wenchen@sjtu.edu.cn; ziweiliu@sjtu.edu.cn; qingqingwu@sjtu.edu.cn).

Jun Li is with the School of Electronic and Optical Engineering, Nanjing University of Science Technology, Nanjing 210094, China (e-mail: jun.li@njust.edu.cn).

Digital Object Identifier 10.1109/LWC.2023.3301496

the channel between the IRS and its serving BS/users since IRS is passive. Secondly, previous IRS studies concentrate on continuous phase shifts at reflecting components, which are difficult to realize practically due to hardware limitations. Hence, we focus on discrete phase shifts of the IRS without reliance on sub-channel CSI.

Deep reinforcement learning (DRL) has been widely used to solve resource allocation problems in wireless networks. The optimization problems are transformed into the design of the Markov decision process (MDP). The major advantage of DRL is that the mobility of the wireless network (e.g., time varying channel, terminal mobility, and real-time control of IRS reflective elements) can be resolved during the process when the agent keeps interacting with the wireless environment. There have been various attempts to operate the IRS based on machine learning. In [4], multi-agent RL algorithm is firstly employed in multiple IRSs-assisted multi-user (MU) systems. In [5], a model-free control of IRS based on received pilot is accomplished with a modified version of double deep Q-network (DDQN) called DRL with extremum seeking control (ESC). However, the proposed schemes only compare to other non-DRL algorithms under different values of Rician factor, which is less persuasive. Secondly, the DRL with ESC doesn't compare with the single DRL without ESC in the experiments, so the performance improvement of ESC is not verified. Thirdly, it is incomprehensible that the large action space doesn't achieve better results than the small action space in the simulation. Inspired by those, we carry out discrete phase shift design based on DRL in a TDD MU-MIMO system free from the instantaneous sub-channel CSI.

The primary innovations of this letter can be summarized as follows:

- We introduced a novel approach for achieving discrete control of IRS in TDD multi-user MIMO systems. The key benefit of our design is that the IRS deployment enhances communication quality without relying on the instantaneous or statistical CSI of sub-channels.
- Based on the characteristics of our model, we have made appropriate modifications to the PPO algorithm, resulting in the creation of a modified version named PPO-GRU. The new algorithm incorporates three significant modifications, including:
 - 1) We integrated Gated Recurrent Unit (GRU), an improved type of Recurrent Neural Network (RNN), into the original PPO network structures of both the actor and critic. This modification allows the actor and critic to handle two types of state information,

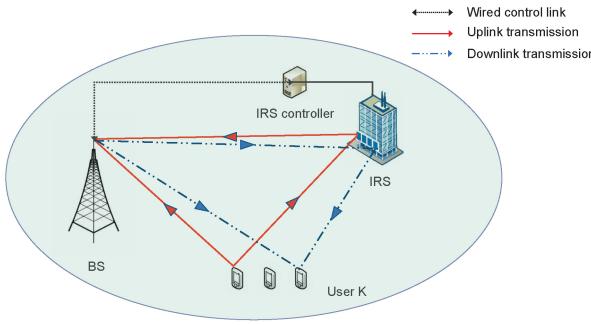


Fig. 1. The TDD multi-user MIMO scenario.

TABLE I
SUMMARY OF MAIN NOTATION

Notation	Description
N_B	Number of BS antennas
N_R	Number of IRS elements
N_K	Number of users
\mathbf{H}_{UB}	Channel gain between users and BS
\mathbf{H}_{UR}	Channel gain between users and IRS
\mathbf{H}_{RB}	Channel gain between IRS and BS
Φ	The reflection coefficients matrix of the IRS
K	The Rician factor
R	Downlink data sum rate
\mathcal{F}	A set of discrete angle

namely channel gains and angles, and deal with their correlation in the time domain within TDD systems.

- 2) Normalization: dynamic mean and variance values are maintained during simulation for all encountered states or advantages. The current state or advantage is then normalized accordingly.
- 3) We incorporated a strategy entropy term into the actor's loss function, ensuring the strategy's entropy remains as large as possible while optimizing the actor's loss.

Notations: A column vector is represented as a boldface lowercase letter, and a matrix is defined with a boldface capital letter. \odot represents the Hadamard product of a matrix, $(\cdot)^H$ is the conjugate transpose operation. For a set A , $|A|$ denotes the number of elements in the set. For a complex valued vector x , $|x|$ denotes L_1 norm. The operator $\text{diag}(\cdot)$ represents the diagonal matrix of a vector. The random variable x following the complex Gaussian distribution with zero-mean and unit variance is represented as $x \sim \mathcal{CN}(0, 1)$.

II. SYSTEM MODEL AND PROBLEM FORMULATION

As shown in Fig. 1, the communication procedure can be divided into uplink channel estimation and downlink data transmission.

A. Channel Estimation

At uplink transmissions, multiple users transmit orthogonal pilot signals to BS simultaneously. The received signal matrix can be represented as

$$\mathbf{Y} = \mathbf{X}\mathbf{H} + \mathbf{N}, \quad (1)$$

where $\mathbf{X} \in \mathbb{C}^{N_K \times N_K}$ is the pilot pattern, $\mathbf{N} \in \mathbb{C}^{N_K \times N_B}$ is the additive white Gaussian noise, which element follows the circularly symmetric complex Gaussian (CSCG) distribution $\mathcal{CN}(0, 1)$.

In an IRS-assisted wireless communication system, the uplink channel gain \mathbf{H} is given by

$$\mathbf{H} = \mathbf{H}_{UB} + \mathbf{H}_{UR}\Phi\mathbf{H}_{RB}, \quad (2)$$

where $\Phi = \text{diag}\{\gamma_1 e^{j\beta_1}, \gamma_2 e^{j\beta_2}, \dots, \gamma_{N_R} e^{j\beta_{N_R}}\}$ is the reflection coefficients matrix. γ_i and β_i depict the amplitude and phase shift reflecting coefficient of IRS element i respectively. We consider the phase shift of each IRS reflecting element restricted to a finite number of discrete value $\mathcal{F} = \{0, \Delta\theta, \dots, \Delta\theta(N_R - 1)\}$, where $\Delta\theta = 2\pi/N_R$.

Let $\mathbf{H}_{UB} \in \mathbb{C}^{N_K \times N_B}$, $\mathbf{H}_{UR} \in \mathbb{C}^{N_K \times N_R}$, $\mathbf{H}_{RB} \in \mathbb{C}^{N_R \times N_B}$ portray the channel gain from users to the BS, the channel gain from users to the IRS, and the channel gain from the IRS to the BS respectively. All the channels are modeled as the Rician channel [6], we take \mathbf{H}_{UB} as an example,

$$\mathbf{H}_{UB} = \sqrt{\frac{K}{K+1}} \mathbf{H}_{UB,LoS} + \sqrt{\frac{1}{K+1}} \mathbf{H}_{UB,NLoS}, \quad (3)$$

where K is the Rician factor, $\mathbf{H}_{UB,LoS}$ denotes the deterministic line-of-sight (LoS) component, and $\mathbf{H}_{UB,NLoS}$ signifies the fading non-line-of-sight (NLoS) component.

To get rid of the dependence on sub-channel CSI, we use the minimum mean square error (MMSE) to estimate the channel, which can be formulated as [7]

$$\hat{\mathbf{H}} = \mathbf{Y}\mathbf{X}^H \left(\mathbf{X}\mathbf{X}^H + \sigma_N^2 \mathbf{I} \right)^{-1}. \quad (4)$$

B. Data Transmission

At downlink transmission, zero-forcing (ZF) precoding is performed according to the reciprocity between uplink and downlink channel. The precoding matrix $\hat{\mathbf{A}}$ is represented as

$$\hat{\mathbf{A}} = [\hat{a}_1, \hat{a}_2, \dots, \hat{a}_K], \quad (5)$$

where \hat{a}_k is the k th power normalized vector of $(\hat{\mathbf{H}}^H \hat{\mathbf{H}})^{-1} \hat{\mathbf{H}}^H$.

The received signal at the k th user can be expressed as

$$y_k = \mathbf{a}_k^H \mathbf{h}_k x_k + \sum_{j \neq k}^{N_K} \mathbf{a}_j^H \mathbf{h}_k x_j + n_k, \quad (6)$$

where x_k is the signal to be sent to the k th user, and $n_k \sim \mathcal{CN}(0, \sigma_k^2)$ is the additive white Gaussian noise. Consequently, the signal-to-interference-plus-noise ratio (SINR) at the k th user can be written as

$$SINR_k = \frac{|\mathbf{a}_k^H \mathbf{h}_k|^2}{\sum_{j \neq k}^K |\mathbf{a}_j^H \mathbf{h}_k|^2 + \sigma_k^2}. \quad (7)$$

The achievable downlink data sum rate can be obtained as

$$R = \sum_{k=1}^{N_K} \log_2 (1 + SINR_k). \quad (8)$$

We only consider the fully reflective IRS, so our optimization problem is formulated as

$$\begin{aligned} \max_{\Phi(t)} \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^T R(t), \\ \text{s.t. } \gamma_i(t) = 1, \quad \forall i \in \{1, 2, \dots, N_R\}, \\ \beta_i(t) \in \mathcal{F}, \quad \forall i \in \{1, 2, \dots, N_R\}. \end{aligned} \quad (9)$$

III. DEEP REINFORCEMENT LEARNING-BASED SOLUTION

After obtaining the uplink CSI of the current time slot, the IRS controller adjusts the current phase based on the CSI and the phase of the previous time slot to improve the downlink transmission rate.

A. MDP Formula

- 1) Environment: We treat the whole wireless communication system except for the IRS controller as the environment, and the working mechanism of the environment is incomprehensible to it.
- 2) Agent: The IRS controller is served as the agent, which changes the configuration of IRS based on the performance feedback from the environment and the phase of each IRS element in the past.
- 3) State: We define the state S_t as the combination of the channel estimation and IRS phase at time slot $t - 1$,

$$S_t = \{S_t^\Phi, S_t^H\}, \quad (10a)$$

$$S_t^\Phi \triangleq [\Re\{\Phi_{t-1}\}, \Im\{\Phi_{t-1}\}], \quad (10b)$$

$$S_t^H \triangleq [\Re\{\hat{H}_{t-1}\}, \Im\{\hat{H}_{t-1}\}]. \quad (10c)$$

- 4) Action: The action is defined as the amount of change in phase from Φ_{t-1} ,

$$\Phi_t = \Phi_{t-1} \odot \Delta\Phi_t. \quad (11)$$

The phase shift $\Delta\Phi$ is limited to the subset (or full set) of N_R point discrete Fourier transform (DFT) vectors $v(k)$,

$$v(k) = \left[1, e^{\frac{j\pi k}{N_R}}, \dots, e^{\frac{j\pi(N_R-1)k}{N_R}} \right]. \quad (12)$$

For an example, when the size of the action space $|A| = 2n + 1$, we could set action space $A = \{v(-n), v(-n+1), \dots, v(0), v(1), \dots, v(n)\}$. IRS Controller will select $v(k) \in A$ depending on its current policy function, then the phase shift at time slot t will be

$$\Delta\Phi_t = \text{diag}\{v(k)\}. \quad (13)$$

Firstly, it has been demonstrated in [8] that utilizing an IRS with 2-bit phase shifters can achieve the same asymptotic squared power gain as the ideal scenario with continuous phase shifts. Therefore, discrete phase control based on DFT matrix is already sufficiently effective [9], [10]. Secondly, increasing the number of IRS elements only increases the dimensionality of the $v(k)$ vector. The agent still selects from A even when the size of the discrete phase set is fixed at $|A| = 2n + 1$.

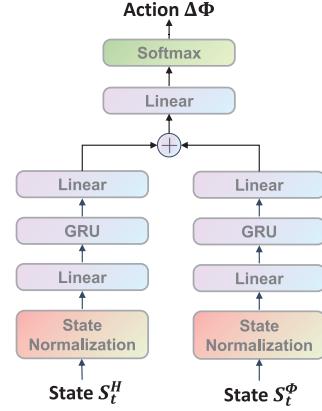


Fig. 2. Revised PPO actor network framework.

As a result, our design is capable of achieving good convergence performance even with a larger number of reflecting elements.

- 5) Reward: We adopt the downlink data sum rate in (8) as the reward.

B. IRS Control Using Improved PPO Algorithm

Our modifications to PPO can be summarized as modifying the network architecture, normalizing state and advantage, and modifying the loss function equation. We will explain each of these modifications in the order listed.

Given that the state in our MDP contains two types of information, namely phase and channel gain, and is correlated in the time domain, we have made modifications to the actor and critic network structures in the traditional PPO algorithm to accommodate this model. The modified actor and critic network structures are similar in nature, we present the actor network as an example in Fig. 2.

Prior to inputting the states into the neural network, we first perform separate normalization for the two types of states. We demonstrate the normalization process using S_t^Φ . At time slot t , all the values of the state $S_t^\Phi = [S_1^\Phi, S_2^\Phi, \dots, S_t^\Phi]$ are recorded, and the normalized state is calculated using the following equation:

$$\hat{S}_t^\Phi = \frac{S_t^\Phi - \mu}{\sigma}, \quad (14)$$

where \hat{S}_t^Φ represents the normalized state of phase, μ is the mean of all the state values at time slot t , and σ is the standard deviation.

In our system model, S_t^Φ and S_t^H typically correlate with a prior data point or a data point spanning a time period. The LoS component within the channel gain experiences minimal variation if the user's position changes only slightly between time slots $t - 1$ and t , ensuring that the receiving antenna stays stationary. Furthermore, Φ_t is a time-dependent sequence data, as its value relies on the preceding phases $\Phi_1, \Phi_2, \dots, \Phi_{t-1}$ as described in (11). Consequently, we incorporate two separate GRUs following two linear layers. The features extracted from the linear layer serve as inputs for the GRUs, enabling them to capture long-term dependencies in this sequence data. This

enhances the model's accuracy and generalization capabilities. The results extracted from S_t^Φ and S_t^H after passing through a three-layer network are added together and then input into another linear layer. In this way, the selected discrete phase of IRS is based on the consideration of the two types of state information.

Advantage is also normalized in our approach, which we refer to as mini batch normalization. After calculating the advantages using General Advantage Estimation (GAE) for a batch [8], instead of directly normalizing the entire batch's advantages, we normalize the advantages of the current mini-batch before using it to update the policy in each iteration. Compared to the original PPO, our improved algorithm requires additional control over two hyperparameters: batch size D and sample mini-batch N_D . The loss function of actor L_{actor} with mini batch advantage normalization is given by

$$L_{actor} = \min \left(\frac{\pi_\theta(a | s)}{\pi_{\theta_k}(a | s)} \hat{A}^{\pi_{\theta_k}}(s, a), g(\epsilon, \hat{A}^{\pi_{\theta_k}}(s, a)) \right), \quad (15)$$

where

$$g(\epsilon, \hat{A}^{\pi_{\theta_k}}(s, a)) = \begin{cases} (1 + \epsilon) \hat{A}^{\pi_{\theta_k}}(s, a) & \hat{A}^{\pi_{\theta_k}}(s, a) \geq 0 \\ (1 - \epsilon) \hat{A}^{\pi_{\theta_k}}(s, a) & \hat{A}^{\pi_{\theta_k}}(s, a) < 0 \end{cases}, \quad (16)$$

in which ϵ is a hyperparameter which roughly controls the variation between the new policy and the old one, and $\hat{A}(s, a)$ is the normalized advantage function.

Thirdly, we modify the loss function expression above. Referring to the definition of entropy in information theory and probability statistics, the entropy of a strategy is represented as

$$\mathcal{H}(\pi(\cdot | s_t)) = - \sum_{a_t} \pi(a_t | s_t) \log(\pi(a_t | s_t)). \quad (17)$$

The greater the entropy of a strategy, the more evenly distributed the probabilities of selecting each action are. To improve the exploration capability of the algorithm, we add a term for strategy entropy to the actor's loss L_{actor} , multiply it by a coefficient δ , and optimize L_{actor} while maximizing the strategy's entropy. The modified loss function L_{actor}' is given by

$$L_{actor}' = L_{actor} + \delta * \mathcal{H}(\pi(\cdot | s_t)). \quad (18)$$

where δ is the entropy coefficient.

IV. SIMULATION RESULTS

A. Simulation Settings

We establish our model in a three-dimensional Cartesian coordinate system. The BS is located at the coordinate [0,0,0], and IRS is placed at the coordinate [5,5,5]. The UEs whose height is ranging from 1.5m to 1.8m are uniformly distributed in a circle area with radius equal to 10m. The Rician factor $K = 10$. The LoS component varies every 20 seconds by randomly selecting the user positions within the circle, while the NLoS component varies every second.

Algorithm 1 Proximal Policy Optimization Based IRS Control in TDD Multi-User MIMO Systems

Input: Initial IRS controller policy parameters θ_0 , initial value function parameters ϕ_0 .

for $k = 0, 1, 2, \dots$ **do**

 Collect the trajectories into a set $\mathcal{D}_k = \{\tau_i\}$ by running current policy $\pi_k = \pi(\theta_k)$ in the environment.

 Randomly select a mini batch of trajectories. Compute rewards-to-go $\hat{R}_t = \sum_{t'=t}^T R(s_{t'}, a_{t'}, s_{t'+1})$.

 Compute advantage estimates \hat{A}_t using GAE method and do the mini batch normalization.

 Update the policy by minimizing the loss function defined in (18).

 Update value function by minimizing the mean-squared error between value function $V_\phi(s_t)$ and \hat{R}_t ,

$$\phi_{k+1} = \arg \min_{\phi} \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T (V_\phi(s_t) - \hat{R}_t)^2.$$

end for

The channel matrix \mathbf{H}_{UB} , \mathbf{H}_{UR} and \mathbf{H}_{RB} are generated in a similar manner, we take \mathbf{H}_{RB} as an example. The LoS channel gain between IRS and BS is

$$\mathbf{H}_{RB, LoS} = \mathbf{v}_B \mathbf{v}_R^H, \quad (19)$$

where the steering vectors are formulated as

$$\mathbf{v}_R = \mathbf{v}(\Psi_R, N_{R,y}) = [1, e^{j\pi\Psi_R}, \dots, e^{j(N_{R,y}-1)\pi\Psi_R}]^T,$$

$$\mathbf{v}_B = \mathbf{v}(\Psi_B, N_{B,x}) = [1, e^{j\pi\Psi_B}, \dots, e^{j(N_{B,x}-1)\pi\Psi_B}]^T.$$

According to [11], the directional cosines Ψ_R, Ψ_B are represented as

$$\Psi_R = \mathbf{e}_R^T \mathbf{e}_{BR}, \quad (20a)$$

$$\Psi_B = \mathbf{e}_B^T \mathbf{e}_{BR}. \quad (20b)$$

We place the uniform linear array (ULA) of the BS at the coordinate [1,0,0], and assume that the reflector array of IRS can be regarded as a ULA placed at [0,1,0], so $e_R = [0, 1, 0]^T$ and $e_B = [1, 0, 0]^T$. The NLoS components follow the complex Gaussian distribution $\mathcal{CN}(0, 1)$. \mathbf{H}_{UB} and \mathbf{H}_{UR} are calculated in the same way.

B. Comparisons With Benchmarks

We set the number of users $N_K = 2$, the number of BS antennas $N_B = 2$, the action space size of IRS controller $|A| = 5$, and the number of the elements of IRS $N_R = 32$. Other network parameters settings are listed in Table II.

In Figs. 3 and 4, we compare the proposed PPO-GRU scheme with four benchmarks: random reflection, multi-armed bandit (MAB), DDQN in [5] and original PPO. All the schemes use the same action sets but handle the information differently. Random reflection cannot analyze and utilize all kinds of information, resulting in the poorest performance. MAB is a simpler version of DQN which builds the connection between reward and action by calculating the reward distribution of all the arms. However, it fails to fully utilize the

TABLE II
NETWORKS PARAMETERS

Total Training steps T	10000000
Batch size D	2048
Sample mini-batch N_D	64
Discount factor γ	0.99
GAE parameter λ	0.95
Learning rate of actor network α_h	0.0003
Learning rate of the critic network α_p	0.0003
PPO clip parameter ϵ	0.2
PPO entropy coefficient δ	0.01

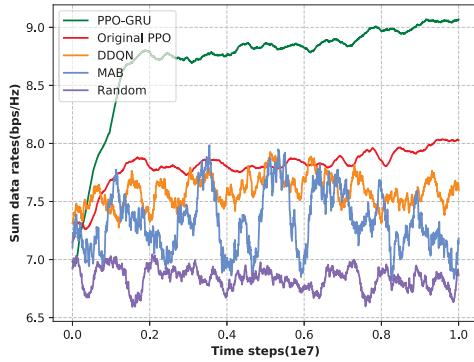


Fig. 3. Performance comparisons of different algorithms.

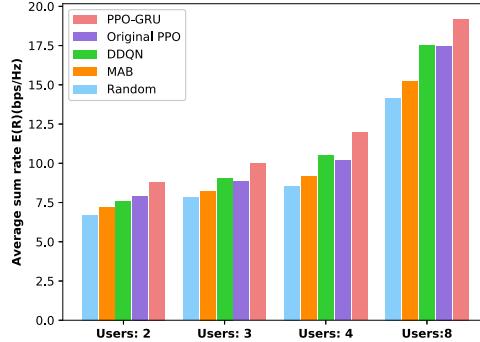


Fig. 4. Comparisons of algorithm performance with different numbers of users.

state information in our model. Compared to other DRL algorithms, such as the original PPO and DDQN, our PPO-GRU achieves better performance and faster convergence. Our modifications can enhance the model's accuracy and generalization capabilities, enabling the PPO-GRU to capture the changing dynamics of the wireless environment and adapt the agent's policy accordingly.

In Fig. 5, we study the impact of different action spaces A on the performance of the proposed PPO scheme. When $|A| = 3$, the algorithm converges at the fastest speed but ends up with the lowest data rate. When $|A| = 11$, the agent performs poorly at the beginning of the training but results in the best performance. However, it will take too much time to converge if we increase the number of users. When $|A| = 5$, it achieves nearly good performance and converges much faster. So a moderate size of action space is best for practical deployment.

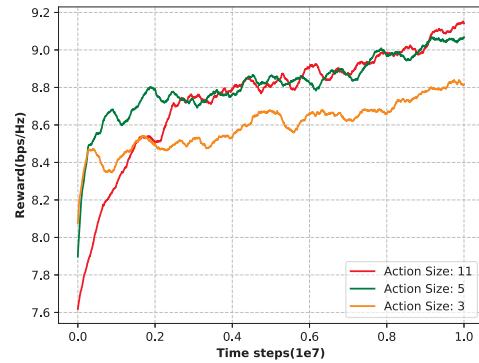


Fig. 5. Performance of PPO with different action sets.

V. CONCLUSION

This letter investigates the problem of maximizing the average data sum rate in IRS-assisted TDD MU-MIMO networks under the constraints of discrete phase shifts. To address this challenging problem, we propose an improved PPO algorithm. Simulation results demonstrate that our modified PPO algorithm outperforms previous algorithms in various scenarios. Moreover, we show that a well-designed action space can achieve both high training efficiency and good performance.

REFERENCES

- [1] M. Hua, Q. Wu, C. He, S. Ma, and W. Chen, "Joint active and passive beamforming design for IRS-aided radar-communication," *IEEE Trans. Wireless Commun.*, vol. 22, no. 4, pp. 2278–2294, Apr. 2023.
- [2] Z. Li, W. Chen, Q. Wu, K. Wang, and J. Li, "Joint beamforming design and power splitting optimization in IRS-assisted SWIPT NOMA networks," *IEEE Trans. Wireless Commun.*, vol. 21, no. 3, pp. 2019–2033, Mar. 2022.
- [3] G. Chen, Q. Wu, W. Chen, D. W. K. Ng, and L. Hanzo, "IRS-aided wireless powered MEC systems: TDMA or NOMA for computation offloading?" *IEEE Trans. Wireless Commun.*, vol. 22, no. 2, pp. 1201–1218, Feb. 2023.
- [4] J. Zhang et al., "Collaborative intelligent reflecting surface networks with multi-agent reinforcement learning," *IEEE J. Sel. Topics Signal Process.*, vol. 16, no. 3, pp. 532–545, Apr. 2022.
- [5] W. Wang and W. Zhang, "Intelligent reflecting surface configurations for smart radio using deep reinforcement learning," *IEEE J. Sel. Areas Commun.*, vol. 40, no. 8, pp. 2335–2346, Aug. 2022.
- [6] A. Goldsmith, *Wireless Communications*. Cambridge, U.K.: Cambridge Univ. Press, 2005.
- [7] H. Poor, *An Introduction to Signal Detection and Estimation*, 2nd ed. Berlin, Germany: Springer, 1994.
- [8] Q. Wu and R. Zhang, "Beamforming optimization for intelligent reflecting surface with discrete phase shifts," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, Brighton, U.K., May 2019, pp. 7830–7833.
- [9] T. L. Jensen and E. De Carvalho, "An optimal channel estimation scheme for intelligent reflecting surfaces based on a minimum variance unbiased estimator," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, Barcelona, Spain, 2020, pp. 5000–5004.
- [10] Z. Sun and Y. Jing, "On the performance of multi-antenna IRS-assisted NOMA networks with continuous and discrete IRS phase shifting," *IEEE Trans. Wireless Commun.*, vol. 21, no. 5, pp. 3012–3023, May 2022.
- [11] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel, "Highdimensional continuous control using generalized advantage estimation," in *Proc. Int. Conf. Learn. Represent.*, 2016. [Online]. Available: <https://arxiv.org/abs/1506.02438>
- [12] W. Wang and W. Zhang, "Jittering effects analysis and beam training design for UAV millimeter wave communications," *IEEE Trans. Wireless Commun.*, vol. 21, no. 5, pp. 3131–3146, May 2022.