

# Meta Learning for Meta-Surface: A Fast Beamforming Method for RIS-Assisted Communications Adapting to Dynamic Environments

Qinpei Luo, Boya Di

School of Electronics Engineering and Computer Science, Peking University, Beijing, China  
 {luoqinpei, diboya}@pku.edu.cn

**Abstract**—Recently reconfigurable intelligent surface (RIS) has been proposed as a promising technique to enhance the capacity of wireless networks by reshaping the electromagnetic characteristics of the environment. However, given numerous RIS elements, it is non-trivial to design an efficient beamforming scheme especially for the real-time mobile applications that require fast response to varying environments. In this paper, aiming to maximize the sum rate of a multi-user system via the RIS-enabled beamforming design, a meta-critic network is proposed to recognize the environment change and automatically perform the self-updating of the learning model. We also develop a stochastic Explore and Reload procedure to alleviate the high-dimensional action space issue. Simulation results demonstrate that the proposed scheme can converge to a higher sum rate more rapidly compared to the state-of-the-art methods in dynamic settings. The robustness of our proposed scheme against different RIS sizes is also verified.

## I. INTRODUCTION

As a novel enabling technology for the next-generation communications, meta surfaces is able to improve spectrum efficiency, among which reconfigurable intelligent surface (RIS) has attracted great attention owing to its ability of desirable signal reflection. Users can be well-served through programming the phase shifts of the RIS elements. To deal with the large number of RIS elements, machine learning based beamforming schemes have been explored, among which reinforcement learning (RL) has served as a potential tool to depict the interaction process between the surface and the environment. However, traditional RL methods may be time-consuming due to its re-training process, making it not feasible especially for the real-time dynamic environments.

In this paper, we propose a meta-critic deep deterministic policy gradient (MC-DDPG) scheme for the RIS-based beamforming adapting to dynamic user's locations. A meta-critic is designed which serves as an automotive tool for real-time model parameter generation in new environments by learning from multiple scenario-specific tasks. Simulation results show that **with only a small amount of cascaded channel information, MC-DDPG outperforms the traditional RL method and an iterative algorithm in terms of the sum rate and the convergence speed in dynamic environments. The robustness of the MC-DDPG scheme against different RIS sizes is also verified, which reveals the feasibility of large-scale MIMO enabled by RIS.**

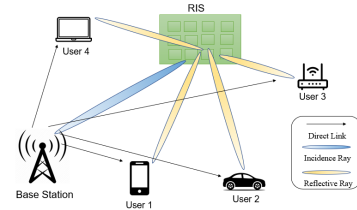


Fig. 1. System model of the RIS-assisted multi-user system

## II. SYSTEM MODEL

Fig. 1 shows a downlink multi-user MISO wireless communication system with a  $M$ -antenna base station (BS) serving  $K$  users each equipped with one antenna. The Line-of-Sight (LoS) channel between the BS and users are often unstable and suffering from severe fading. An RIS consisting of  $N$  elements is deployed between the users and the BS to reflect the transmit signals towards the users.

The direct channel between the BS and  $K$  users can be denoted as  $\mathbf{H}_{BU} \in \mathbb{C}^{K \times M}$ . The BS-RIS link and the RIS-user link can be denoted by  $\mathbf{H}_{BI} \in \mathbb{C}^{N \times M}$  and  $\mathbf{H}_{IU} \in \mathbb{C}^{K \times N}$ , respectively. We assume that each channel of  $\mathbf{H}_{BU}$ ,  $\mathbf{H}_{BI}$ ,  $\mathbf{H}_{IU}$  follows the Markov process with the transition of time slots. The equivalent end-to-end channel from the BS to user  $k$  can be given as

$$\mathbf{H}_k = \mathbf{H}_{IU,k} \mathbf{\Theta} \mathbf{H}_{BI} + \mathbf{H}_{BU,k}, k \in \mathcal{K}_u, \quad (1)$$

in which  $\mathbf{\Theta} \in \mathbb{C}^{N \times N} = \text{diag}([e^{j\theta_1}, \dots, e^{j\theta_N}])$ ,  $[e^{j\theta_1}, \dots, e^{j\theta_N}]$  being the phase shifts configuration of RIS elements.

## III. PROBLEM FORMULATION AND ALGORITHM DESIGN

### A. Sum Rate Maximization Problem

We consider the sum rate maximization problem in  $T$  time slots, each of which has a duration of  $\Delta T$ . The signal-to-noise (SNR) and data rate of user  $k$  can be expressed as

$$\gamma_{k,t} = \frac{|(\mathbf{H}_{IU,k} \mathbf{\Theta} \mathbf{H}_{BI} + \mathbf{H}_{BU,k}) \mathbf{V}_{k,t} m_k|^2}{|(\mathbf{H}_{IU,k} \mathbf{\Theta} \mathbf{H}_{BI} + \mathbf{H}_{BU,k}) \sum_{j=1, \neq k}^K \mathbf{V}_{j,t} m_j|^2 + \sigma_{k,t}^2}, \quad (2)$$

$$R_{k,t} = |\Delta T \log(1 + \gamma_{k,t})|. \quad (3)$$

$V_{j,t}$  refers to the digital beamforming vector from the BS to the  $j$ -th user,  $m_k$  denotes the symbol BS sends to user  $j$ , and  $n_{k,t}$  represents Gaussian noise which follows  $N(0, \sigma_{k,t}^2)$ .

For a fixed digital beamforming scheme such as ZF [1] or MMSE, we can get a sub-optimal solution of  $V_t$  directly, the sum rate maximization problem can be formulated as

$$\max_{\Theta_t} \sum_{t=1}^T \sum_{k=1}^K R_{k,t}, \quad (4)$$

Given the time-varying characteristics of channels, we then reformulate it as a MDP consisting of the following components: 1) **Action**:  $a_t = \Theta_t, \forall \theta_t \in \Theta_t, \theta_t \in (-\pi, \pi)$ . 2) **State**:  $s_t = \{\mathbf{H}_t, \Theta_{t-1}\}$ , where  $\mathbf{H} = \mathbf{H}_{IU}\Theta\mathbf{H}_{BI} + \mathbf{H}_{BU}$ . 3) **Reward**:  $r_t = \eta \sum_{k=1}^K \gamma_{k,t}$ , where  $\eta$  is a coefficient.

### B. MC-DDPG Algorithm Design

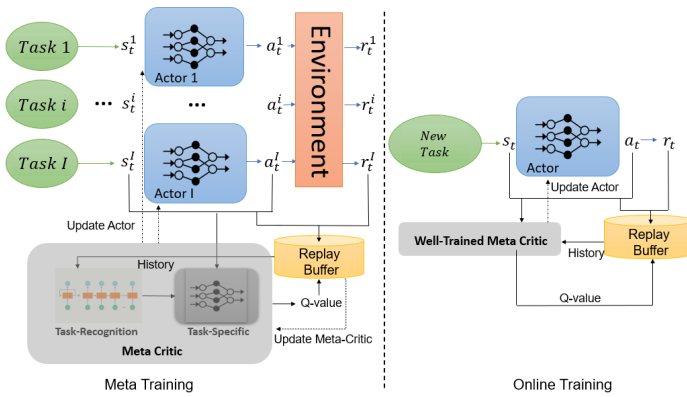


Fig. 2. Framework of MC-DDPG

The design of our proposed MC-DDPG is shown in Fig. 2, consisting of the meta learning phase and online learning phase. The experience from all actors will be sent to one meta critic for updating. The task it refers to a process of the BS maximizing the sum rates of all users. For different tasks, the channel states and locations of users are different.

The meta learning phase can be described as follows: First, for each task  $i$ , its current state  $s_t^i$  is fed to actor  $i$  so the actor can generate an action  $a_t^i$  from policy  $\pi(a|s_t^i)$ . Task  $i$  operates the action and receives the reward  $r_t^i$  from the environment, then the transition tuples  $\langle s_t^i, a_t^i, r_t^i \rangle$  will be stored in the replay buffer. Meanwhile, the meta critic collects the history information of task  $i$ , i.e.,  $\mathcal{H}_t = \{(s_k, s_{k+1}, a_k, r_k)\}, k \in [t - \bar{t}, t - 1]$  from the replay buffer together with the state-action pair  $\langle s_t^i, a_t^i \rangle$  to give a task specific Q-value for updating the actor networks. The meta critic is also updated by the trajectories of all tasks in the replay buffer. In the online learning phase, for a newly-coming real-time task, the update of the actor network is the same with that of the meta learning phase while the critic is kept static.

## IV. SIMULATION RESULTS

For the simulation, we set the number of antennas of BS  $M = 8$ , number of users  $K = 4$ . We compared it to

### Algorithm 1 MC-DDPG Algorithm

- 1: **Meta Training**;
- 2: **input**: Multiple task samples.
- 3: **Initialize**: (For each task  $i$ ) Critic Networks and actor network; Target Networks; Replay Buffer;
- 4: **for**  $eps$  in range( $MaxEpisode$ ) **do**
- 5:   Sample  $I$  tasks and initialize states.
- 6:   **for**  $t$  in range( $MaxStep$ ) **do**
- 7:     **for** each task  $i$  **do**
- 8:       Select action and get reward and next state.
- 9:       Store transition tuple into the replay buffer
- 10:       Sample a batch of data and update the Meta Critic.
- 11:       Update the actor networks with delay.
- 12:       Soft-update target networks with delay.
- 13:   **Output**: Well-trained meta critic.
- 14: 

---
- 15: **Online Training**;
- 16: **input**: A new task; Well-trained meta critic.
- 17: **Initialize**: Policy network; Replay Buffer.
- 18: **for**  $eps$  in range( $MaxEpisode$ ) **do**
- 19:   **Initialize** system state.
- 20:   **for**  $t$  in range( $MaxStep$ ) **do**
- 21:     Select action and get reward and next state.
- 22:     Store transition tuple into the replay buffer.
- 23:     Sample a batch of data and update the actor network.
- 24: **Output**: Trained policy of actor.

two benchmark algorithms: 1) Twin delayed deep deterministic policy gradient (TD3) [2] without the meta critic and 2) Zero-Force Exhausting (ZF Exhaust) [1], where the digital beamforming bases on the ZF method, and the RIS phase shift optimization is performed via the exhaustion attack method. We assume that the task is updated every 300 episodes, each of which consists of 20 time slots.

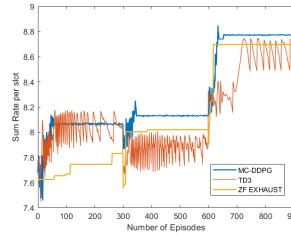


Fig. 3. Sum rate performance with respect to varying users' locations

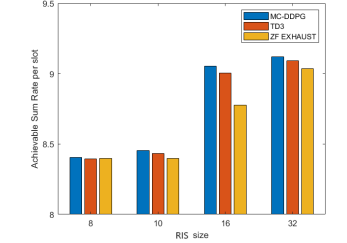


Fig. 4. Achievable sum rate v.s. the number of RIS elements

Fig. 3 shows the performance of the proposed scheme when mobile users move rapidly. The proposed MC-DDPG can rapidly converge to a higher sum rate compared to the benchmarks. As the MC-DDPG can converge within 100 episodes and we set each user's position changes 0.01m each episode, we remark that it can support the user mobility at a minimum speed of 36 km/h. Fig. 4 shows the sum rate varying with the number of RIS elements  $N$ . We observe that the proposed MC-DDPG converges to a higher sum rate as  $N$  increases from 64 to 1024 compared to two benchmarks, which shows its capability of supporting large-scale RIS.

## REFERENCES

- [1] B. Di, H. Zhang, L. Song, Y. Li, Z. Han, and H. V. Poor, "Hybrid beamforming for reconfigurable intelligent surface based multi-user communications: Achievable rates with limited discrete phase shifts," *IEEE J. Selected Areas Commun.*, vol. 38, no. 8, pp. 1809–1822, Aug. 2020.
- [2] S. Fujimoto, H. Hoof, and D. Meger, "Addressing function approximation error in actor-critic methods," in *Int. Conf. Machine Learning (ICML)*, Jul. 2018, pp. 1587–1596.