

# DRL-Based Joint Trajectory Planning and Beamforming Optimization in Aerial RIS-Assisted ISAC System

Xipeng Chen<sup>1</sup>, Xiaowen Cao<sup>1</sup>, Lifeng Xie<sup>2</sup>, and Yejun He<sup>1</sup>

<sup>1</sup> State Key Laboratory of Radio Frequency Heterogeneous Integration

Sino-British Antennas and Propagation Joint Laboratory, MOST

Guangdong Engineering Research Center of Base Station Antennas

Shenzhen Key Laboratory of Antennas and Propagation

College of Electronics and Information Engineering, Shenzhen University, 518060, China

<sup>2</sup> Department of Broadband Communication, Peng Cheng Laboratory, Shenzhen, China

Email: 2164486229@qq.com, caoxwen@szu.edu.cn, xielf@pcl.ac.cn, heyejun@126.com

**Abstract**—In this paper, we consider a reconfigurable intelligent surface (RIS) enabled integrated sensing and communication (ISAC) system, where the RIS is mounted on an unmanned aerial vehicle (UAV) to enhance signal quality and connectivity coverage by exploring the flexible mobility and adjusting the phase and amplitude of reflected signals. Towards this end, a communication rate max-min problem is formulated by jointly optimizing the beamforming vector, RIS phase shift and UAV trajectory under the constraints of power consumption and sensing requirement. However, due to the coupled variables, this problem is a non-convex problem and hard to be optimally solved. Hence, we reformulate the primary problem as a sequential decision-making problem and exploit a deep reinforcement learning (DRL)-based solution to find a tractable solution. Numerical results validate the effectiveness and the superiority of the proposed algorithm compared with the benchmark schemes.

**Index Terms**—Reconfigurable intelligent surface, unmanned aerial vehicle, integrated sensing and communication, Deep reinforcement learning.

## I. INTRODUCTION

As the sixth generation (6G) networks develop, various intelligent applications have emerged like smart cities, smart transportation, and industrial Internet of Things, which requires high-precision sensing capabilities and high-throughput wireless communication. However, the surge in the demand for connectivity and the proliferation of devices including sensors and smart devices, intensify the congestion within the available frequency bands, leading to a degradation on network performance [1].

In view of this, the integrated sensing and communications (ISAC) has emerged as a promising approach to enhance the sensing capability and improve spectrum utilization [2], which has recognized by International Mobile Telecommunications (IMT)-2030 as one of the core usage scenarios for 6G. It allows simultaneous data transmission and wireless sensing on the same spectrum. On the other

hand, owing to the capability to create a more favorable propagation environment, reconfigurable intelligent surface (RIS) has also emerged as an innovative technique in 6G system [3]. As RIS could assist to establish additional non-line-of-sight (NLoS) links for exploiting more degrees of freedom (DoFs) to promote communication performance, the introduction of RIS technology into ISAC system may offer a solution to overcome limitations caused by traditional propagation environment, which has attracted extensive research [4]–[6]. Especially, in a radar-communication coexistence system, the authors in [4] formulated a joint beamforming and communication covariance matrix optimization problem to maximize the radar signal-to-interference-plus-noise ratio (SINR). In a dual-functional radar communication (DFRC) system, the researchers in [5] reduced the Cramér-Rao bound (CRB) on the estimation error and expanded it to the extended target models. Besides, the authors in [6] considered two scenarios with/without the interference at the receiver end to enhance target detection probability.

While current studies have tackled the joint optimization of RIS-aided ISAC in multiple ways, the efficiency might be limited by the fixed location of the RIS in the network. The rapid development of aerial platforms, particularly unmanned aerial vehicles (UAVs), has introduced a dynamic element that can serve as temporary base stations (BS) or relays, offering a flexible deployment solution to enhance ISAC performance [7], which has attracted many attentions. For instance, the authors in [8] jointly designed the UAV maneuver and the beamforming to maximize the weighted sum-rate, while a joint ground user-UAV association and real-time trajectory problem was solved to satisfy the sensing or communication performance in [9]. On the other side, to fully unlock the potential of RIS, researchers have explored the concept of Aerial RIS (ARIS) with flexible deployment by introducing UAV-mounted RIS. In [10], an

AIRS-assisted secure transmission design was proposed to maximize the worst-case sum secrecy rate, while the authors in [11] considered an ARIS-assisted edge computing scheme to facilitate offloading computing tasks from ground user to access point. Additionally, it is worthy to notice that the UAV's trajectory optimization is always a hot research topic but often poses several challenges under specific application and environmental conditions due to its limitation in weight and battery. Compared with traditional solutions, deep reinforcement learning (DRL) has been deemed as a dynamic and adaptive solution to solve such trajectory optimization problems [12], by leveraging neural networks to learn optimal policies through interactions with the environment.

Despite the existing progress on UAV-enabled ISAC and ARIS system, it's notable that the intersection of ARIS and ISAC has largely remained unexplored, which thus motivates our work. In this paper, we propose an ARIS-assisted ISAC (ARIS-ISAC) system. The objective is to maximize the sum of the minimum achievable rate through the joint optimization of UAV trajectory, RIS phase shift and beamforming to guarantee the communication performance of the users. The mobility of the UAV makes the optimal beamforming and RIS phase dynamically change, which further makes the optimization design more complicated. Thus, we exploit the DRL-based algorithm to solve the primary problem by reformulating it as a sequential decision-making problem.

## II. SYSTEM MODEL

As shown in Fig. 1, an ARIS-assisted ISAC downlink system is proposed. Specifically, an ISAC BS with  $M$  transmitting/receiving antennas serves  $K$  single-antenna users and simultaneously detects  $J$  targets, where user set and target set are denoted as  $\mathcal{K} = \{1, 2, \dots, K\}$  and  $\mathcal{J} = \{1, 2, \dots, J\}$ . In order to lessen the power consumption of UAV and improve the flexibility of RIS, a RIS containing  $N$  elements is assembled on the UAV to assist the downlink communication and sensing simultaneously.

According to whether direct link is blocked by obstacles, users are divided into the blocked users with set  $\mathcal{K}_B$  and the unblocked users with set  $\mathcal{K}_U$  with  $\mathcal{K} = \mathcal{K}_B \cup \mathcal{K}_U$ . It is assumed that the BS-target link is unblocked due to the fact that the pathloss of the BS-ARIS-target link exists the multiple effect.

In addition, the UAV flies in 2D space. And the UAV departs from a initial point  $\mathbf{q}_i = [x_i, y_i]$  and flies at a fixed altitude  $H$ . The ISAC period  $T$  is divided into  $L$  time slots with length being  $\tau$  and the UAV is at  $\mathbf{q}_t = [x_t, y_t]$  at the  $t$ -th time slot. We also assume that the UAV operates within a specific area  $\mathcal{D}$ , that is,  $\mathbf{q}_t \in \mathcal{D}, \forall t$ . Due to the limited velocity of the UAV, the velocity is subject to

$$\frac{\|\mathbf{q}_t - \mathbf{q}_{t-1}\|}{\tau} \leq \nu_{\max}, \quad (1)$$

where  $\nu_{\max}$  represents the maximum velocity of the UAV.

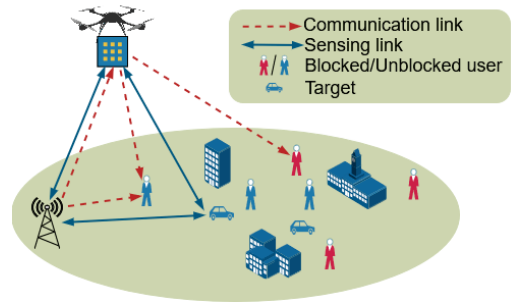


Fig. 1. The ARIS-assisted ISAC system.

The flight energy consumption of UAV at time slot  $t$  is modeled as [13]

$$e_t^{\text{UAV}} = \tau \left( P_0 \left[ 1 + \frac{3(V_t^{\text{xy}})^2}{U_{\text{tip}}} \right] + P_1 \left( \sqrt{1 + \frac{(V_t^{\text{xy}})^4}{4V_0^4}} - \frac{(V_t^{\text{xy}})^2}{2V_0^2} \right)^{\frac{1}{2}} + C_0 (V_t^{\text{xy}})^3 \right). \quad (2)$$

where  $V_t^{\text{xy}} = \frac{\|\mathbf{q}_t - \mathbf{q}_{t-1}\|}{\tau}$  is the velocity of the UAV.

### A. Transmitting Signal

The signal transmitted by the ISAC BS at time slot  $t$  is

$$\mathbf{x}_t = \sum_{k=1}^K \mathbf{w}_{k,t}^c s_{k,t}^c + \sum_{j=1}^J \mathbf{w}_{j,t}^s s_{j,t}^s, \quad (3)$$

where  $\mathbf{w}_{k,t}^c \in \mathbb{C}^{M \times 1}$  and  $\mathbf{w}_{j,t}^s \in \mathbb{C}^{M \times 1}$  denote the beamforming vectors for communication with the user  $k$  and sensing with the target  $j$  at time slot  $t$ , respectively. The according data streams are  $s_{k,t}^c$  and  $s_{j,t}^s$ .

Let  $\mathbf{W}_t^c = [\mathbf{w}_{1,t}^c, \dots, \mathbf{w}_{K,t}^c] \in \mathbb{C}^{M \times K}$ ,  $\mathbf{W}_t^s = [\mathbf{w}_{1,t}^s, \dots, \mathbf{w}_{J,t}^s] \in \mathbb{C}^{M \times J}$ ,  $\mathbf{s}_t^c = [s_{1,t}^c, \dots, s_{K,t}^c]^T \in \mathbb{C}^{K \times 1}$ , and  $\mathbf{s}_t^s = [s_{1,t}^s, \dots, s_{J,t}^s]^T \in \mathbb{C}^{J \times 1}$ . Thus, the equation (3) could be revised as

$$\mathbf{x}_t = \mathbf{W}_t^c \mathbf{s}_t^c + \mathbf{W}_t^s \mathbf{s}_t^s. \quad (4)$$

Note that we assume that the ISAC symbols satisfy  $\mathbb{E}\{\mathbf{s}_t^c \mathbf{s}_t^{cH}\} = \mathbf{I}$  and  $\mathbb{E}\{\mathbf{s}_t^s \mathbf{s}_t^{sH}\} = \mathbf{I}$  while  $\mathbf{s}_t^c$  and  $\mathbf{s}_t^s$  are statistically independent and uncorrelated, i.e.,  $\mathbb{E}\{\mathbf{s}_t^c \mathbf{s}_t^{sH}\} = \mathbf{0}$ . In this paper,  $(\cdot)^T$  and  $(\cdot)^H$  denote transpose and conjugate transpose operations, respectively.  $\mathbf{I}$  and  $\mathbf{0}$  denote the identity matrix and all zero matrix, respectively.

### B. Communication Model

The signal received by the blocked user  $k \in \mathcal{K}_B$  at time slot  $t$  is represented as

$$y_{k,t}^c = \mathbf{h}_{k,t}^c \Phi_t \mathbf{G}_t \mathbf{x}_t + n_{k,t}, \quad (5)$$

where  $\mathbf{h}_{k,t}^c \in \mathbb{C}^{1 \times N}$  is the channel coefficient from the ARIS to the blocked user  $k$  and  $\mathbf{G}_t \in \mathbb{C}^{N \times M}$  denotes the channel of the BS-ARIS link. In addition,  $\Phi_t = \text{diag}(e^{j\theta_{1,t}}, e^{j\theta_{2,t}}, \dots, e^{j\theta_{N,t}}) \in \mathbb{C}^{N \times N}$  is the RIS phase matrix with  $\theta_{i,t} \in [0, 2\pi), \forall i, t$  and  $n_{k,t} \sim \mathcal{CN}(0, \sigma_k^2)$  is the noise.

The signal received by the unblocked user  $k \in \mathcal{K}_U$  at time slot  $t$  is represented as

$$y_{k,t}^{uc} = (\mathbf{h}_{k,t}^{uc} \Phi_t \mathbf{G}_t + \mathbf{d}_{k,t}^{uc}) \mathbf{x}_t + n_{k,t}. \quad (6)$$

where  $\mathbf{h}_{k,t}^{uc} \in \mathbb{C}^{1 \times N}$  is the channel from the ARIS to the unblocked user  $k$  and  $\mathbf{d}_{k,t}^{uc} \in \mathbb{C}^{1 \times M}$  the channel from the BS to the unblocked user  $k$ .

The SINR of the user  $k$  at time slot  $t$  is expressed as (7) on the top of next page. Accordingly, the communication rate of the user  $k$  at slot  $t$  is

$$\Upsilon_{k,t} = \log_2(1 + \gamma_{k,t}^{com}). \quad (8)$$

### C. Sensing Model

With the assumption that the interferences among all received signal are perfectly eliminated, the echo signal from the target  $j$  at time slot  $t$  is denoted as

$$\mathbf{y}_{j,t}^s = \mathbf{V}_t (\mathbf{h}_{j,t}^s \Phi_t \mathbf{G}_t + \mathbf{d}_{j,t}^s)^H \quad (9)$$

where  $\mathbf{h}_{j,t}^s \in \mathbb{C}^{1 \times N}$  is the channel from the ARIS to the target  $j$ ,  $\mathbf{d}_{j,t}^s \in \mathbb{C}^{1 \times M}$  denotes the channel from the BS to the target  $j$ ,  $\mathbf{V}_t \in \mathbb{C}^{M \times M}$  is the receiving beamforming matrix, and  $\mathbf{n}_{j,t} \in \mathbb{C}^{M \times 1}$  is the noise.

Thus, the signal-to-noise ratio (SNR) of target  $j$  at time slot  $t$  is

$$\gamma_{j,t}^{sen} = \frac{|\mathbf{V}_t \mathbf{f}_{j,t}^s H \mathbf{f}_{j,t}^s \mathbf{w}_{j,t}^s|^2}{|\mathbf{V}_t \mathbf{n}_{j,t}|^2}. \quad (10)$$

where  $\mathbf{f}_{j,t}^s = \mathbf{h}_{j,t}^s \Phi_t \mathbf{G}_t + \mathbf{d}_{j,t}^s \in \mathbb{C}^{1 \times M}$ . In order to guarantee the sensing accuracy of each targets, radar SNR is constrained to

$$\gamma_{j,t}^{sen} \geq \gamma_{th}^{sen}, \forall j \in \mathcal{J}, \quad (11)$$

where  $\gamma_{th}^{sen}$  is the threshold of radar SNR.

### D. Problem Formulation

The objective is to maximize the sum of the minimum rate within the whole ISAC period through the joint optimization of the UAV trajectory, RIS phase, and BS transmitting/receiving beamforming vectors under the constraint of ensuring the sensing accuracy of each target. Accordingly, the optimization problem is formulated as (12), where  $P_{\max}^{BS}$  denotes the maximum BS transmitting power and  $e_{\max}^{UAV}$  is the maximum UAV energy consumption. With respect to problem (P1), the optimization variables are highly coupled. Besides, the unimodular constrain in (12f) is intractable. Consequently, problem (P1) is difficult to solve by traditional optimization techniques. As an online DRL method, the proximal policy optimization (PPO) algorithm is adopted

as the optimization technique in this paper.

$$(P1) : \max_{\{\mathbf{w}, \mathbf{q}, \Phi_t, \mathbf{V}_t\}} \sum_{t=1}^T \min_k \Upsilon_{k,t} \quad (12)$$

$$\text{s.t. } \gamma_{j,t}^{sen} \geq \gamma_{th}^{sen}, \forall j \quad (12a)$$

$$\text{Tr}(\mathbf{W}_t^c \mathbf{W}_t^{cH}) + \text{Tr}(\mathbf{W}_t^s \mathbf{W}_t^{sH}) \leq P_{\max}^{BS}, \forall t \quad (12b)$$

$$\sum_{t=1}^T e_t^{UAV}(\mathbf{q}) \leq e_{\max}^{UAV} \quad (12c)$$

$$\frac{\|\mathbf{q}_t - \mathbf{q}_{t-1}\|}{\tau} \leq \nu_{\max}, \forall t \quad (12d)$$

$$\mathbf{q}_t \in \mathcal{D}, \forall t \quad (12e)$$

$$|\Phi_t| = 1, \forall t \quad (12f)$$

$$\|\mathbf{V}_t\| = 1, \forall t, \quad (12g)$$

## III. JOINT OPTIMIZATION VIA PPO-BASED APPROACH

In this section, problem (P1) is represented as the Markov decision process (MDP) and solved by the PPO-based algorithm based on the clipping function.

### A. Problem Reformulation based on MDP

As well known, the goal in a MDP is to find the optimal policy, which is a strategy that prescribes the best action to take in each state to maximize the expected cumulative long-term reward. We first reformulate the primary optimization problem as a MDP. To be specific, the ISAC BS is regarded as the agent with the proposed ARIS-ISAC system being the environment. The agent interacts with the environment to generate a 3-tuple  $\langle \mathcal{S}, \mathcal{A}, \mathcal{R} \rangle$ , where  $\mathcal{S}$ ,  $\mathcal{A}$  and  $\mathcal{R}$  are the set of states, actions and immediate reward, respectively. And the detailed descriptions are as follows.

1) *State Space*  $\mathcal{S}$ : The state  $s_t$  at the slot  $t$  consists of the UAV position and the channel state information. For simplicity, the channel information is expressed as  $\mathbf{H}_t$  which is equivalent to the set  $\{\mathbf{d}_{k,t}^{uc}, \mathbf{h}_{k,t}^{uc}, \mathbf{h}_{k,t}^c, \mathbf{d}_{j,t}^s, \mathbf{h}_{j,t}^s, \mathbf{G}_t\}$ . Since the complex variables are hard to be supported in the neural network, the channel information is expressed as the channel gain  $|\mathbf{H}_t|^2$ . The state  $s_t$  is given by

$$s_t = \{|\mathbf{H}_t|^2, \mathbf{q}_t\}. \quad (13)$$

2) *Action Space*  $\mathcal{A}$ : The agent takes an action  $a_t$  after observing the state  $s_t \in \mathcal{S}$  of the environment. The action set  $a_t$  includes the transmit/receive beamforming vectors, the phase shift matrix and UAV displacement. The phase shift matrix is simplified as the phase angle  $\Theta_t$ . The UAV displacement is denoted by the flight distance  $d_t$  as well as the flight angle  $\theta_t$ . The action  $a_t$  taken at the state  $s_t$  is expressed by

$$a_t = \{\text{Re}\{\mathbf{W}_t\}, \text{Im}\{\mathbf{W}_t\}, \text{Re}\{\mathbf{V}_t\}, \text{Im}\{\mathbf{V}_t\}, \Theta_t, d_t, \theta_t\}. \quad (14)$$

3) *Reward Function*  $\mathcal{R}$ : The reward function is the instant reward. As a result, the reward function  $r_t$  is defined as

$$r_t = \min_k \Upsilon_{k,t}. \quad (15)$$

$$\gamma_{k,t}^{com} = \begin{cases} \frac{\left| \mathbf{h}_{k,t}^c \Phi_t \mathbf{G}_t \mathbf{w}_{k,t}^c \right|^2}{\sum_{j=1}^J \left| \mathbf{h}_{k,t}^c \Phi_t \mathbf{G}_t \mathbf{w}_{j,t}^s \right|^2 + \sum_{i=1, i \neq k}^K \left| \mathbf{h}_{k,t}^c \Phi_t \mathbf{G}_t \mathbf{w}_{i,t}^c \right|^2 + \sigma_k^2} & k \in \mathcal{K}_B \\ \frac{\left| \left( \mathbf{h}_{k,t}^{uc} \Phi_t \mathbf{G}_t + \mathbf{d}_{k,t}^{uc} \right) \mathbf{w}_{k,t}^c \right|^2}{\sum_{j=1}^J \left| \left( \mathbf{h}_{k,t}^{uc} \Phi_t \mathbf{G}_t + \mathbf{d}_{k,t}^{uc} \right) \mathbf{w}_{j,t}^s \right|^2 + \sum_{i=1, i \neq k}^K \left| \left( \mathbf{h}_{k,t}^{uc} \Phi_t \mathbf{G}_t + \mathbf{d}_{k,t}^{uc} \right) \mathbf{w}_{i,t}^c \right|^2 + \sigma_k^2} & k \in \mathcal{K}_U \end{cases} \quad (7)$$

However, the goal of deep reinforcement learning is to find the optimal policy for the sake of maximizing expected discounted reward. Thus, the discount factor  $\gamma$  is introduced and the cumulative discounted reward at the slot  $t$  is expressed as

$$R_t = \sum_{t'=t}^T \gamma^{t'-t} r_{t'}, \quad (16)$$

whose expectation is the value function  $V(s_t) = \mathbb{E}[R_t]$ .

### B. PPO-based Algorithm

Based on the definition of the MDP problem, the framework of the PPO-based scheme is shown in Fig. 2. Different

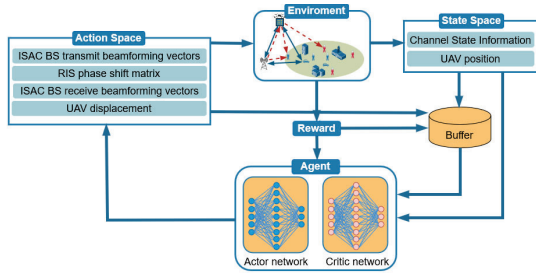


Fig. 2. The structure of the PPO-based algorithm.

from the policy gradient, PPO agent utilizes the old policy  $\pi_{\theta_{old}}$  to interact with the environment to collect samples, while the new policy  $\pi_{\theta}$  is gained from training, which can improve sample efficiency [14]. The objective function is expressed as

$$\mathcal{J}(\theta) = \mathbb{E}_{\pi_{\theta_{old}}} \left[ \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} A^{\theta_{old}}(s_t, a_t) \right]. \quad (17)$$

The advantage function  $A^{\theta_{old}}(s_t, a_t)$  is a crucial metric to quantify whether an action is better or worse than the average action, which is given by

$$A^{\theta_{old}}(s_t, a_t) = r_t + \gamma V^{\theta_{old}}(s_{t+1}) - V^{\theta_{old}}(s_t). \quad (18)$$

To ensure the similarity between new and old policy, the trust region is introduced. Therefore, the objective function with the clipping function is denoted as

$$\mathcal{J}^{clip}(\theta) = \mathbb{E}_{\pi_{\theta_{old}}} \left[ \min \left( \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} A^{\theta_{old}}(s_t, a_t), \text{clip} \left( \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}, 1 - \varepsilon, 1 + \varepsilon \right) A^{\theta_{old}}(s_t, a_t) \right) \right]. \quad (19)$$

where  $\text{clip}(x, y, z)$  means that  $x = y$  when  $x < y$  and  $x = z$  when  $x > z$ . The parameter  $\theta$  of the actor network is updated by

$$\theta = \arg \max_{\theta} \mathcal{J}^{clip}(\theta). \quad (20)$$

The critic network is parameterized as  $\phi$  which is updated by

$$\phi = \arg \min_{\phi} \frac{1}{2} [V^{\theta}(s_t) - R_t]^2. \quad (21)$$

The PPO-based approach is summarized in Algorithm 1.

#### Algorithm 1 The PPO-Based Algorithm

**Input:** Channels  $|\mathbf{H}_t|^2$  and UAV positions  $\mathbf{q}_t$   
**Output:** Actions  $\text{Re}\{\mathbf{W}_t\}$ ,  $\text{Im}\{\mathbf{W}_t\}$ ,  $\text{Re}\{\mathbf{V}_t\}$ ,  $\text{Im}\{\mathbf{V}_t\}$ ,  $\Theta_t$ ,  $d_t$ ,  $\theta_t$

**Initialize:** Initialize  $s_0$ ,  $\theta = \theta_{old}$  and  $\phi$

- 1: **for** each episode **do**
- 2:   Receive initial state  $s_0$
- 3:   **for** step = 1, 2, ...,  $L$  **do**
- 4:     Take the action  $a_t$  by the actor network  $\theta_{old}$
- 5:     Store the transition into the experience pool
- 6:     **if** arrives mini-batch size  $B$  **then**
- 7:       Compute the advantage function
- 8:       **for** epoch = 1, 2, ...,  $m$  **do**
- 9:         Update  $\theta$ ,  $\phi$ ,  $\theta_{old}$
- 10:       **end for**
- 11:     **end if**
- 12:     **if** UAV flies out of the area  $\mathcal{D}$  **then**
- 13:       **break**
- 14:     **end if**
- 15:   **end for**
- 16: **end for**

## IV. SIMULATION RESULTS

We set the simulation parameters and provide comparisons and then present simulations to verify the effectiveness of PPO-based algorithm.

### A. Simulation Scenario and Parameter Setting

In our simulation setup, it is assumed that 4 users (including 3 unblocked users and 1 blocked user) and 2 targets are dispersed. The UAV flies from the initial point (50 m, 30 m) at the fixed height 50 m in the given area  $100 \times 100 \text{ m}^2$  and the BS is at (30 m, 100 m) on the ground. The maximum velocity of the UAV is 10 m/s. There are two hidden layers whose sizes are  $64 \times 64$  in the neural network.



The following four baseline schemes are considered.

**i) Fixed trajectory:** The UAV flies at the fixed speed 6 m/s and flies towards x-axis before 10 steps as well as towards y-axis after 10 steps as shown in Fig. 4.

**ii) Maximum ratio transmission (MRT):** The MRT method precodes the signal in the direction of the channel vector, therefore, the MRT precoding vector to the  $k$ -th user is given by  $\mathbf{w}_{k,t}^c = \frac{\mathbf{f}_{k,t}^H}{\|\mathbf{f}_{k,t}\|}$ , where  $\mathbf{f}_{k,t}$  is the whole channel of the BS-user link and BS-ARIS-user link. It is assumed that the targets are regarded as the users and the MRT precoding vector to the  $j$ -th target is  $\mathbf{w}_{j,t}^s = \frac{\mathbf{f}_{j,t}^H}{\|\mathbf{f}_{j,t}\|}$ , where  $\mathbf{f}_{j,t}$  is the whole channel of the BS-target link and BS-ARIS-target link.

**iii) Random RIS:** The phase shift is chosen from  $(0, 2\pi]$  uniformly.

**iv) Fixed RIS:** A fixed phase shift is adopted.

For simplicity, it is assumed that the channel characteristics are invariant within each time slot and the channel is modeled as Rician channel. According to distance-dependent loss model, the path loss is

$$C = C_0 \left( \frac{d}{d_0} \right)^{-\alpha}, \quad (22)$$

where  $C_0 = -30$  dB is the reference path loss at  $d_0 = 1$  m and  $\alpha$  is the path loss exponent. The channel is modeled as a Rician channel given by

$$\mathbf{H} = \sqrt{C} \left( \sqrt{\frac{\beta}{1+\beta}} \mathbf{H}^{LoS} + \sqrt{\frac{1}{1+\beta}} \mathbf{H}^{NLoS} \right), \quad (23)$$

where  $\mathbf{H}^{LoS}$  and  $\mathbf{H}^{NLoS}$  are LoS component and NLoS components as well as  $\beta$  is the Rician factor. Generally speaking, the NLoS components are modeled as Rayleigh fading and the LoS components are gained from the steering vector. The channel setups are as follows. The path loss exponents and Rician factors of the BS-ARIS, BS-user, BS-target, ARIS-user and ARIS-target link are set as  $\alpha_{BA} = 2$ ,  $\alpha_{BU} = 3$ ,  $\alpha_{BT} = 3$ ,  $\alpha_{AU} = 2.8$ ,  $\alpha_{AT} = 2.8$ ,  $\beta_{BA} = \infty$ ,  $\beta_{BU} = 0$ ,  $\beta_{BT} = 0$ ,  $\beta_{AU} = 10$ , and  $\beta_{AT} = 10$ , respectively. The key simulation parameter settings are shown in I.

### B. Performance of Proposed Algorithm and Baseline

We first plot the smooth reward of different schemes when  $M = 2$ ,  $N = 16$  and  $P_{\max}^{BS} = 10$  W as shown in Fig. 3. It is shown that at the beginning of the training, the reward is little even negative value. This is because some wrong actions make the constrains not satisfied at the beginning of the training, which generates the penalty to the reward. Then it is observed that the proposed algorithm outperforms the other benchmarks. This is due to the fact that the proposed algorithm jointly optimizes the UAV trajectory, RIS phase shifts and beamforming. In addition, all schemes convergence at the stable value.

The UAV trajectory is depicted in Fig. 4 with  $M = 2$ ,  $N = 16$  and  $P_{\max}^{BS} = 10$  W. It can be observed that the UAV trends to fly closer to the blocked user. For the reason

Table I: Simulation Parameter Settings

Parameters	Values
Max training steps	800000
Maximum steps per episode $L$	20
Mini-batch size $B$	20
Discount factor $\gamma$	0.99
Smooth factor $\lambda$	0.95
Clipping factor $\epsilon$	0.2
Learning rate $lr$	$10^{-4}$
Communication noise	-160 dBm
Sensing noise	-140 dBm
Max energy consumption of UAV	39 dB
Sensing threshold	10 dB
Aerodynamics parameters $P_0, P_1, U_{tip}, V_0$ , and $C_0$	80 W, 31.43 W, 120 m/s 40 m/s and 0.0046 kg/m

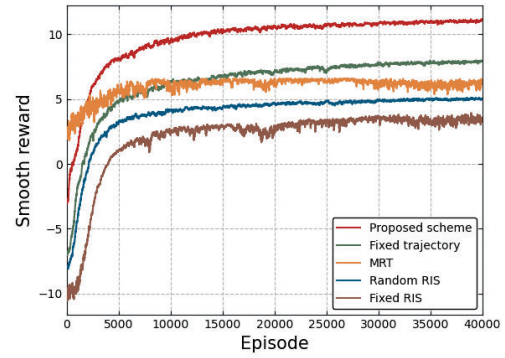


Fig. 3. The convergence of different scheme. (The smooth episode reward is obtained by averaging the episode reward curve using a sliding window of size 60.)

that the achievable rate of the blocked user is generally less than the unblocked user without a direct link. The movement of the UAV guarantees the performance of the worst user to maximize the sum of the minimum achievable rate.

The rate versus the transmitting power is compared as shown in Fig. 5 when  $M = 2$  and  $N = 16$ . It is expected that the rate of all schemes increases with the increasing

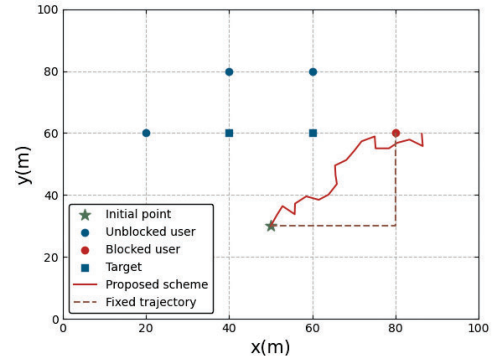


Fig. 4. The 2D UAV-trajectory.

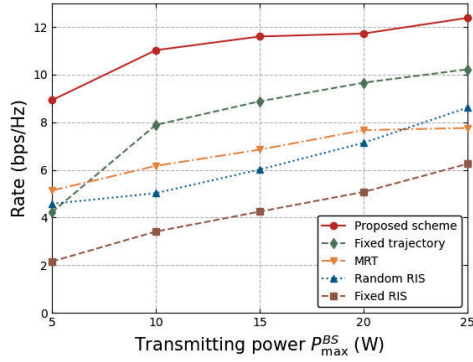


Fig. 5. The rate of the different schemes versus the transmitting power  $P_{\max}^{BS}$ .

power. In addition, the proposed scheme outperforms in rate maximization compared to other benchmark schemes, which shows the effectiveness of the joint optimization of trajectory and beamforming design. Then the fixed-trajectory scheme is observed to achieve higher rate than other scheme without beamforming optimization when the transmission power is large, which indicates its importance.

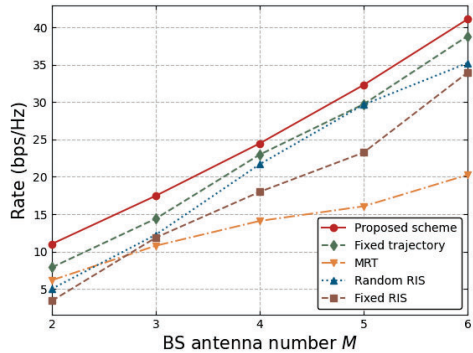


Fig. 6. The rate of the different schemes versus the BS antenna numbers  $M$ .

In Fig. 6, we further compare the performance of each solution in terms of BS antenna numbers when  $N = 16$  and  $P_{\max}^{BS} = 10$  W. It can be observed that the rate of all solutions improves along the increasing BS antenna numbers. This is because more antennas can send signal energy more concentrated to the specific user and gains higher beamforming gain.

## V. CONCLUSIONS

We proposed a novel ARIS-ISAC system, where the ARIS was deployed to serve the communication users and assist radar sensing. A joint optimization problem over beamforming, RIS phase and UAV trajectory was formulated with the aim of maximizing the sum of the minimum achievable

rate during the UAV flight. We introduced the PPO-based algorithm for the problem. Simulation results were presented to demonstrate the effectiveness of the PPO-based algorithm and the ARIS could ensure communication performance of the users via the movement of the UAV.

## ACKNOWLEDGEMENT

This work was supported in part by National Key Research and Development Program of China under Grant 2023YFE0107900, in part by the National Natural Science Foundation of China under Grant 62071306, and in part by the Shenzhen Science and Technology Program under Grants JCYJ20200109113601723, JSGG20210420091805014, and JSGG20210802154203011.

## REFERENCES

- [1] F. Liu, C. Masouros, A. P. Petropulu, H. Griffiths, and L. Hanzo, "Joint radar and communication design: Applications, state-of-the-art, and the road ahead," *IEEE Trans. Wireless Commun.*, vol. 68, no. 6, pp. 3834-3862, Jun. 2020.
- [2] Z. Wei, H. Qu, Y. Wang, X. Yuan, H. Wu, Y. Du, K. Han, N. Zhang, and Z. Feng, "Integrated sensing and communication signals towards 5G-A and 6G: A survey," *IEEE IoT J.*, vol. 10, no. 13, pp. 11068-11092, Jul. 2023.
- [3] H. Zhang, B. Di, K. Bian, Z. Han, H. V. Poor, and L. Song, "Toward ubiquitous sensing and localization with reconfigurable intelligent surfaces," *Proc. IEEE*, vol. 110, no. 9, pp. 1401-1422, Sept. 2022.
- [4] M. Rihan, A. Zappone and S. Buzzi, "Robust RIS-Assisted MIMO Communication-Radar Coexistence: Joint Beamforming and Waveform Design," *IEEE Trans. Commun.*, vol. 71, no. 11, pp. 6647-6661, Nov. 2023.
- [5] X. Song, J. Xu, F. Liu, T. X. Han, and Y. C. Eldar, "Intelligent Reflecting Surface Enabled Sensing: Cramr-Rao Bound Optimization," *IEEE Trans. Signal Process.*, vol. 71, pp. 2011-2026, 2023.
- [6] G. Cheng, Y. Fang, J. Xu, and D. W. K. Ng, "Optimal Coordinated Transmit Beamforming for Networked Integrated Sensing and Communications," *IEEE Trans. Wireless Commun.*, 2024.
- [7] J. Mu, R. Zhang, Y. Cui, N. Gao, and X. Jing, "UAV Meets Integrated Sensing and Communication: Challenges and Future Directions," *IEEE Commun. Mag.*, vol. 61, no. 5, pp. 62-67, May. 2023.
- [8] Z. Lyu, G. Zhu and J. Xu, "Joint Maneuver and Beamforming Design for UAV-Enabled Integrated Sensing and Communication," *IEEE Trans. Wireless Commun.*, vol. 22, no. 4, pp. 2424-2440, Apr. 2023.
- [9] J. Wu, W. Yuan and L. Bai, "On the Interplay Between Sensing and Communications for UAV Trajectory Design," *IEEE IoT J.*, vol. 10, no. 23, pp. 20383-20395, Dec. 2023.
- [10] W. Wei, X. Pang, J. Tang, N. Zhao, X. Wang, and A. Nallanathan, "Secure Transmission Design for Aerial IRS Assisted Wireless Networks," *IEEE Trans. Commun.*, vol. 71, no. 6, pp. 3528-3540, Jun. 2023.
- [11] B. Duo, M. He, Q. Wu, and Z. Zhang, "Joint Dual-UAV Trajectory and RIS Design for ARIS-Assisted Aerial Computing in IoT," *IEEE IoT J.*, vol. 10, no. 22, pp. 19584-19594, Nov. 2023.
- [12] Y. Qin, Z. Zhang, X. Li, W. Huangfu, and H. Zhang, "Deep Reinforcement Learning Based Resource Allocation and Trajectory Planning in Integrated Sensing and Communications UAV Network," *IEEE Trans. Wireless Commun.*, vol. 22, no. 11, pp. 8158-8169, Nov. 2023.
- [13] C. Deng, X. Fang and X. Wang, "Beamforming Design and Trajectory Optimization for UAV-Empowered Adaptable Integrated Sensing and Communication," *IEEE Trans. Wireless Commun.*, vol. 22, no. 11, pp. 8512-8526, Nov. 2023.
- [14] X. Liu, H. Zhang, K. Long, M. Zhou, Y. Li, and H. V. Poor, "Proximal policy optimization-based transmit beamforming and phase-shift design in an IRS-aided ISAC system for the THz band," *IEEE J. Sel. Areas Commun.*, vol. 40, no. 7, pp. 2056-2069, Jul. 2022.