

Intelligent Trajectory Design for RIS-NOMA aided Multi-robot Communications

Xinyu Gao, *Graduate Student Member, IEEE*, Xidong Mu, *Member, IEEE*,
Wenqiang Yi, *Member, IEEE*, and Yuanwei Liu, *Senior Member, IEEE*

Abstract—A novel reconfigurable intelligent surface-aided multi-robot network is proposed, where multiple mobile robots are served by an access point (AP) through non-orthogonal multiple access (NOMA). The goal is to maximize the sum-rate of whole trajectories for the multi-robot system by jointly optimizing trajectories and NOMA decoding orders of robots, phase-shift coefficients of the RIS, and the power allocation of the AP, subject to predicted initial and final positions of robots and the quality of service (QoS) of each robot. To tackle this problem, an integrated machine learning (ML) scheme is proposed, which combines long short-term memory (LSTM)-autoregressive integrated moving average (ARIMA) model and dueling double deep Q-network (D³QN) algorithm. For initial and final position prediction for robots, the LSTM-ARIMA is able to overcome the problem of gradient vanishment of non-stationary and non-linear sequences of data. For jointly determining the phase shift matrix and robots' trajectories, D³QN is invoked for solving the problem of action value overestimation. Based on the proposed scheme, each robot holds an optimal trajectory based on the maximum sum-rate of a whole trajectory, which reveals that robots pursue long-term benefits for whole trajectory design. Numerical results demonstrated that: 1) LSTM-ARIMA model provides high accuracy predicting model; 2) The proposed D³QN algorithm can achieve fast average convergence; and 3) RIS-NOMA networks have superior network performance compared to RIS-aided orthogonal counterparts.

Index Terms—RIS, NOMA, LSTM-ARIMA algorithm, D³QN algorithm, multi-robot system.

I. INTRODUCTION

Nowadays, it is commonly convinced that robots are far less capable when operating independently, and the real power lies in the cooperation of multiple robots. Therefore, multi-robot systems in a shared environment have attracted significant attention in terms of various emerging applications, e.g., cargo delivery, automatic patrol, and emergency rescue [2]. Among these scenarios, robots are required to coordinate with each other to achieve some well-defined goals, e.g., moving from one given position to another. However, with the increasing complexity of the application environment, the large local computational resources are also consumed when collaborating to handle tasks in the multi-robot system. Due to the cooperation requirement and high computation complexity for the trajectory design in multi-robot systems, wireless communications with advanced multiple access techniques are important for multi-robot systems [3].

As an emerging technique, non-orthogonal multiple access (NOMA) [4]–[6] adopts a flexible successive interference cancellation (SIC) receiver for robust multiple access. It improves spectrum efficiency by opportunistically exploring users' channel conditions. However, there still exists a shortage of spectrum in some communication regions, e.g., blind zone. To tackle this problem, reconfigurable reflecting surfaces (RIS) [7]–[12] is a potential candidate to improve the spectrum efficiency, which is passive equipment that can proactively reflect the signal to the users. Specifically, the employment of the RIS is able to create a virtual line of sight (LOS) between access points (AP) and robots when robots are located in the communication blind zone. In view of the advantages brought by the RIS and NOMA techniques, RIS-NOMA is regarded as a potential solution to efficiently handle the trajectory design problems of the multi-robot system.

A. Related Works

1) *RIS-NOMA Networks*: RIS-NOMA technique becomes appealing in recent years, and it has been applied in various scenarios, including spectrum efficiency and user connectivity improvement [13]–[17], energy consumption decrease [18], etc. In [19], alternating successive convex approximation (SCA) and semi-definite relaxation (SDR) based energy-efficient algorithms were proposed to yield a good tradeoff between the sum-rate maximization and total power consumption minimization on RIS-aided NOMA networks, by maximizing the system energy efficiency by jointly optimizing the transmit beamforming at the base station (BS) and the reflecting beamforming at the RIS. To establish stable and high-quality communication links between AP and the robotic users, an indoor robot navigation system was investigated in [20], where RIS was employed to enhance the connectivity and NOMA was adopted to improve the communication efficiency between the AP and robotic users. An SDR-based solution in [21] was proposed to address a joint power control at the users and beamforming design at the IRS for maximizing the sum-rate of all users in a RIS-aided uplink NOMA system. Also, the impact of the number of reflecting elements on the sum-rate was revealed. The authors in [22] proposed a novel framework of resource allocation in multi-cell RIS-aided NOMA networks, which was capable of being enhanced with the aid of the RIS, and the proper location of the RIS is also guarantee the trade-off between spectrum and energy efficiency. In order to examine the effectiveness of RIS in the NOMA system with respect to transmitting power consumption, the authors in [23] explored the relationship

Part of this work has been presented at the IEEE International Conference on Communications, 14–23 June, 2021 [1].

X. Gao, X. Mu, W. Yi, and Y. Liu are with the School of Electronic Engineering and Computer Science, Queen Mary University of London, London E1 4NS, U.K. (e-mail: {x.gao, xidong.mu, w.yi, yuanwei.liu}@qmul.ac.uk).

between individual user's transmit power and the variables of phase shifts, and solved the phase shift determination problem via the sequential rotation algorithm.

2) *Communication-connected Multi-robot Trajectory Design*: Leveraging the wireless communication technology is capable of providing services for the communication-connected multi-robot systems that require communication quality, so as to provide optimal trajectory design in special environments [24]. To avoid path conflicts between robots and repeated exploration and information collected from the same region by different robots, a novel algorithm was proposed [25]. Among them, decisions to select locations for exploration and information collection were guided by a utility function that combines Gaussian Process-based distributions for information entropy and communication signal strength, along with a distributed coordination protocol. The authors of [26] formulated the problem of multi-robot informative path planning under continuous connectivity constraints as an integer program leveraging the ideas of bipartite graph matching and minimal node separators. In [27], a framework for planning and perception for multi-robot exploration in large and unstructured three-dimensional environments was presented. A Gaussian mixture model for global mapping was employed in the proposed method to model complex environment geometries while maintaining a small memory footprint which enables distributed operation with a low volume of communication. In [28], the authors developed a communication-based navigation framework for multi-robot systems in unknown areas that solely exploit the sensing information and shared data among the agents. Additionally, the millimeter-wave [29] and multi-input multi-output (MIMO) [30] technologies were integrated in the multi-robot trajectory design. In contrast to conventional multi-robot path planning, the authors of [31] defined a type of multi-robot association-path planning problem aiming to jointly optimize the robots' paths and the robots-AP associations. Using geometrically motivated assumptions and characteristics of MIMO, the authors in [32] derived transmitter spacing rules that can be easily added to existing path plans to improve backhaul throughput for data offloading from the robot team, with minimal impact on other system objectives.

B. Motivations and Contributions

Although mentioned studies have revealed the potential abilities of communication technology in multi-robot systems, the research on the communication-connected multi-robot long-term trajectory design is in its early stage. For instance, the previous research mainly focuses on the trajectory design constrained by the geographic environment, ignoring the important role of communication quality in multi-robot trajectory design. Since the ultimate goal for the multi-robot system is to achieve specific goals such as simultaneous localization and mapping or moving to some target areas with minimum time/energy consumption [33], [34], the achievable maximum performance (i.e., capacity) of the whole system is a basis for the multi-robot system to efficiently achieve these goals. Additionally, the exploration of the long-term benefits of the dynamic movement of mobile robots in networks is also

neglected in previous research contributions. The limitations and challenges are summarized as follows:

- **The characterization of the channel model**: The signal may be blocked in some areas by the obstacles, while the channel model also changes abruptly. As result, the position-dependent channel model indicates that the channel model characterization is challenging.
- **The determination of initial and final positions**: In some special scenarios (e.g., fixed-point cruise), different initial and final positions bring various trajectory designs and obtain miscellaneous achievable sum-rate. The trajectory with a high sum-rate is able to make the robot exchange much environment information with the controller and further effectively promote subsequent operation. Therefore, it's important to find the optimal initial and final positions corresponding to the trajectories with maximum sum-rate.
- **The influence of the resource allocation strategy**: To realize optimal trajectories, the robots need to continuously evaluate the previous network performance and resource allocation strategy at each timestep to reap the rewards and determine the next action. Therefore, the robots continuously interacting with the environment is also challenging.
- **The design of the ML-based algorithm**: Since the positions of robots are determined by the previous positions and environments, conventional algorithms fail to handle a time-varying-based MDP, and vanilla ML algorithms have limitations to achieve good performance. Thus, compared to these algorithms, the vanilla ML algorithms need to be further improved for obtaining better performance than the vanilla algorithms.

In response to the above limitations and challenges, we proposed a novel ML framework, which integrated Long short-term memory (LSTM)-autoregressive integrated moving average (ARIMA) and dueling double deep Q-network (D³QN) algorithms to give full play to the conventional ML algorithms. LSTM-ARIMA algorithm is able to efficiently handle non-stationary [35] and non-linear sequences of any size of data and completely solve the problem of gradient vanishment¹. D³QN algorithm has the capability of solving the problem of action value overestimation in double DQN and dueling DQN algorithms and improving the accuracy of action selection for each state. After providing the range of the initial and final positions by the LSTM-ARIMA algorithm, the D³QN algorithm is able to characterize the channel models and further provide the optimal initial and final positions, trajectories, and phase shift design for the robots. The main contributions of this paper are summarized as follows:

- We propose a new framework for RIS-NOMA-aided multi-robot networks. The RIS is employed to enhance communication efficiency, by proactively reflecting the incident signals, while NOMA is invoked for improving

¹For LSTM, the memory neural networks are introduced to train a model, where the non-stationary data with dramatic fluctuation can be handled. For ARIMA, it is to utilize the difference operation to convert the non-stationary sequences to stationary sequences first and then handle the stationary sequences based on the AR model.

the spectrum efficiency of the multi-robot system. Based on the framework, an optimization problem is formulated to obtain the maximum sum-rate of the whole trajectories for all robots. by jointly optimizing trajectories for robots, reflecting coefficient matrix of RIS, successive interference cancelation (SIC) decoding order for NOMA, power allocation at the AP, subject to the quality of service (QoS) for the robots.

- We adopt the LSTM-ARIMA algorithm for training a more accurate model to guide the optimal trajectory design for the robots. LSTM-ARIMA algorithm is able to handle non-stationary and non-linear sequences of any size of data, as well as completely solve the problem of gradient vanishment. For channels, at each point/time, the experimental observations of small-scale fading of channels are different, meaning they have different sequence properties (e.g., mean and variance). Hence, the observations at different points/times are non-stationary. The LSTM-ARIMA model is able to unify small-scale fading properties by difference operation to further mimic the small-scale fading distribution in the considered environment.
- We demonstrate that the proposed D³QN algorithm is capable of providing the optimized robots' trajectories, phase shifts of RIS, and optimal initial and final positions for the robots, according to the obtained results of the LSTM-ARIMA algorithm. The D³QN algorithm is an online reinforcement learning (RL) algorithm to train an off-policy, which fully reaps the merits of double DQN and dueling DQN to introduce the target network and split the network structure to obtain the accurate Q value. After numerous repetitive training, the performance of the D³QN algorithm is proved to outperform double DQN and dueling DQN algorithms.
- We demonstrate that the proposed D³QN algorithm efficiently solves the trajectory design problem of the multi-robot system. Specifically, the D³QN algorithm is able to achieve the fastest convergence speed than conventional double DQN and dueling DQN algorithms. Furthermore, when applying the D³QN, dueling DQN and double DQN algorithms on trajectory design upon different elements of RIS, the proposed D³QN algorithm is able to find a shorter total path than the conventional dueling DQN and double DQN algorithms. Compared to the RIS-OMA technique, the RIS-NOMA technique is able to achieve a higher communication sum-rate with a shorter traveling distance for robots, which verifies the effectiveness of the NOMA technique. The optimal decoding order has better performance than the fixed decoding order and random decoding order, which also shows that the decoding order is a considerable factor in NOMA networks.

C. Organizations

The rest of this paper is organized as follows. Section II presents the system model for the considered RIS-aided multi-robot NOMA networks, and the passive beamforming and trajectory design problems are formulated. In Section III, we propose the LSTM-ARIMA model and D³QN algo-

rithm, which is employed to predict the initial position and the final position, jointly planning trajectories and designing the beamforming. Section IV presents numerical results to verify the effectiveness of the proposed machine learning-based optimization algorithms for joint trajectories planning and passive beamforming design, as well as the performance of the algorithms. Finally, Section V concludes this paper.

Notations: Scalars, vectors, and matrices are denoted by lower-case, bold-face lower-case, and bold-face upper-case letters, respectively. $\mathbb{C}^{K \times N}$ denotes the space of $K \times N$ complex-valued vectors. The conjugate transpose of vector \mathbf{a} is denoted by \mathbf{a}^H . $\text{diag}(\mathbf{a})$ denotes a diagonal matrix with the elements of vector \mathbf{a} on the main diagonal. $|\mathbf{a}|$ denotes the norm of vector \mathbf{a} . $*$ denotes the dot multiplication operation. $\log_2(\mathbf{A})$ represent a logarithmic function with a constant base of 2 for matrix \mathbf{A} .

II. SYSTEM MODEL AND PROBLEM FORMULATION

A. System Model

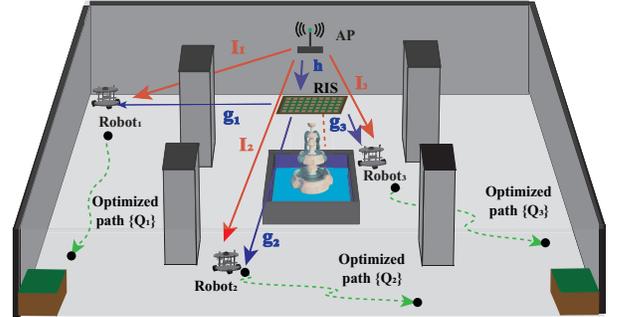


Fig. 1: Illustration of the RIS aided multi-robots cruise system for the indoor environment.

As shown in Fig. 1, we focus our attention on a downlink RIS-aided multi-robot NOMA network, where one single-antenna AP serves \mathcal{X} single-antenna mobile wheeled robots, assisting by a RIS with K passive reflecting elements. The passive reflecting elements in the RIS can be partitioned into M sub-surfaces, while each sub-surfaces consists of $\bar{K} = K/M$ elements. Assume that the two-dimensional (2D) horizontal ground space for robots moving is approximately smooth without undulation, indicating that the height h_r of the robot (the height of the antenna) is regarded as a constant value. The other values' (e.g., velocity, turning radius) adjustment involved in the robot's mechanical motion are also considered to be accurate. Additionally, the receiver is located on the top of the robots. Define a three-dimensional (3D) Cartesian coordinate system, the positions of AP and robot i are denoted as (x_A, y_A, h_A) and (x_i, y_i, h_r) , respectively. Note that for users' fairness, the RIS is arranged in the center of the ceiling in the motion space, where the position can be denoted as (x_I, y_I, h_I) .

In view of the deployment of the RIS, the composite received signal of each robot is the combination of two components, which are derived from the AP-robot direct link, and the AP-RIS-robot reflecting link. Denote the baseband equivalent channels from the AP to i -th robot, the AP to RIS and the RIS to i -th robot as \mathbf{l}_i , $\mathbf{h}^H \in \mathbb{C}^{1 \times M}$, and $\mathbf{g}_i \in \mathbb{C}^{M \times 1}$, $i = \{1, 2, \dots, \mathcal{X}\}$, respectively. For the distance-reference path

loss, it can be modeled as $L_u = Cd_u^{-\gamma}$, $u = \{Ai, Ri, AI\}$ [36], where C denotes the path loss at the reference distance of 1 meter, $d_{[\cdot]}$ is distance for individual link, γ represents the path loss factor, respectively. Since the positions of robots are time-sensitive, the positions of i -th robot can be re-denoted as $q_i(t) = (x_i(t), y_i(t), h_r)$. For small-scale fading, we assume a Rician fading channel models for the BS to i -th robot, the RIS to i -th robot channels and the AP to RIS channel. Thus, the three individual channels for AP to i -th robot, RIS to i -th robot and AP to RIS can be expressed as

$$\mathbf{l}_i(q_i(t)) = L_{Ai}(q_i(t)) \left\{ \sqrt{\frac{1}{\alpha_{Ai}(q_i(t)) + 1}} \cdot [\sqrt{\alpha_{Ai}(q_i(t))} \mathbf{l}_i^{\text{LOS}}(q_i(t)) + \hat{\mathbf{l}}_i^{\text{NLOS}}] \right\}, \quad (1)$$

$$\mathbf{g}_i(q_i(t)) = L_{Ri}(q_i(t)) \left\{ \sqrt{\frac{1}{\alpha_{Ri}(q_i(t)) + 1}} \cdot [\sqrt{\alpha_{Ri}(q_i(t))} \tilde{\mathbf{g}}_i^{\text{LOS}}(q_i(t)) + \hat{\mathbf{g}}_i^{\text{NLOS}}] \right\}, \quad (2)$$

$$\mathbf{h} = L_{AI} \left\{ \sqrt{\frac{1}{\alpha_{AI} + 1}} [\sqrt{\alpha_{AI}} \tilde{\mathbf{h}}^{\text{LOS}} + \hat{\mathbf{h}}^{\text{NLOS}}] \right\}, \quad (3)$$

where $\alpha_u, u = \{Ai, Ri, AI\}$, $z^{\text{LOS}}, z = \{\tilde{\mathbf{l}}_i, \tilde{\mathbf{g}}_i, \tilde{\mathbf{h}}\}$, $w^{\text{NLOS}}, w = \{\hat{\mathbf{l}}_i, \hat{\mathbf{g}}_i, \hat{\mathbf{h}}\}$ denote the Rician factor, deterministic LOS component and random non-line-of-sight (NLOS) Rayleigh fading components, respectively. Here, we are interested in the expected values to obtain the NLOS component. When $\alpha_{Ai}(q_i(t)) = 0$ or $\alpha_{Ri}(q_i(t)) = 0$, the AP to i -th robot link or the RIS to i -th robot link is blocked. Otherwise, $\alpha_{Ai}(q_i(t)) = \bar{\alpha}$ or $\alpha_{Ri}(q_i(t)) = \bar{\alpha}$, where $\bar{\alpha}$ is a constant value. With respect to RIS, denote $\Phi(t) = \text{diag}(\phi_1(t), \phi_2(t), \dots, \phi_M(t))$ as the reflection coefficients matrix of the RIS, where $\phi_m(t) = \beta_m(t)e^{j\theta_m(t)}$, $k = \{1, 2, 3, \dots, M\}$ denotes the reflection coefficient of m -th sub-surface of the RIS. Among them, the $|\beta_m(t)| = 1$ and $\theta_m(t) \in [0, 2\pi)$ denote the amplitude and phase of m -th sub-surface in the RIS. Thus, the effective channel from the AP to the robot i is given by

$$\mathbf{H}_i(q_i(t)) = \mathbf{h}^H \Phi(t) \mathbf{g}_i(q_i(t)) + \mathbf{l}_i(q_i(t)), i = \{1, 2, \dots, \mathcal{X}\}. \quad (4)$$

According to the fairness principle, the interference among robots should be considered while one AP serves \mathcal{X} robots simultaneously. We consider the NOMA technology to mitigate interference among robots and for sharing the same time/frequency resources to all the robots. In the NOMA technology, the superposition coding (SC) method is applied at the AP. Let $S_i = \sqrt{p_i} s_i$ denote the transmitted signal for the robot i , $i = \{1, 2, \dots, \mathcal{X}\}$. s_i represents the transmitted information symbol for the robot i . It is worth noting that S_i is satisfied $\mathbb{E}[|S_i|^2] = p_i \leq P_i, i = \{1, 2, \dots, \mathcal{X}\}$, with p_i and P_i denoting the transmitted power and its maximum value of the robot i , respectively. The successive interference cancelation (SIC) method is applied for each robot to remove the interference. The robots with stronger channel power gain decode signals of other robots with weaker channel power gain priority over decoding their own signal. The larger power will be allocated to the weak user first while the smaller power will allocate to the strong user. Denote $O(i)$ as the

decoding order of the robot i . For any two robots i and j , $i \neq j$, $i, j = \{1, 2, \dots, \mathcal{X}\}$, if the decoding order satisfying $O(i) < O(j)$, the received signal of robot i can be modeled as

$$Y_i(q_i(t)) = \mathbf{H}_i(q_i(t)) S_i(t) + \sum_{O(j) > O(i)} \mathbf{H}_i(q_i(t)) S_j(t) + n, \quad (5)$$

where the $n \sim \mathcal{CN}(0, \sigma^2)$ denotes the additive white Gaussian noise (AWGN) with average power σ^2 . For each robot i , $i = \{1, 2, \dots, \mathcal{X}\}$, the achievable rate can be denoted as R_i . To simplify the problem, denote the equivalent channel of AP-RIS-robot link $\mathbf{h}^H \Phi(t) \mathbf{g}_i(q_i(t)) = (v(t))^H \psi(q_i(t))$, where $\psi(q_i(t)) = \text{diag}\{\mathbf{h}^H\} \mathbf{g}_i(q_i(t))$, $v(t) = [v_1(t), v_2(t), \dots, v_M(t)]^H$, and $v_m = e^{j\theta_m}$. so the signal-to-interference-plus-noise ratio (SINR) of robot i is given by

$$\tau_i(q_i(t)) = \frac{|(v(t))^H \psi(q_i(t)) + \mathbf{l}_i(q_i(t))|^2 p_i(t)}{\sum_{O(j) > O(i)} |(v(t))^H \psi(q_i(t)) + \mathbf{l}_i(q_i(t))|^2 p_j(t) + \sigma^2}, \quad (6)$$

where the σ^2 denotes the variance of the AWGN. Then, according to $R = \log_2(1 + \text{SINR})$, the achievable communication rate at robot i can be expressed as

$$R_i(q_i(t)) = \log_2 \left(1 + \frac{|(v(t))^H \psi(q_i(t)) + \mathbf{l}_i(q_i(t))|^2 p_i(t)}{\sum_{O(j) > O(i)} |(v(t))^H \psi(q_i(t)) + \mathbf{l}_i(q_i(t))|^2 p_j(t) + \sigma^2} \right). \quad (7)$$

We consider perfect SIC in our system model, to guarantee rate fairness between two robots, the conditions $R_{i \rightarrow j} \geq R_{i \rightarrow i}$ should be satisfied under given decoding order $O(i) < O(j)$. The $R_{i \rightarrow j}$ and $R_{i \rightarrow i}$ can be expressed as follows:

$$R_{i \rightarrow j}(q_i(t)) = \log_2 \left(1 + \frac{|(v(t))^H \psi(q_j(t)) + \mathbf{l}_j(q_j(t))|^2 p_i(t)}{\sum_{O(k) > O(i)} |(v(t))^H \psi(q_j(t)) + \mathbf{l}_j(q_j(t))|^2 p_k(t) + \sigma^2} \right), \quad (8)$$

$$R_{i \rightarrow i}(q_i(t)) = \log_2 \left(1 + \frac{|(v(t))^H \psi(q_i(t)) + \mathbf{l}_i(q_i(t))|^2 p_i(t)}{\sum_{O(k) > O(i)} |(v(t))^H \psi(q_i(t)) + \mathbf{l}_i(q_i(t))|^2 p_k(t) + \sigma^2} \right). \quad (9)$$

B. Problem Formulation

Our goal is to maximize the sum-rate of all robot trajectories by jointly optimizing trajectories for robots, reflection coefficient matrix of RIS, SIC decoding order for NOMA, power allocation at the AP, subject to the QoS for each robot. Hence, the optimization problem is formulated as

$$\max_{v, \Omega, \{p_i\}, \varepsilon_i, \xi_i, \mathbf{Q}} \sum_{i=1}^{\mathcal{X}} R_i(\mathbf{Q}_i) \quad (10)$$

$$\text{s.t. } R_i(q_i(t)) \geq \bar{R}, \quad \forall t \in [0, T], \quad (10a)$$

$$|v_m(t)| = 1, \quad \forall m \in \{1, 2, \dots, M\}, \\ \forall t \in [0, T], \quad (10b)$$

$$\Omega \in \mathbf{\Pi}, \quad (10c)$$

$$\sum_{i=1}^{\mathcal{X}} p_i(t) \leq \mathcal{P}, \quad (10d)$$

$$q_i(0) = \varepsilon_i, q_i(T) = \xi_i \quad (10e)$$

$$|\dot{q}_i(t)| = V, \quad \forall t \in [0, T], \quad (10f)$$

$$q_i(t) \in \mathbf{Q}_i, \quad \forall t \in [0, T], \quad (10g)$$

where the \bar{R} , ε_i , ξ_i , $\mathbf{\Pi}$, and $\mathbf{Q} = [\mathbf{Q}_1, \mathbf{Q}_2, \dots, \mathbf{Q}_X]$ denote the minimal required communication rate of all the robots, the possible initial positions for the robots, the possible final positions for the robots, the set of all the possible decoding orders, and the set of trajectories for all the robots, respectively. Constraint (10a) and constraint (10b) are the QoS requirements for robot i and the restraint for the RIS reflection coefficients. Constraint (10c) is the decoding conditions for the NOMA scheme. Constraint (10d) is the total transmission power constraint. Constraint (10e) and constraint (10f) characterize the initial position and final position for each robot and constant velocity for all robots. However, the main difficulty in solving the problem (10) can be summarized in the following reasons. Firstly, the individual instantaneous channel is modeled according to the position of the robot, while it is highly coupled with transmit power allocation and RIS phase shift. Hence, these variables have to be optimized simultaneously instead of decoupling these variables and optimizing them successively. Secondly, the optimal robot trajectories need to be explored by choosing different initial and final positions, while the trajectory for each robot is not always the shortest path from ε_i to ξ_i , which constrained by the simultaneous arriving conditions for all the robots. Thirdly, the current positions of robots are determined by the previous positions and the environment, which follows the MDP. Fourthly, the experimental observations of small-scale fading of channels are different, the properties of small-scale fading should be stable to guide the trajectories design. The conventional optimization algorithms may fail to be promoted to MDP scenarios. However, our goal is to provide a policy for the whole system, which pursues a long-term effect. Therefore, conventional optimization methods are not proper to be employed to solve these difficulties. The machine learning methods can be invoked to optimize studied RIS-assisted NOMA networks and obtain the optimized trajectories. Additionally, for the implementation of optimal SIC decoding order, the strong users decode the weak users' signals, while the weak users decode the received signal directly without SIC. Therefore, the optimization problem with SIC decoding order constraint becomes a combinatorial optimization problem, where the optimal solution is difficult to obtain by employing the ML algorithm in the current stage [37]. Therefore, in this paper, we apply the exhaustive search method [38] to explore the optimal SIC decoding order.

III. PROPOSED SOLUTIONS

In this section, we propose an ML-based scheme to solve the problem (10), as shown in Fig. 2. To characterize the non-stationary property in practical robotic communications, we model the sequence property of small-scale fading at different times/locations with the same Rician fading but different parameters in this work as a case study. It is worth mentioning that other non-stationary models can be solved via the same method. Regarding channel statistic change, we assume it changes every coherence time and the decision will be made within each coherence time. Additionally, assuming that the initial position and the final position have a one-to-one

correspondence, we also aim to explore the optimal initial-final pair at one time slot using ML algorithm. Secondly, in order to find the optimal trajectories and initial-final positions pair for robots, an improved ML-based algorithm, i.e., D³QN, is invoked for jointly planning trajectories and designing the beamforming.

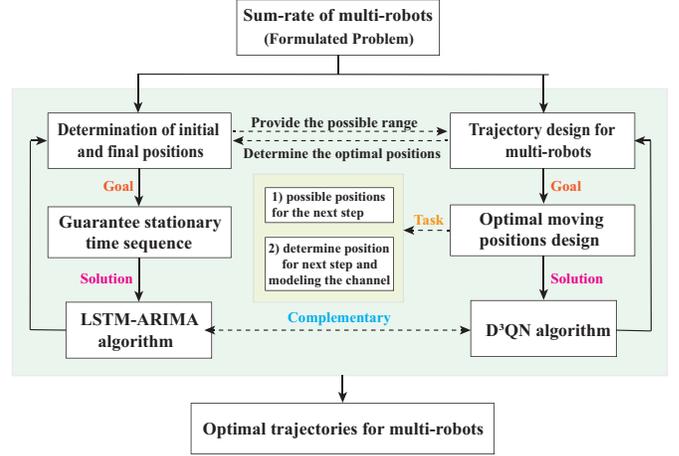


Fig. 2: Schematic for solving formulated problem.

A. LSTM-ARIMA Model for Prediction

As a variant of recurrent neural network (RNN), the LSTM network has an excellent performance in efficiently handling non-stationary and non-linear sequences of data. But for long sequences, LSTM cannot completely solve the problem of gradient vanishment. The ARIMA model is not puzzled by this problem and provides a valid solution for a linear sequence of data. However, it is a time prediction model which essentially captures linear relationships, while nonlinear relationships cannot be involved. Therefore, we consider combining the advantages of the two models and propose a novel LSTM-ARIMA model for prediction.

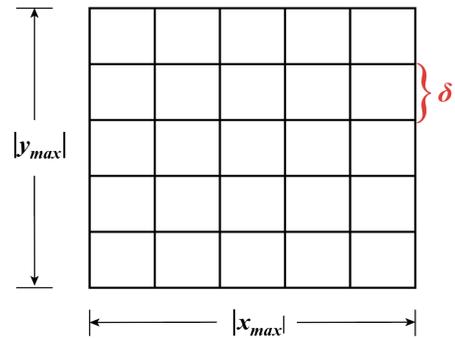


Fig. 3: Discretization of the geographic map.

1) *Data Pre-processing*: In order to characterize the geographic map of the environment, as shown in Fig. 3, the map is discretized into $(x_{max}y_{max})/\delta^2$ cells, where x_{max} , y_{max} and δ denote maximum bound of x-axis, maximum bound of y-axis and resolution of the map. To simplify the problem, we assume that δ is small enough which makes the size of cells can be approximated to the center of cells. Randomly generate [39] the position pairs $\mathbf{S}_o = \{S_1, \dots, S_{4N}, S_{4N+1}, \dots, S_{5N}, S_{5N+1}, \dots, S_{6N}\}$ as the historical positions data set, where $\mathbf{S}_{train} =$

$\{S_1, S_2, \dots, S_{4N}\}$ is selected as the training samples and $S_{test} = \{S_{4N+1}, S_{4N+2}, \dots, S_{5N}\}$ as the test samples for LSTM and ARIMA models. In order to reduce the influence of value range on network performance, it is necessary to normalize the values of training samples to $[0, 1]$. By invoking Min-Max normalization method [40], the normalized samples can be expressed as

$$\tilde{S}_{\bar{N}, \bar{n}} = \frac{S_{\bar{N}, \bar{n}} - S_{min}}{S_{max} - S_{min}}, \quad \bar{n} = 1, 2, \quad \bar{N} \in \{1, 2, \dots, 5N\} \quad (11)$$

where $\tilde{S}_{\bar{N}, \bar{n}}$, $S_{\bar{N}, \bar{n}}$, S_{max} and S_{min} are the normalized values of positions pairs, the original values of positions pairs, the maximum values of all positions pairs, and the minimum values of all positions pairs, respectively. Note that, $\bar{n} = 1$ and $\bar{n} = 2$ represent the initial position and final position, while S_{max} and S_{min} are determined by all initial and final positions.

2) *Prediction Based on LSTM Model*: We take the normalized data $\tilde{\mathbf{S}} = \{\tilde{S}_1, \tilde{S}_2, \dots, \tilde{S}_{4N}\}$ and denote $\bar{\mathbf{S}}^{\text{LSTM}} = \{\bar{S}_1^{\text{LSTM}}, \bar{S}_2^{\text{LSTM}}, \dots, \bar{S}_N^{\text{LSTM}}\}$ as the input position pairs and predicted position pairs. According to the LSTM model, denote ν_i, ν_f, ν_o, H , and U as input gate, forget gate, output gate, size of input layer, and size of hidden memory, respectively. Thus, at each timestep t , the definition of those parameters can be calculated as

$$\nu_i^t = g\left(\sum_h \omega_{i,h} \tilde{\mathbf{S}}_i^t + \sum_u \omega_{i,u} \mathbf{c}_i^{t-1} \mathbf{b}_u^{t-1}\right), \quad (12)$$

$$\nu_f^t = g\left(\sum_h \omega_{f,h} \tilde{\mathbf{S}}_f^t + \sum_u \omega_{f,u} \mathbf{c}_f^{t-1} \mathbf{b}_u^{t-1}\right), \quad (13)$$

$$\nu_o^t = g\left(\sum_h \omega_{o,h} \tilde{\mathbf{S}}_o^t + \sum_u \omega_{o,u} \mathbf{c}_o^{t-1} \mathbf{b}_u^{t-1}\right), \quad (14)$$

where $g(\cdot)$, \mathbf{c}^{t-1} and $\mathbf{b}_u^t = \nu_f^t \mathbf{b}_u^{t-1} + \nu_i^t g(\sum_h \omega_{i,u} \tilde{\mathbf{S}}_i^t + \sum_u \omega_{i,u} \mathbf{c}_i^{t-1})$ denote the activation function, output of each hidden layer, state of each layer, respectively. Then, the predicted denormalized output $\bar{\mathbf{S}}^{\text{LSTM}}$ of the model can be calculated by

$$\bar{S}_n^{\text{LSTM}} = \nu_o \tanh(\mathbf{b}_u^{t-1}) * (S_{max} - S_{min}) + S_{min}, \quad n \in \{1, 2, \dots, N\}. \quad (15)$$

To prevent LSTM training from overfitting, the adjustment of weight on the output fully connected layer is increased by L2 normalization regularization [41].

3) *Prediction Based on ARIMA Model*: ARIMA is an effective statistical model for time series analysis and forecasting, which is integrated of auto regression (AR) model, moving average (MA) model and difference model. The AR model and MA model can be expressed as

$$\bar{S}_n^{\text{AR}} = \mu_n + \sum_{z=1}^a \kappa_z \tilde{S}_{n-z}, \quad n \in \{1, 2, \dots, N\}, \quad (16)$$

$$\bar{S}_n^{\text{MA}} = \zeta_n + \sum_{z=1}^b \varphi_z \zeta_{n-z} = \sum_{z=1}^b \varphi_z L \zeta_n, \quad n \in \{1, 2, \dots, N\}, \quad (17)$$

where $\mu_n, a, \kappa_z, \zeta_n, b, \varphi_z$, and L denote the constant value, order of autoregressive model, autoregressive coefficient, model errors, order of the moving average model, moving average

coefficient, and lag operator, respectively. Then, in order to guarantee the stability of input data, a difference operator Δ^d is employed, which can be given by

$$\Delta^d \bar{S}_n = (1 - L)^d (\bar{S}_n^{\text{AR}} + \bar{S}_n^{\text{MA}}), \quad n \in \{1, 2, \dots, N\}. \quad (18)$$

Thus, the the predicted denormalized output $\bar{\mathbf{S}}^{\text{ARIMA}}$ of ARIMA model can be expressed as

$$\bar{S}_n^{\text{ARIMA}} = \left[\sum_{z=1}^a \kappa_z (1 - L)^d (\bar{S}_n^{\text{AR}} + \bar{S}_n^{\text{MA}}) + \sum_{z=1}^b \varphi_z \zeta_{n-z} + \mu_n + \zeta_n \right] * (S_{max} - S_{min}) + S_{min}, \quad n \in \{1, 2, \dots, N\}. \quad (19)$$

4) *Fusion Based on Critic Weight Method*: In order to further improve the performance of the model, the critic weight method [42] is exploited to assign weights to the predicted values of the two models. An evaluate indicator \mathcal{I} is proposed to calculate the RMSE of the two models, which can be expressed as

$$\mathcal{I} = \sqrt{\frac{1}{N} \sum_{n=1}^N (|\bar{S}_n - S_{4N+n}|)}. \quad (20)$$

When \mathcal{I} belows the threshold ς , the trained model can be accepted, the predicted data $\bar{\mathbf{S}} = \{\bar{S}_1, \bar{S}_2, \dots, \bar{S}_N\}$ by LSTM model and ARIMA model are obtained by inputing $\mathbf{S}_{ini} = \{S_{5N+1}, S_{5N+2}, \dots, S_{6N}\}$. Accordingly, the predicted data $\bar{\mathbf{S}}$ after model fusion can be calculated as

$$\bar{S}_n = \bar{w}_n \bar{S}_n^{\text{LSTM}} + \tilde{w}_n \bar{S}_n^{\text{ARIMA}}, \quad n \in \{1, 2, \dots, N\}. \quad (21)$$

When weights \bar{w}_n and \tilde{w}_n are determined, the channel properties set \mathbf{Z} , initial position ε_i and final position ξ_i can be opted from the prediction pairs $\bar{\mathbf{S}}_{ini}$. The detailed pseudo code is shown in **Algorithm 1**. The problem (10) can be reformulated as

$$\max_{v^d, \Omega, \{p_i\}, \varepsilon_i, \xi_i, \mathbf{Q}^d} \sum_{i=1}^{\mathcal{X}} R_i(\mathbf{Q}_i^d) \quad (22)$$

$$\text{s.t.} \quad R_i(q_i^{(n)}) \geq \bar{R}, \quad \forall n \in \{0, 1, \dots, N\}, \quad (22a)$$

$$|v_m^{(n)}| = 1, \quad \forall m \in \{1, 2, \dots, M\}, \quad \forall n \in \{0, 1, \dots, N\}, \quad (22b)$$

$$q_i^{(0)} = \varepsilon_i, q_i^{(N)} = \xi_i, \quad \forall \varepsilon_i, \xi_i \in \bar{\mathbf{S}}_{ini}, \quad (22c)$$

$$|q_i^{(n)}| = V, \quad \forall n \in \{0, 1, \dots, N\}, \quad (22d)$$

$$q_i^{(n)} \in \mathbf{Q}_i^d, \quad \forall n \in \{0, 1, \dots, N\}, \quad (22e)$$

$$(10c) - (10e), \quad (22f)$$

where $\mathbf{Q}_i^d, q_i^{(n)}, v_m^{(n)}$, and N denote the trajectory based on discreted map, position in the trajectory, RIS constraint on each position, and total number of positions in each trajectory, respectively.

Remark 1. The position set predicted by LSTM model constitutes an unstable sequence, and ARIMA model is able to make d -th order difference to the unstable sequence through its difference function, turning the sequence into a stable sequence.

5) *Complexity for LSTM-ARIMA Model*: The complexity of the proposed LSTM-ARIMA model is mainly related to the LSTM model, ARIMA model and the learning process.

Algorithm 1 LSTM-ARIMA algorithm for prediction

Input:

LSTM network structure, ARIMA model structure.

Return: The channel properties \mathbf{Z} , The predicted position pairs $\bar{\mathbf{S}}_{ini}$.

- 1: **Initialize:** Parameters of LSTM network, Parameters of ARIMA algorithm, \mathcal{X} robots, geographic map.
 - 2: Randomly generate $\mathbf{S} = \{S_1, S_2, \dots, S_{4N}, S_{4N+1}, \dots, S_{5N}\}$ as historical data for robots.
 - 3: Split \mathbf{S} into $\mathbf{S}_{train} = \{S_1, S_2, \dots, S_{4N}\}$ and $\mathbf{S}_{test} = \{S_{4N+1}, S_{4N+2}, \dots, S_{5N}\}$.
 - 4: Normalize historical data by invoking Min-Max normalization method: $\tilde{S}_{N,\bar{n}} = \frac{S_{N,\bar{n}} - S_{min}}{S_{max} - S_{min}}$.
 - 5: Input normalized $\tilde{\mathbf{S}} = \{\tilde{S}_1, \tilde{S}_2, \dots, \tilde{S}_{4N}\}$ and $\tilde{\mathbf{S}} = \{\tilde{S}_{4N+1}, \tilde{S}_{4N+2}, \dots, \tilde{S}_{5N}\}$.
 - 6: Calculate ν_i, ν_f, ν_o of LSTM network.
 - 7: **for** $n = 1$ to N **do**
 - 8: Predict $\bar{S}_n^{LSTM} = \nu_o \tanh(\mathbf{b}_u^{t-1}) * (S_{max} - S_{min}) + S_{min}$, $n \in \{1, 2, \dots, N\}$.
 - 9: Calculate $\Delta^d \bar{S}_n = (1 - L)^d (\bar{S}_n^{AR} + \bar{S}_n^{MA})$.
 - 10: Predict $\bar{S}_n^{ARIMA} = [\sum_{z=1}^a \kappa_z (1 - L)^d (\bar{S}_n^{AR} + \bar{S}_n^{MA}) + \sum_{z=1}^b \varphi_z \zeta_{n-z} + \mu_n + \zeta_n] * (S_{max} - S_{min}) + S_{min}$.
 - 11: Calculate \bar{w}_n and \tilde{w}_n by employing critic weight method.
 - 12: Calculate $\bar{S}_n = \bar{w}_n \bar{S}_n^{LSTM} + \tilde{w}_n \bar{S}_n^{ARIMA}$.
 - 13: **end for**
 - 14: Calculate $\mathcal{I} = \sqrt{\frac{1}{N} \sum_{n=1}^N (|\bar{S}_n - S_{4N+n}|)}$.
 - 15: **if** $\mathcal{I} > \varsigma$ **then**
 - 16: Go to line 1
 - 17: **end if**
 - 18: **if** $\mathcal{I} \leq \varsigma$ **then**
 - 19: Accept trained LSTM-ARIMA fusion model.
 - 20: Input \mathbf{S}_{ini} to the model and obtain \mathbf{Z} and $\bar{\mathbf{S}}_{ini}$.
 - 21: **end if**
-

LSTM model maintains a total of four parameters, which are input gate ν_i , output gate ν_o , forget gate ν_f and candidate states $\nu_c = \nu_o \tanh(\mathbf{b}_u^{t-1})$. The total size of parameters are $4\text{len}(\nu_f)(\text{len}(\nu_i) + \text{len}(\nu_c) + \text{len}(\nu_o))$, where the length of ν_o is the same with ν_f , ν_c is determined by ν_i, ν_f and ν_o , where $\text{len}(\cdot)$ denotes the length of parameter. Thus, the computation complexity of LSTM model is $4(\tilde{n}\tilde{m} + \tilde{n}^2 + \tilde{n})$, where $\tilde{n} = \text{len}(\nu_f)$, $\tilde{m} = \text{len}(\nu_i)$. The computation of ARIMA model is determined by the length of sequence, which can be expressed as \tilde{N} . In terms of learning process, the computational complexity is related to the number of timesteps T and episodes N_e . Thus, the computation complexity can be expressed as $O(4(\tilde{n}\tilde{m} + \tilde{n}^2 + \tilde{n})\tilde{N}TN_e)$.

B. Dueling double deep Q-network Algorithm for Trajectories Planning and Passive Beamforming Design

Our goal is to obtain the whole trajectory with a maximum sum-rate, while acquiring a local maximum sum-rate does not guarantee that the overall maximum sum-rate can be achieved. The double DQN algorithm [42] improves the accuracy of the

Q-value obtained by the DQN algorithm, which employs the current network to help estimate the Q-network. Specifically, it evaluates the action in the next state of the current state and chooses the optimal action, and then adopts it for the current state to estimate a new Q value. However, the optimal action selected and other unselected actions in the next state may have the same impact on the next state. On the basis of the double DQN algorithm, we introduce the dueling DQN [43] algorithm to evaluate the advantages of actions selected in the next state. In this section, we propose an ML-based algorithm, namely, D³QN-based algorithm, which is able to guide the robots to interact with the environment to train an optimal policy, and further for trajectory planning and the phase shifts of the RIS, as well as the power allocation from the AP to the robots. According to the Rician distribution and Shannon formula, the perfect channel estimation can be executed after the position, phase shift, and power allocation strategy are determined by the ML policy to obtain the reward.

1) *D³QN-based Algorithm for Trajectories Planning and RIS Design:* In the D³QN-based model, the AP acts as an agent, which is able to control both the power allocation policy from the AP to robots, the phase shift of RIS, and the robots' positions. At each timestep, the AP observes the state of the RIS-aided system and carries out an optimal action based on the policy, which is determined by Q-function. Following the selected action, the AP receives the reward when the current state of the model is transmitted to the next state, where the reward can be calculated by obtaining sum-rate from three robots.

The state space $\mathbf{E} = \{e_{\bar{t}}\}$ at each epoch of the RIS-aided multi-robot networks is defined into three parts: the current phase shift $\{\theta_m\} \in [0, 2\pi)$ of passive reflecting elements in the RIS at each position, the current position $q_i = (x_i, y_i, h_r)$ of the robot i , and the current group of allocation power $\{p_i\}$ from the AP to all the robots². Thus, the state space \mathbf{E} can be expressed as

$$\mathbf{E} = [\{\theta_m\} \quad \{q_i\} \quad \{p_i\}]. \quad (23)$$

The position q_i on the trajectories for robot i is determined by the previous position except initial-final positions generated from $\bar{\mathbf{S}}_{ini}$. The primary state space complexity is calculated as $(M + 2\mathcal{X})$. Additionally, the total number of positions N_i of each robot should meet the condition $N_i > \max_i (|x_{\xi_i} - x_{\xi_i}| + |y_{\xi_i} - y_{\xi_i}| - 1)$, which ensures that all robots can reach the final position from the initial position. The action space $\mathbf{F} = \{f_{\bar{t}}\}$ at each epoch of the RIS-aided multi-robot networks is defined into three parts: the available quantity of phase shifts $\{\frac{2\pi n_0}{2^{B_0}}, n_0 = 0, 1, 2, \dots, 2^{B_0} - 1\}$, the distance with moving direction $\mathbf{D}_i = \{d_r, d_l, d_0, d_u, d_d\}$ for robot i , the available quantity of power allocation $\{p^1, p^2, \dots, p^a\}_i$ for robot i . $B_0, d_g, \mathbf{g} = \{r, l, 0, u, d\}$, a denote the resolution for the RIS phase shift, the right-left-stillness-up-down direction with 1 unit pace, and the total number of the available power allocated to the robots. Note that action "stillness" in one state is only applicable to at most $\mathcal{X} - 1$ robots to choice. Thus, the action state \mathbf{F} can be expressed as

$$\mathbf{F} = [\{\{\frac{2\pi n_0}{2^{B_0}}\}_m\} \quad \{\mathbf{D}_i\} \quad \{\{p^1, p^2, \dots, p^a\}_i\}] \quad (24)$$

²The SIC decoding order is assumed perfect in this paper.

Accordingly, the primary action space complexity is calculated as $(2^{B_0} \cdot M + 5\mathcal{X} + a\mathcal{X})$. The reward is a considerable factor in the optimization of trajectory planning and passive beamforming design, which can be determined by observing the different sum-rate of all robots between two adjacent positions. In general, the robots move from the current position to the next position at each epoch. However, in order to avoid collisions among robots, if the distance between any two robots equals δ , the robots are regarded as located in adjacent cells, and there exists at least one action that is unable to be selected. According to the definition of reinforcement learning, the rewards brought by the unavailable action can be defined as -10. The reward function can be calculated as

$$\mathcal{R} = \begin{cases} -10, & \forall i_1, i_2 \in \{1, 2, \dots, \mathcal{X}\}, \|q_{i_1} - q_{i_2}\| = \delta, \\ \left[\sum_{i=1}^{\mathcal{X}} R_i(q_i(\mathbf{n}+1)) \right]_{(e)} - \left[\sum_{i=1}^{\mathcal{X}} R_i(q_i(\mathbf{n})) \right]_{(e)}, & \text{Otherwise,} \end{cases} \quad (25)$$

where $[\cdot]_e$ denote the sum-rate of trajectories of all robots at e -th epoch. Thus, it is observed from (25), that maximizing the long-term sum rewards makes the dedication to maximizing the optimized trajectories and passive beamforming of the RIS-aided multi-robot system. According to the double DQN model, there are two networks with the same structure: the current Q-network and the target Q-network. The agent selects actions with the maximum Q value from the current Q-network, which can be expressed as

$$f_{\bar{\tau}}^{\max}(e_{\bar{\tau}}, \psi_{\bar{\tau}}) = \arg \max_{f_{\bar{\tau}}} Q[\Gamma(e_{\bar{\tau}}), f_{\bar{\tau}}; \psi_{\bar{\tau}}], \quad (26)$$

and the Q-value in the target Q-network can be calculated by

$$Q'[\Gamma(e_{\bar{\tau}}), f_{\bar{\tau}}^{\max}(e_{\bar{\tau}}, \psi_{\bar{\tau}}); \bar{\psi}_{\bar{\tau}}] = \mathcal{R} + \eta Q[\Gamma(e_{\bar{\tau}}), f_{\bar{\tau}}^{\max}(e_{\bar{\tau}}, \psi_{\bar{\tau}}); \bar{\psi}_{\bar{\tau}}], \quad (27)$$

where $\Gamma(\cdot)$, $\psi_{\bar{\tau}}$, and $\bar{\psi}_{\bar{\tau}}$ are the feature vector of state, the parameters for the current Q-network and target Q-network, respectively. The update method for $\psi_{\bar{\tau}}$ and $\bar{\psi}_{\bar{\tau}}$ are independent. The updated for ψ_t can be expressed as

$$\begin{aligned} \psi'_t &= \psi_t + \eta_0 \{ Q'[\Gamma(e_{\bar{\tau}}), f_{\bar{\tau}}^{\max}(e_{\bar{\tau}}, \psi_{\bar{\tau}}); \bar{\psi}_{\bar{\tau}}] \\ &\quad - Q[\Gamma(e_{\bar{\tau}}), f_{\bar{\tau}}; \psi_{\bar{\tau}}] \} \nabla_{\psi_t} Q[\Gamma(e_{\bar{\tau}}), f_{\bar{\tau}}; \psi_{\bar{\tau}}], \\ &= \psi_t + \eta_0 \{ \mathcal{R} + \eta Q[\Gamma(e_{\bar{\tau}}), f_{\bar{\tau}}^{\max}(e_{\bar{\tau}}, \psi_{\bar{\tau}}); \bar{\psi}_{\bar{\tau}}] \\ &\quad - Q[\Gamma(e_{\bar{\tau}}), f_{\bar{\tau}}; \psi_{\bar{\tau}}] \} \nabla_{\psi_t} Q[\Gamma(e_{\bar{\tau}}), f_{\bar{\tau}}; \psi_{\bar{\tau}}], \end{aligned} \quad (28)$$

where the $\nabla_{\{\cdot\}}$ is the gradient operator. The loss function can be expressed as follows:

$$\begin{aligned} \text{Loss}(\psi_{\bar{\tau}}) &= \mathbb{E} \{ \{ Q'[\Gamma(e_{\bar{\tau}}), f_{\bar{\tau}}^{\max}(e_{\bar{\tau}}, \psi_{\bar{\tau}}); \bar{\psi}_{\bar{\tau}}] \\ &\quad - Q[\Gamma(e_{\bar{\tau}}), f_{\bar{\tau}}; \psi_{\bar{\tau}}] \}^2 \}, \\ &= \mathbb{E} \{ \{ \mathcal{R} + \eta Q[\Gamma(e_{\bar{\tau}}), f_{\bar{\tau}}^{\max}(e_{\bar{\tau}}, \psi_{\bar{\tau}}); \bar{\psi}_{\bar{\tau}}] \\ &\quad - Q[\Gamma(e_{\bar{\tau}}), f_{\bar{\tau}}; \psi_{\bar{\tau}}] \}^2 \}. \end{aligned} \quad (29)$$

For $\bar{\psi}_{\bar{\tau}}$, it can be updated every \mathbb{T} epochs according to $\psi_{\bar{\tau}}$, which can be expressed as

$$\bar{\psi}_{\bar{\tau}} = \psi_{\bar{\tau}}|_{\mathbb{T}}. \quad (30)$$

However, when evaluating the potential actions for the current positions, there are two or more actions that bring the same effect on the same states. In this case, the robot may sample the wrong action, which will affect the subsequent actions, resulting in the inability to finally obtain the optimal

trajectory. In order to reduce the impact of this case, we employ the dueling architecture in the target Q-network to evaluate state and action, respectively. The target Q-value obtained by equation (27) can be rewritten as:

$$\begin{aligned} Q'[\Gamma(e_{\bar{\tau}}), f_{\bar{\tau}}^{\max}(e_{\bar{\tau}}, \psi_{\bar{\tau}}); \bar{\psi}_{\bar{\tau}}^C, \bar{\psi}_{\bar{\tau}}^{\text{D1}}, \bar{\psi}_{\bar{\tau}}^{\text{D2}}] \\ = Q_e(\Gamma(e_{\bar{\tau}}); \bar{\psi}_{\bar{\tau}}^C, \bar{\psi}_{\bar{\tau}}^{\text{D2}}) + Q_f[\Gamma(e_{\bar{\tau}}), f_{\bar{\tau}}^{\max}(e_{\bar{\tau}}, \psi_{\bar{\tau}}); \bar{\psi}_{\bar{\tau}}^C, \bar{\psi}_{\bar{\tau}}^{\text{D1}}], \end{aligned} \quad (31)$$

where the $Q_e(\cdot)$, $Q_f(\cdot)$, $\bar{\psi}_{\bar{\tau}}^C$, $\bar{\psi}_{\bar{\tau}}^{\text{D1}}$, and $\bar{\psi}_{\bar{\tau}}^{\text{D2}}$ denote the state evaluate function, action evaluate function, parameters for convolutional layers, parameters for dense layer to present advantage function, and parameters for dense layer to present value function, respectively. $\bar{\psi}_{\bar{\tau}}$ denotes the parameters of double DQN network architecture, including three convolutional layers and two dense layers in series. In contrast, $\bar{\psi}_{\bar{\tau}}^C$ denotes the parameters of same three convolutional layers in dueling DQN structure, while $\bar{\psi}_{\bar{\tau}}^{\text{D1}}$ and $\bar{\psi}_{\bar{\tau}}^{\text{D2}}$ indicate the parameters of two dense layers in parallel. Note that, we only consider to split state and action of target Q-network in double DQN networks. Additionally, prioritized experience replay (PER) has been employed in dueling DQN network architecture. In the experience pool, PER extracts samples with a large absolute value of $|Q'[\Gamma(e_{\bar{\tau}}), f_{\bar{\tau}}^{\max}(e_{\bar{\tau}}, \psi_{\bar{\tau}}); \bar{\psi}_{\bar{\tau}}] - Q'[\Gamma(e'_{\bar{\tau}}), f_{\bar{\tau}}^{\max}(e'_{\bar{\tau}}, \psi_{\bar{\tau}}); \bar{\psi}_{\bar{\tau}}^C, \bar{\psi}_{\bar{\tau}}^{\text{D1}}, \bar{\psi}_{\bar{\tau}}^{\text{D2}}]|$ for training, which is capable of speeding up the learning of the Q function. More importantly, PER can enhance the identifiability of $Q_e(\cdot)$ and $Q_f(\cdot)$, we introduce a loss function, where the target Q-value can be rewritten as

$$\begin{aligned} Q'[\Gamma(e_{\bar{\tau}}), f_{\bar{\tau}}^{\max}(e_{\bar{\tau}}, \psi_{\bar{\tau}}); \bar{\psi}_{\bar{\tau}}^C, \bar{\psi}_{\bar{\tau}}^{\text{D1}}, \bar{\psi}_{\bar{\tau}}^{\text{D2}}] \\ = Q_e(\Gamma(e_{\bar{\tau}}); \bar{\psi}_{\bar{\tau}}^C, \bar{\psi}_{\bar{\tau}}^{\text{D2}}) + Q_f[\Gamma(e_{\bar{\tau}}), f_{\bar{\tau}}^{\max}(e_{\bar{\tau}}, \psi_{\bar{\tau}}); \bar{\psi}_{\bar{\tau}}^C, \bar{\psi}_{\bar{\tau}}^{\text{D1}}] \\ - \frac{1}{\mathcal{A}} \sum_{(f_e)' \in \mathcal{A}} Q_f[\Gamma(e_{\bar{\tau}}), (f_e)'(e_{\bar{\tau}}, \psi_{\bar{\tau}}); \bar{\psi}_{\bar{\tau}}^C, \bar{\psi}_{\bar{\tau}}^{\text{D1}}], \end{aligned} \quad (32)$$

where the \mathcal{A} , $(f_e^{\max})'(\cdot)$ are the action space of current Q-network and an sampled action from action space of current Q-network, respectively. The updated for $\{\bar{\psi}_{\bar{\tau}}^C, \bar{\psi}_{\bar{\tau}}^{\text{D1}}, \bar{\psi}_{\bar{\tau}}^{\text{D2}}\}$ can be expressed as (33), and the loss function can be calculated as (34).

Remark 2. *The advantage function refers to the advantage of the action value function compared to the value function of the current state. If the advantage function is greater than 0, it means that the action is worse than the average action, otherwise the current action is not as good as the average action. To avoid situations where the advantage function is equal to 0, we need to centralize the advantage function.*

Thus, according to the D³QN model mentioned above, through interaction with the environment, the agent is able to find the optimal policy for the robot trajectories planning and phase designing in the RIS. The detailed pseudo code is shown in **Algorithm 2**.

Remark 3. *The optimal policy of the Agent is always to choose the best action in any given state, while the best action often has smaller Q-values than the non-optimal ones in most cases. Such a problem is called the over-estimations of action value (Q-value). To overcome this problem, all the optional actions based on the current state need to be additionally evaluated as much as possible.*

$$\begin{aligned}
\{\overline{\psi}_t^C, \overline{\psi}_t^{D1}, \overline{\psi}_t^{D2}\}' &= \{\overline{\psi}_t^C, \overline{\psi}_t^{D1}, \overline{\psi}_t^{D2}\} + \eta_0 \{Q'[\Gamma(e_t'), f_t^{max}(e_t', \psi_t); \overline{\psi}_t^C, \overline{\psi}_t^{D1}, \overline{\psi}_t^{D2}] - Q[\Gamma(e_t), f_t; \psi_t] \nabla_{\psi_t} Q[\Gamma(e_t), f_t; \psi_t], \\
&= \{\overline{\psi}_t^C, \overline{\psi}_t^{D1}, \overline{\psi}_t^{D2}\} + \eta_0 \{ \mathcal{R} + \eta \{ Q_e(\Gamma(e_t'); \overline{\psi}_t^C, \overline{\psi}_t^{D2}) + Q_f[\Gamma(e_t'), f_t^{max}(e_t', \psi_t); \overline{\psi}_t^C, \overline{\psi}_t^{D1}] \\
&\quad - \frac{1}{\mathcal{A}} \sum_{(f_e)' \in \mathcal{A}} Q_f[\Gamma(e_t'), (f_e)'(e_t', \psi_t); \overline{\psi}_t^C, \overline{\psi}_t^{D1}] \} - Q[\Gamma(e_t), f_t; \psi_t] \nabla_{\psi_t} Q[\Gamma(e_t), f_t; \psi_t], \quad (33)
\end{aligned}$$

$$\begin{aligned}
\text{Loss}(\{\overline{\psi}_t^C, \overline{\psi}_t^{D1}, \overline{\psi}_t^{D2}\}) &= \mathbb{E} \{ \{ Q_e(\Gamma(e_t'); \overline{\psi}_t^C, \overline{\psi}_t^{D2}) + Q_f[\Gamma(e_t'), f_t^{max}(e_t', \psi_t); \overline{\psi}_t^C, \overline{\psi}_t^{D1}] \\
&\quad - \frac{1}{\mathcal{A}} \sum_{(f_e)' \in \mathcal{A}} Q_f[\Gamma(e_t'), (f_e)'(e_t', \psi_t); \overline{\psi}_t^C, \overline{\psi}_t^{D1}] - Q[\Gamma(e_t), f_t; \psi_t] \}^2 \}. \quad (34)
\end{aligned}$$

Algorithm 2 D³QN algorithm for trajectories planning and beamforming design

Input:

DQN network structure, LSTM network structure, ARIMA structure.

Return: Parameters of D³QN network, optimal initial-final pair $\overline{\mathbf{S}}_{\text{op}} = \{\varepsilon'_{\text{op}}, \xi'_{\text{op}}\}$.

- 1: **Initialize:** 2D moving space \mathbf{M} for robots, reply memory \mathcal{D} , the number of episodes E , minimal batch size \mathcal{M} , episodes indicator \overline{T} , parameters update frequency for target Q network $\mathbb{T} < \overline{T}$, temporary reward buffer vector $\mathfrak{R} = \{\mathfrak{R}_1, \mathfrak{R}_2, \dots, \mathfrak{R}_N\}$, and temporary initial-final position pairs buffer vector $\overline{\mathbf{S}} = \{\overline{\mathbf{S}}_1, \overline{\mathbf{S}}_2, \dots, \overline{\mathbf{S}}_N\}$, η , ϵ , \mathbf{Q} , \mathbf{Q}' , \mathbf{R} , \mathbf{E} , \mathbf{F} , ψ_t , $\overline{\psi}_t^C$, $\overline{\psi}_t^{D1}$, and $\overline{\psi}_t^{D2}$, Ψ , $\{p_i\}$, $\{\{p^1, p^2, \dots, p^a\}_i\}$, N .
 - 2: Explore the positions of AP, RIS, and boundaries in \mathbf{M} .
 - 3: Train LSTM-ARIMA model for initial-final position pairs prediction.
 - 4: Generate $\{\varepsilon'_i, \xi'_i\}$ by $\overline{\mathbf{S}}_{\text{ini}}$ for all the robots.
 - 5: **for** episode t from 1 to \overline{T} **do**
 - 6: The agent randomly selects $e_t \in \mathbf{E}$ and $\Gamma(e_t)$.
 - 7: Input $\Gamma(e_t)$ to current Q-network and obtain Q-values with all actions.
 - 8: Sample f_t of each e_t by invoking ϵ -greedy policy.
 - 9: Determine the decoding order based on the current state by NOMA method.
 - 10: Execute f_t and obtain $\Gamma(e_t')$ of new state, observe \mathcal{R} .
 - 11: **if** The reward achieve the terminated condition $|\mathcal{R} - \mathcal{R}_0| \leq \hat{\mathcal{R}}$ **then**
 - 12: Update for D³QN-network at current episode is terminal, **is_terminated = true**.
 - 13: **end if**
 - 14: Store transition $(\Gamma(e_t), f_t, \mathcal{R}, \Gamma(e_t'), \mathbf{is_terminated})$ in \mathcal{D} .
 - 15: Update e_t' as e_t .
 - 16: Sample \mathcal{M} of transition $(\Gamma(e_t)^l, f_t^l, \mathcal{R}^l, \Gamma(e_t')^l, \mathbf{is_terminated}^l)$ from \mathcal{D} , $l = 1, 2, 3, \dots, \mathcal{M}$.
 - 17: Calculate $Q'[\Gamma(e_t'), f_t^{max}(e_t', \psi_t); \overline{\psi}_t^C, \overline{\psi}_t^{D1}, \overline{\psi}_t^{D2}] =$

$$\begin{cases} \mathcal{R}, & \text{if } \mathbf{is_terminated} \text{ is true,} \\ Q_e(\Gamma(e_t'); \overline{\psi}_t^C, \overline{\psi}_t^{D2}) + \{ Q_f[\Gamma(e_t'), f_t^{max}(e_t', \psi_t); \overline{\psi}_t^C, \overline{\psi}_t^{D1}] - \frac{1}{\mathcal{A}} \sum_{(f_e)' \in \mathcal{A}} Q_f[\Gamma(e_t'), (f_e)'(e_t', \psi_t); \overline{\psi}_t^C, \overline{\psi}_t^{D1}] \}, & \text{if } \mathbf{is_terminated} \text{ is false,} \end{cases}$$
 - 18: Perform a gradient descent step to calculate all parameters of target Q-network:
$$\begin{aligned}
\{\overline{\psi}_t^C, \overline{\psi}_t^{D1}, \overline{\psi}_t^{D2}\}' &= \{\overline{\psi}_t^C, \overline{\psi}_t^{D1}, \overline{\psi}_t^{D2}\} + \eta_0 \{ \mathcal{R} + \eta \{ Q_e(\Gamma(e_t'); \overline{\psi}_t^C, \overline{\psi}_t^{D2}) + Q_f[\Gamma(e_t'), f_t^{max}(e_t', \psi_t); \overline{\psi}_t^C, \overline{\psi}_t^{D1}] \\
&\quad - \frac{1}{\mathcal{A}} \sum_{(f_e)' \in \mathcal{A}} Q_f[\Gamma(e_t'), (f_e)'(e_t', \psi_t); \overline{\psi}_t^C, \overline{\psi}_t^{D1}] \} - Q[\Gamma(e_t), f_t; \psi_t] \nabla_{\psi_t} Q[\Gamma(e_t), f_t; \psi_t].
\end{aligned}$$
 - 19: **if** $\overline{T} \% \mathbb{T} = 1$ **then**
 - 20: Update parameters $\{\overline{\psi}_t^C, \overline{\psi}_t^{D1}, \overline{\psi}_t^{D2}\} = \{\overline{\psi}_t^C, \overline{\psi}_t^{D1}, \overline{\psi}_t^{D2}\}'$
 - 21: **end if**
 - 22: **if** $\{\hat{t} = \overline{T}\}$ or $\{\mathcal{R} \text{ remains the same for } \overline{T} \text{ consecutive episodes}\}$ **then**
 - 23: Update for D³Q network is end, reset \mathbf{E} , \mathbf{F} .
 - 24: Execute $\mathfrak{R}_t = \mathcal{R}$, and $\overline{\mathbf{S}}_t = \{\varepsilon'_i, \xi'_i\}$.
 - 25: **end if**
 - 26: **end for**
 - 27: Obtain $\overline{\mathbf{S}}_{\text{op}} = \{\varepsilon'_{\text{op}}, \xi'_{\text{op}}\}$ from N explored pairs according to \mathfrak{R} .
-

2) *Dueling DQN-based Algorithm for Trajectories Planning and RIS Design*: According to the dueling DQN-based algorithm, the AP acts as an agent, where the equipped controller has the capability of determining power allocation policy from the AP to robots, the phase shift of RIS, and the robots' positions. The state space $\mathbf{E} = \{e_{\bar{t}}\}$, action space $\mathbf{F} = \{f_{\bar{t}}\}$, and reward function \mathcal{R} in double DQN-based algorithm are the same with D³QN-based algorithm. According to dueling DQN architecture, the Q-function is split as follows:

$$\begin{aligned} Q' & [\Gamma(e'_{\bar{t}}), f_{\bar{t}}(e'_{\bar{t}}, \psi_{\bar{t}}); \bar{\psi}_{\bar{t}}^C, \bar{\psi}_{\bar{t}}^{D1}, \bar{\psi}_{\bar{t}}^{D2}] \\ & = Q_e(\Gamma(e'_{\bar{t}}); \bar{\psi}_{\bar{t}}^C, \bar{\psi}_{\bar{t}}^{D2}) + Q_f[\Gamma(e'_{\bar{t}}), f_{\bar{t}}(e'_{\bar{t}}, \psi_{\bar{t}}); \bar{\psi}_{\bar{t}}^C, \bar{\psi}_{\bar{t}}^{D1}] \\ & \quad - \frac{1}{\mathcal{A}} \sum_{(f_e)' \in \mathcal{A}} Q_f[\Gamma(e'_{\bar{t}}), (f_e)'(e'_{\bar{t}}, \psi_{\bar{t}}); \bar{\psi}_{\bar{t}}^C, \bar{\psi}_{\bar{t}}^{D1}], \end{aligned} \quad (35)$$

and the loss function can be expressed as

$$\begin{aligned} \text{Loss}(\{\bar{\psi}_{\bar{t}}^C, \bar{\psi}_{\bar{t}}^{D1}, \bar{\psi}_{\bar{t}}^{D2}\}) \\ & = \mathbb{E}\{\{Q_e(\Gamma(e'_{\bar{t}}); \bar{\psi}_{\bar{t}}^C, \bar{\psi}_{\bar{t}}^{D2}) \\ & \quad + Q_f[\Gamma(e'_{\bar{t}}), f_{\bar{t}}(e'_{\bar{t}}, \psi_{\bar{t}}); \bar{\psi}_{\bar{t}}^C, \bar{\psi}_{\bar{t}}^{D1}] \\ & \quad - \frac{1}{\mathcal{A}} \sum_{(f_e)' \in \mathcal{A}} Q_f[\Gamma(e'_{\bar{t}}), (f_e)'(e'_{\bar{t}}, \psi_{\bar{t}}); \bar{\psi}_{\bar{t}}^C, \bar{\psi}_{\bar{t}}^{D1}] \\ & \quad - Q[\Gamma(e_{\bar{t}}), f_{\bar{t}}; \psi_{\bar{t}}]\}^2\}. \end{aligned} \quad (36)$$

3) *Double DQN-based Algorithm for Trajectories Planning and RIS Design*: According to double DQN-based algorithm, the AP acts as an agent, where the power allocation policy from the AP to robots, phase shift of RIS, and the robots' positions are determined by the controller. The state space $\mathbf{E} = \{e_{\bar{t}}\}$, action space $\mathbf{F} = \{f_{\bar{t}}\}$, and reward function \mathcal{R} in double DQN-based algorithm are the same with D³QN-based algorithm, while the Q-function and loss function follows (27) and (29). Moreover, the parameters of target Q-network are updated follows (30).

4) *Complexity for D³QN Algorithm*: The complexity of performance of the D³QN algorithm depends on the convolution layers and learning process. Thus, the complexity of convolution layers can be expressed as $\sum_{\bar{d}=1}^{\bar{D}} \bar{M}_{\bar{d}}^2 \bar{K}_{\bar{d}}^2 \bar{C}_{\bar{d}-1} \bar{C}_{\bar{d}}$, where \bar{D} , \bar{d} , \bar{K} , and \bar{C}_d denote the total number of convolution layers, the number of rows of feature two-dimension data, length of convolution kernel, and the number of convolution kernels of \bar{d} -th layer, respectively. For learning process, according to the reinforcement learning method, denote \bar{F} , \bar{E} and \bar{T} as the total number of actions that the agent is able to choose, the number of saved state-action pairs, and the number of timesteps. Thus, the computation complexity of D³QN can be expressed as $O((\sum_{\bar{d}=1}^{\bar{D}} \bar{M}_{\bar{d}}^2 \bar{K}_{\bar{d}}^2 \bar{C}_{\bar{d}-1} \bar{C}_{\bar{d}} + |\bar{F}| + \bar{E})\bar{T})$.

IV. NUMERICAL RESULTS

In this section, we provide simulation results to verify the effectiveness of the proposed machine learning-based optimization algorithms for joint trajectories planning and passive beamforming design, as well as the performance of the algorithms. In the simulations, the number of robots \mathcal{X} is denoted as 3, which are randomly located in the initial positions obtained by the LSTM-ARIMA model. The RIS is fixed at the center of the ceiling with a height of 3m, and the standard size of space is 8m and 6m. Additionally, there are

four pillars with regular size $1\text{m} \times 1\text{m} \times 3\text{m}$, two parterres with regular size $1\text{m} \times 1\text{m} \times 1\text{m}$, and a fountain with a regular base size of $1.5\text{m} \times 1.5\text{m} \times 1\text{m}$ in the space. Note that, all the objects' heights mentioned above are more than that of the robot. The maximal transmit power at AP is pre-defined as 10 dBm, while the height of AP is defined as 2m. The number of reflecting elements in RIS K and in sub-surface \bar{K} is defined as 20 and 5, while the total number of sub-surface M is only increased linearly with K . The other simulation parameters are provided in Table. I. The performance of the proposed LSTM-ARIMA and D³QN algorithms, the trajectories for all the robots, and the achievable sum-rate for all the robots are analyzed in the following sections.

TABLE I: Simulation parameters

Parameter	Description	Value
C	Path loss when $d = 1\text{m}$	-30dB
δ^2	Noise power variance	-75dBW
$\bar{\alpha}_{Ai}$	Path loss factor for AP-robot link	3.5
$\bar{\alpha}_{Ri}$	Path loss factor for RIS-robot link	2.8
$\bar{\alpha}_{AI}$	Path loss factor for AP-RIS link	2.2
\bar{T}	The replay memory capacity	1000
\bar{D}	The number of episodes	10000
\bar{N}	The size of minibatch	64
η	Discount factor	0.9
ψ	Learning rate	0.05
ϵ	Probability decision value	0.1

A. Initial-final Position Pairs Predicted by LSTM-ARIMA

We adopted the LSTM-ARIMA model to predict the position pairs, where the position pairs are partitioned into initial positions and final positions groups. The possible initial position and final position are determined in the designated space, whose ranges are denoted as $[1,7]\text{m} \times [5,6]\text{m}$ and $[1.1,6.9]\text{m} \times (0,1)\text{m}$. As shown in Fig. 4, the performance of the LSTM, ARIMA, and LSTM-ARIMA fusion models has been investigated under training samples by the different unit N definitions. For LSTM, the changing trend of the RMSE proves the long sequence degrades its performance. It can be obtained that when the unit achieves $N = 80$, LSTM is able to achieve the best performance than the other unit definition. For ARIMA, it is obtained that its performance keeps improving with N increasing since it is able to handle long sequence issues. However, the historical data mentioned above is commonly regarded as a long unregular sequence, where LSTM and ARIMA have the capability of solving the unregular sequence and long sequence, respectively. Therefore, the dynamic weights \bar{w}_n and \tilde{w}_n are assigned to LSTM and ARIMA, which are dynamic determined for each predicted position pair. After fusion of the advantages of LSTM and ARIMA, when the unit is defined as $N = 90$, LSTM-ARIMA achieves the best performance than other cases. As shown in Fig. 5, the results for predicted initial and final position pairs are provided by the well-trained model ($N = 90$). With the reliable model for initial-final position pairs prediction, the optimal trajectory with maximum sum-rate can be further explored.

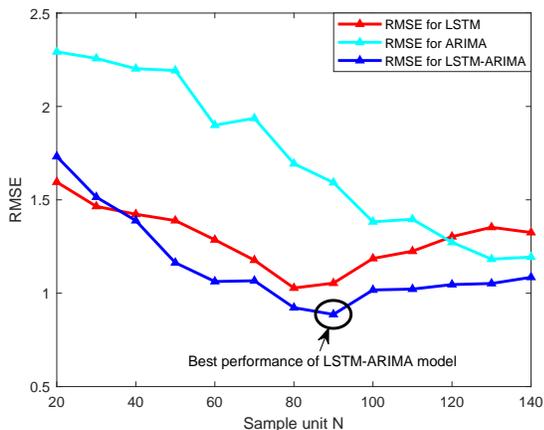


Fig. 4: Prediction performance of different model.

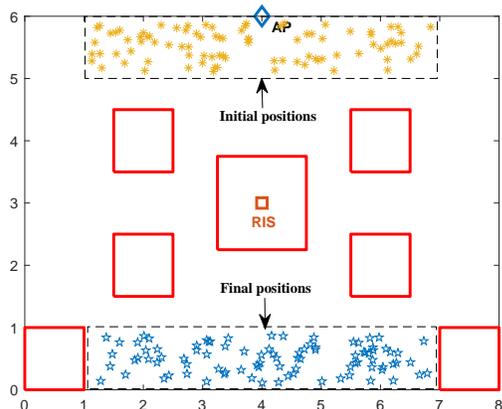


Fig. 5: The range of optimal initial-final position pairs predicted by LSTM-ARIMA Model ($N = 90$).

B. Convergence for D^3QN Algorithm

The performance of the ML-based algorithm occupies a pivotal place in the entire optimization. For the proposed D^3QN algorithm, after selecting the initial-final positions of the robot, we optimize the sum-rate for all robots. In order to analyze the performance of D^3QN , we compared D^3QN algorithm to the double DQN algorithm and the dueling DQN algorithm with $\mathcal{X} = 3$, $\mathcal{P} = 10dbm$. As shown in Fig. 6, the double DQN algorithm and dueling DQN algorithm have no big difference in convergence speed. They can converge when episodes are 689 and 723 respectively. However, the convergence speed of the proposed D^3QN is faster than these two algorithms, reaching 618. It is worth noting that in virtue of ϵ -greedy strategy, the convergence episodes of these three algorithms cannot be guaranteed to be the same during each training. Therefore, the result given in the figure is the average convergence given by 10 repetitive training.

C. Trajectory and Achieved Sum-rate for Robots

In this subsection, we provide the trajectory planning results for the robots. There are three pre-defined conditions for the robot's trajectory design. Firstly, denote the velocity of robots and the map resolution as 0m/s (stillness) or 0.1 m/s (movement), and 0.1 m, respectively, which guarantees each robot is able to move at most one cell at each timestep. Then, the size of a cell has been approximated to the center point of the cell which makes the robot arrive at the center of

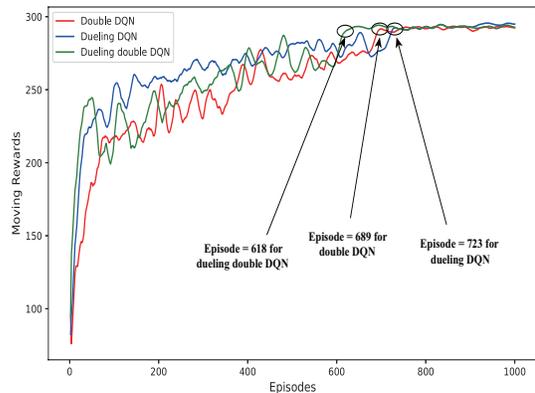


Fig. 6: Performance for D^3QN algorithm comparing different algorithms.

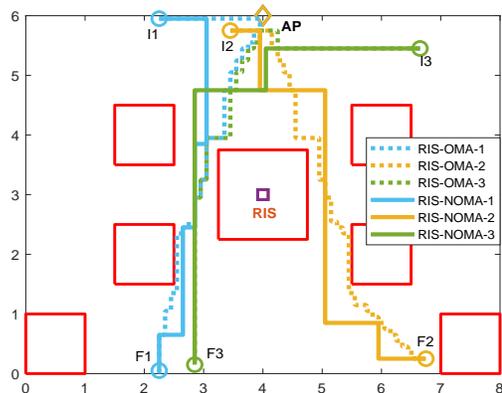


Fig. 7: Trajectories for each robot under RIS-OMA and RIS-NOMA cases, $K = 30$.

cells at each timestep and easily characterizes the received communication information. Finally, the movement direction of a robot is back, forth, stillness, left, and right. In the following, we discuss the results of trajectory design for the robots.

1) *Obtained Trajectories for Robots:* As shown in Fig. 7, the planned paths for all robots are depicted under "RIS-OMA" and "RIS-NOMA" cases, while the number of elements and sub-surfaces of the RIS is denoted as 30 and 6. The "o" with " $I_{\bar{w}}, \bar{w} = \{1, 2, 3\}$ " denotes the initial position for the robots, while the "o" with " $F_{\bar{w}}, \bar{w} = \{1, 2, 3\}$ " represents the final position. Moreover, dotted lines and solid lines with three colors are utilized to describe the paths under NOMA-aided and OMA-aided networks. The initial and final positions are the possible optimal pairs according to the proposed ML framework, which is randomly selected from the obtained range by LSTM-ARIMA and finally determined by the D^3QN algorithm. As shown in Fig. 8 and Fig. 9, the received communication rate is calculated on the robot located positions, which directly reflects the communication quality of the whole trajectories on each robot. In the two figures, the dark blue represents the area where robots are unselected or unreached, while the other colors indicate the received communication rate for robots' trajectories. Intuitively, it can be observed that the obtained trajectories by NOMA-aided and OMA-aided networks both tend to be close to RIS, while the

shape of paths by NOMA-aided networks seems more regular than that of OMA-aided networks. One probable reason is, that for each robot, the OMA technology tends to achieve a reflected LOS-dominated channel model by traversing the cells covered by the RIS. The other reason is we aim to achieve the maximum long-term benefits in a time period instead of at each timestep. In other words, NOMA can improve the efficiency of trajectories design in a limited time instead of OMA. Additionally, all trajectories are designed close to the RIS and AP, which is able to verify the RIS is able to improve the received communication rate of robots

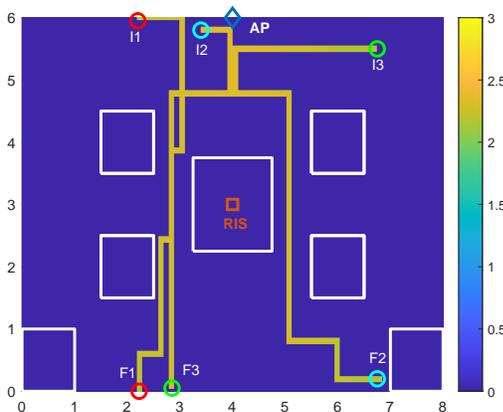


Fig. 8: Communication rate for trajectories based on NOMA, $K = 30$.

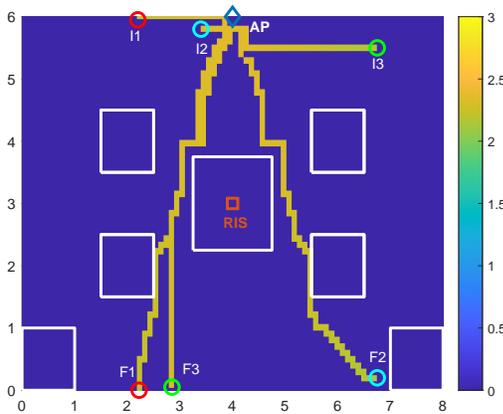


Fig. 9: Communication rate for trajectories based on OMA, $K = 30$.

2) *Impact of RIS on Robot Trajectory*: The employment of RIS brings an impact on trajectory design for the robots, which is demonstrated in Fig. 10 and Fig. 11. In Fig. 10, we compare “RIS-NOMA” and “RIS-OMA” cases, where the designed total path length increased from 24.6m to 28.2m in NOMA networks, while it upgrades from 28.5m to 29.4m in OMA networks. This is because, with adding more elements in the RIS, the planned trajectories under the “RIS-NOMA” case improve the overall communication quality of the environment, which makes robots have more spaces to explore. Additionally, it can be observed that the planned path length in the “RIS-NOMA” case is longer than the path length in “RIS-OMA” case, which indicates the “stillness” are selected by robots at some states. Compared to the RIS deployment cases, the paths of robots in “without RIS” cases are the

longest, where the path length of each robot is the same in the NOMA and OMA networks. This is because the position of AP, positions of obstacles, simultaneously arriving conditions, and communication quality limitations make the length of the communication-sensitive path become different and much longer than the geographic path.

In Fig. 11, mark the “OMA” as a benchmark scheme, the maximal sum-rate for three robots at any positions on their designed trajectories is obtained. For the rate of robot 1 and robot 2, we can observe that they have the same rate at several adjacent positions. This is because they select “stillness” at the current state, which proves they are constrained by the simultaneous arrival condition. These positions are the maximum rate position on the trajectory for each robot. “stillness” guarantees to continue to receive high-quality communication information for robots while other robots still move. In the figure, it is observed that the maximum rate of position for each robot is 2.6m, 2.8m, 3.2m, robot 1 chooses “stillness” 17 times while robot 2 chooses “stillness” 4 times. It can be obtained that the total length of optimized trajectories for all robots is 7.8m, 9.1m, and 9.5m, while for all robots, the maximum sum-rate has been achieved when the total path length of each robot is 2.6m, 2.8m, and 3.1m. Finally, compared to the OMA-aided scheme, the “RIS-NOMA” case outperforms the “RIS-OMA” case.

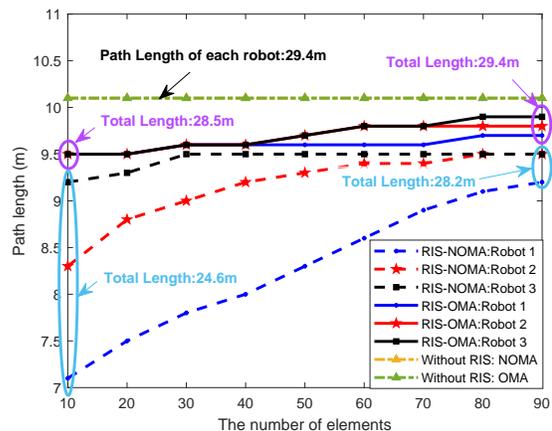


Fig. 10: Path length for all robots with different number of elements in RIS.

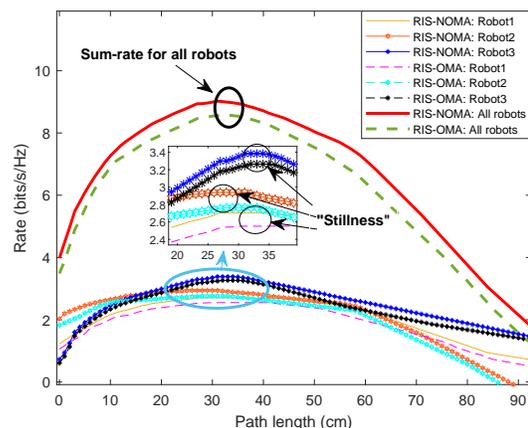


Fig. 11: sum-rate versus path length with $K=30$.

In Fig. 12, the comparison of the proposed algorithm,

double DQN algorithm, dueling DQN algorithm, and DQN is presented, where the different elements in the RIS will influence the obtained total path length. It can be observed that the proposed D^3QN algorithm is able to find a shorter total path than the conventional ML algorithms. This is because the structure of the proposed D^3QN algorithm significantly supplements the limitations of the conventional ML algorithm. Therefore, it can be obtained that the performance of the D^3QN algorithm outperforms the conventional ML algorithms. Additionally, the ML solutions for the RIS-NOMA case all achieve a shorter total path length than the “Without RIS: NOMA” case. The effectiveness of RIS on total path length optimization has been proved.

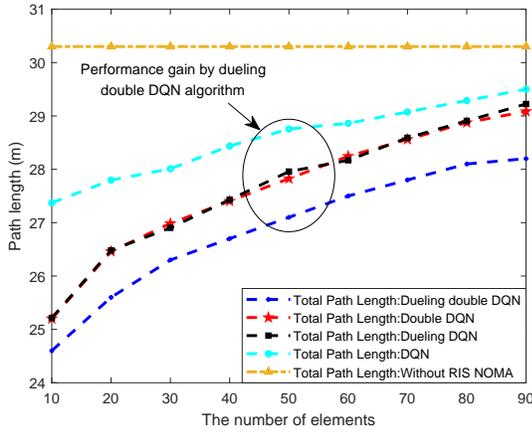


Fig. 12: sum-rate versus path length with $K=30$ by different algorithms.

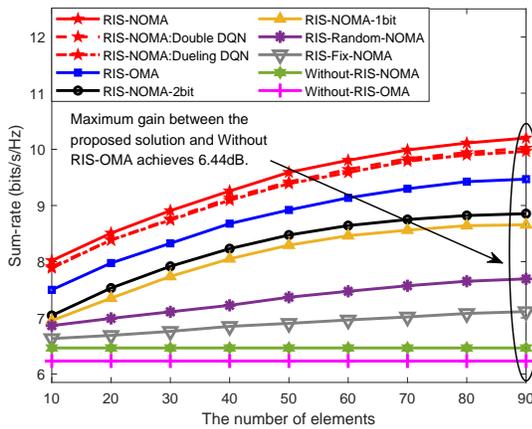


Fig. 13: Sum-rate of whole trajectories for all robots under different number of elements.

In Fig. 13, the sum-rate of all robots’ trajectories is demonstrated, where the different elements in the RIS will influence the obtained communication sum-rate of robots. In order to verify the benefits brought by the RIS, we compare “RIS-NOMA” with benchmarks of “RIS-OMA”, “RIS-NOMA-2bit”, “RIS-NOMA-1bit”, “RIS-Random-NOMA”, “RIS-Fix-NOMA”, “Without-RIS-NOMA”, and “Without-RIS-OMA” cases. It can be observed that the sum-rate performances of all considered RIS-aided schemes increase with the increase of elements. The performances of “RIS-NOMA” significantly outperform the other seven benchmark schemes. The achieved sum-rate by adjusting the RIS phase shift with higher resolu-

tion bits outperforms the lower resolution bits, while the gap between the “RIS-NOMA” and the RIS-aided low phase shift resolution case becomes larger when the elements increase. It is observed that the performance achieved by optimizing RIS phase shifts outperforms the “RIS-Random” case whose reflection coefficients are randomly set. For “RIS-Fix-NOMA”, the phase shift of RIS is fixed and the long-term power allocation is optimized, it can be obtained that a significant performance degradation occurred compared to jointly optimizing the phase shift and power allocation. Finally, it can be observed that all RIS-NOMA cases outperform the RIS-OMA cases, while the maximum gain between the proposed solution and Without RIS-OMA achieves 6.44dB. Additionally, the performance of different algorithms on the sum-rate optimization is also considered, where the proposed algorithm outperforms the double DQN and dueling DQN algorithms. The performance of the double DQN and dueling DQN algorithms is close. Considering the results for path length and the maximum sum-rate of the path, it can be obtained that employing NOMA technology is able to strike a balance between the maximum sum-rate exploration and the shortest path design between the initial position and final position.

In Fig. 14, to explore the influence brought by decoding order, we compare three schemes based on the different numbers of clusters: optimal order, random order, and fixed order. Optimal order denotes the best order selected by the D^3QN -based algorithm. As shown in the Figure, the DQN-based optimal order scheme significantly outperforms the random order scheme and fixed order, which highlights the necessity of exploring the optimal decoding order. Additionally, it can be observed that the gaps among the optimal schemes and the random order schemes present a huge difference with the increase of elements. The fixed order becomes insensitive to sum-rate improvement, which keeps very slightly changing with the elements increase. This is because optimal order schemes provide the optimal power allocated policy for each robot, which is ignored in random order schemes and fixed schemes.

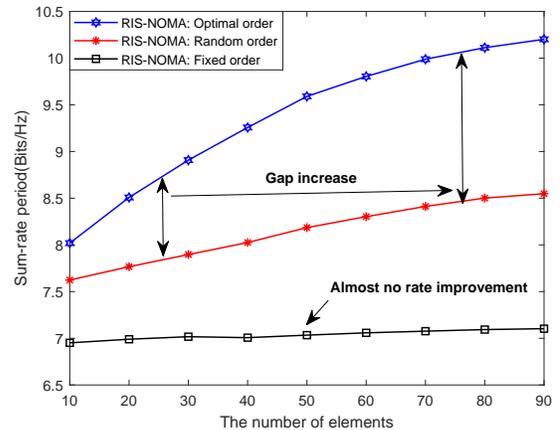


Fig. 14: Sum-rate versus decoding order under different elements numbers of RIS.

Remark 4. The optimal decoding order is found by an exhaustive search algorithm, which costs at least $\mathcal{X}!$ iterations, which

is better to be invoked in the scenario with a small number of decoding users, otherwise, it will bring high complexity to the whole algorithm.

V. CONCLUSION

In this paper, we explored a downlink RIS-aided multi-robot NOMA network. The robot trajectories' sum-rate maximization problem was formulated by jointly optimizing trajectories for robots, reflecting the coefficients matrix, the decoding order, and the power allocation at the AP, subject to the QoS for all the robots. To tackle the formulated problem, a novel machine learning scheme was proposed, which combines the LSTM-ARIMA model and D³QN algorithm. LSTM-ARIMA algorithm is employed to predict the possible initial and final positions for robots, while the D³QN algorithm is invoked to plan optimal trajectories for the robots and design the phase shift matrix, determining the optimal initial and final positions for robots. Numerical results were provided for demonstrating that the proposed RIS-aided NOMA networks achieve significant gains compared to RIS-OMA and without-RIS-aided schemes. The investigated LSTM-ARIMA and D³QN algorithms attained considerable performance compared to the vanilla ML algorithm. In a real application, to improve the efficiency of the communication system, energy is another key element that needs to be considered. Therefore, striking a tradeoff between energy and sum-rate is a potential research topic in future work [45]. Moreover, with respect to the multi-robot system in the indoor environment, deploying a single RIS may not be efficient to establish LOS links due to the dense obstacles. In this case, the multiple distributed RISs can be considered to leverage the richer reflected paths to further improve performance [9], [46].

REFERENCES

- [1] X. Gao, Y. Liu and X. Mu, "Trajectory and Passive Beamforming Design for IRS-aided Multi-Robot NOMA Indoor Networks," *Proc. IEEE Int. Conf. Commun. (ICC)*, Montreal, Canada, 2021, pp. 1-6.
- [2] E. Kagan, N. Shvalb, and I. Ben-Gal, "Autonomous Mobile Robots and Multi-Robot Systems: Motion-Planning, Communication, and Swarming," *John Wiley & Sons*, 2019.
- [3] M. Batalin and G. Sukhatme, "Coverage, Exploration and Deployment by a Mobile Robot and Communication Network," *Telecommun. Syst.*, vol. 26, 181–196, 2004.
- [4] Y. Saito, Y. Kishiyama, A. Benjebbour, T. Nakamura, A. Li, and K. Higuchi, "Non-Orthogonal Multiple Access (NOMA) for Cellular Future Radio Access," *Proc. IEEE 77th Veh. Tech. Conf. (VTC Spring)*, Dresden, Germany, 2013, pp. 1-5.
- [5] L. Dai, B. Wang, Y. Yuan, S. Han, I. Chih-lin, and Z. Wang, "Non-orthogonal multiple access for 5G: solutions, challenges, opportunities, and future research trends," *IEEE Commun. Mag.*, vol. 53, no. 9, pp. 74-81, September 2015.
- [6] Z. Ding, P. Fan, and H. V. Poor, "Impact of user pairing on 5G nonorthogonal multiple access," *IEEE Trans. Veh. Technol.*, vol. 65, no. 8, pp. 6010-6023, Aug. 2016.
- [7] Y. Zhu, G. Zheng, and K. Wong, "Stochastic Geometry Analysis of Large Intelligent Surface-Assisted Millimeter Wave Networks," *IEEE J. Sel. Areas Commun.*, vol. 38, no. 8, pp. 1749-1762, Aug. 2020.
- [8] Q. Wu and R. Zhang, "Towards smart and reconfigurable environment: Intelligent reflecting surface aided wireless network?" *IEEE Commun. Mag.*, vol. 58, no. 1, pp. 106–112, 2020.
- [9] W. Mei, B. Zheng, C. You and R. Zhang, "Intelligent Reflecting Surface-Aided Wireless Networks: From Single-Reflection to Multireflection Design and Optimization," *Proc. IEEE*, vol. 110, no. 9, pp. 1380-1400, Sept. 2022.
- [10] J. Zhao and Y. Liu, "A survey of intelligent reflecting surfaces (IRSs): towards 6G wireless communication networks," *arXiv preprint arXiv:1907.04789*, 2019.
- [11] E. Basar, M. Di Renzo, J. De Rosny, M. Debbah, M. -S. Alouini, and R. Zhang, "Wireless Communications Through Reconfigurable Intelligent Surfaces," *IEEE Access*, vol. 7, pp. 116753-116773, 2019.
- [12] J. Si et al., "Covert transmission assisted by intelligent reflecting surface," *IEEE Trans. Commun.*, vol. 69, no. 8, pp. 5394-5408, 2021.
- [13] Q. Chen, M. Li, X. Yang, R. Alturki, M. D. Alshehri, and F. Khan, "Impact of Residual Hardware Impairment on the IoT Secrecy Performance of RIS-Assisted NOMA Networks," *IEEE Access*, vol. 9, pp. 42583-42592, 2021.
- [14] Y. Li, M. Jiang, Q. Zhang and J. Qin, "Joint Beamforming Design in Multi-Cluster MISO NOMA Reconfigurable Intelligent Surface-Aided Downlink Communication Networks," *IEEE Trans. on Commun.*, vol. 69, no. 1, pp. 664-674, Jan. 2021.
- [15] Z. Ding and H. V. Poor, "A simple design of IRS-NOMA transmission," *IEEE Commun. Lett.*, vol. 24, no. 5, pp. 1119-1123, May 2020.
- [16] X. Mu, Y. Liu, L. Guo, J. Lin, and N. Al-Dhahir, "Exploiting Intelligent Reflecting Surfaces in NOMA Networks: Joint Beamforming Optimization", *IEEE Trans. Wirel. Commun.*, vol. 19, no. 10, pp. 6884-6898, Oct. 2020.
- [17] X. Gao, Y. Liu, X. Liu and L. Song, "Machine Learning Empowered Resource Allocation in IRS Aided MISO-NOMA Networks," *IEEE Trans. Wirel. Commun.*, vol. 21, no. 5, pp. 3478-3492, May 2022.
- [18] Z. Li et al., "Energy Efficient Reconfigurable Intelligent Surface Enabled Mobile Edge Computing Networks With NOMA," *IEEE Trans. Cogn. Commun. Netw.*, vol. 7, no. 2, pp. 427-440, June 2021.
- [19] F. Fang, Y. Xu, Q. -V. Pham, and Z. Ding, "Energy-Efficient Design of IRS-NOMA Networks," *IEEE Trans. Veh. Technol.*, vol. 69, no. 11, pp. 14088-14092, Nov. 2020.
- [20] X. Mu, Y. Liu, L. Guo, J. Lin and R. Schober, "Intelligent Reflecting Surface Enhanced Indoor Robot Path Planning: A Radio Map-Based Approach," *IEEE Trans. Wirel. Commun.*, vol. 20, no. 7, pp. 4732-4747, July 2021.
- [21] M. Zeng, X. Li, G. Li, W. Hao and O. A. Dobre, "Sum Rate Maximization for IRS-Assisted Uplink NOMA," *IEEE Commun. Lett.*, vol. 25, no. 1, pp. 234-238, Jan. 2021.
- [22] W. Ni, X. Liu, Y. Liu, H. Tian, and Y. Chen, "Resource Allocation for Multi-Cell IRS-Aided NOMA Networks," *IEEE Trans. Wirel. Commun.*, vol. 20, no. 7, pp. 4253-4268, July 2021.
- [23] H. Wang, C. Liu, Z. Shi, Y. Fu, and R. Song, "On Power Minimization for IRS-Aided Downlink NOMA Systems," *IEEE Wirel. Commun. Lett.*, vol. 9, no. 11, pp. 1808-1811, Nov. 2020.
- [24] Y. Liu et al., "Robotic Communications for 5G and Beyond: Challenges and Research Opportunities," *IEEE Commun. Mag.*, vol. 59, no. 10, pp. 92-98, October 2021.
- [25] B. Woosley, P. Dasgupta, J. Rogers, and J. Twigg, "Multi-robot information driven path planning under communication constraints," *Auton. Robot.*, vol. 44, pp. 721–737, 2020.
- [26] A. Dutta, A. Ghosh, and O. Kreidl, "Multi-robot Informative Path Planning with Continuous Connectivity Constraints," *Proc. Int. Conf. Robot. Autom. (ICRA)*, Montreal, Canada, 2019, pp. 3245-3251.
- [27] M. Corah, C. O'Meadhra, K. Goel, and N. Michael, "Communication-Efficient Planning and Mapping for Multi-Robot Exploration in Large Environments," *IEEE Robot. Autom. Lett.*, vol. 4, no. 2, pp. 1715-1721, April 2019.
- [28] E. Olcay, F. Schuhmann, and B. Lohmann, "Collective navigation of a multi-robot system in an unknown environment," *Rob. Auton. Syst.*, vol. 132, pp. 103604, 2020.
- [29] J. Walton, M. Wallace, and S. Howard, "Multiple-access multiple-input multiple-output (MIMO) communication system," *U.S. Patent 7,636,573*, issued Dec. 22, 2009.
- [30] T. S. Rappaport, R. W. Heath Jr, R. C. Daniels, and J. N. Murdock, "Millimeter wave wireless communications," *Pearson Education*, 2015.
- [31] C. Tatino, N. Pappas, and D. Yuan, "Multi-Robot Association-Path Planning in Millimeter-Wave Industrial Scenarios," *IEEE Netw. Lett.*, vol. 2, no. 4, pp. 190-194, Dec. 2020.
- [32] A. Pogue, S. Hanna, A. Nichols, X. Chen, D. Cabric, and A. Mehta, "Path Planning Under MIMO Network Constraints for Throughput Enhancement in Multi-robot Data Aggregation Tasks," *Proc. Int. Conf. Robot. Artif. Intell. (IROS)*, 2020, pp. 11824-11830.
- [33] B. Zhou et al., "Performance Limits of Visible Light-Based Positioning for Internet-of-Vehicles: Time-Domain Localization Cooperation Gain," *IEEE Trans. Intell. Transp. Syst.*, vol. 22, no. 8, pp. 5374-5388, Aug. 2021.
- [34] B. Zhou, A. Liu, and V. Lau, "Visible Light-Based User Position, Orientation and Channel Estimation Using Self-Adaptive Location-Domain Grid Sampling," *IEEE Trans. Wirel. Commun.*, vol. 19, no. 7, pp. 5025-5039, July 2020.
- [35] G. Nason, Stationary and non-stationary time series. in *Stat. Volcanol.*, 60, 2006.

- [36] Q. Wu and R. Zhang, "Intelligent reflecting surface enhanced wireless network via joint active and passive beamforming," *IEEE Trans. Wirel. Commun.*, 2019.
- [37] Z. Ding et al., "Application of Non-Orthogonal Multiple Access in LTE and 5G Networks," *IEEE Commun. Mag.*, vol. 55, no. 2, pp. 185-191, February 2017.
- [38] J. Zhang, L. Zhu, Z. Xiao, X. Cao, D. O. Wu and X. -G. Xia, "Optimal and Sub-Optimal Uplink NOMA: Joint User Grouping, Decoding Order, and Power Control," *IEEE Wirel. Commun. Lett.*, vol. 9, no. 2, pp. 254-257, Feb. 2020.
- [39] R. Eckhardt, S. Ulam, and J. Von Neumann, "the monte carlo method," *Los Alamos Sci.*, no. 15, pp. 131, 1987.
- [40] S. Patro, and K. Sahu, "Normalization: A preprocessing stage," *arXiv preprint arXiv:1503.06462*, 2015.
- [41] T. V. Laarhoven, "L2 regularization versus batch and weight normalization," *arXiv preprint arXiv:1706.05350*, 2017.
- [42] D. Diakoulaki, G. Mavrotas, and L. Papayannakis, "Determining objective weights in multiple criteria problems: The critic method," *Comput. Oper. Res.*, vol. 22, no. 7, pp. 763-770, 1995.
- [43] H. Van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double q-learning," *Proc. Conf. AAAI Artif. Intell.*, Phoenix, USA, 2016, vol. 30, no. 1.
- [44] Z. Wang, T. Schaul, and M. Hessel, "Dueling network architectures for deep reinforcement learning," *Proc. 33rd Int. Conf. Mach. Learn. (ICML)*, New York City, USA, 2016, pp. 1995-2003.
- [45] Z. Li et al., "Multiobjective optimization based sensor selection for TDOA tracking in wireless sensor network," *IEEE Trans. Veh. Tech.*, vol. 68, no. 12, pp. 12360-12374, 2019.
- [46] W. Mei and R. Zhang, "Multi-Beam Multi-Hop Routing for Intelligent Reflecting Surfaces Aided Massive MIMO," *IEEE Trans. Wirel. Commun.*, vol. 21, no. 3, pp. 1897-1912, March 2022.