# Energy-Efficient Communication Networks via Multiple Aerial Reconfigurable Intelligent Surfaces: DRL and Optimization Approach

Pyae Sone Aung, Yu Min Park, Yan Kyaw Tun, *Member, IEEE,* Zhu Han, *Fellow, IEEE,* and Choong Seon Hong, *Senior Member, IEEE,*

arXiv:2207.03149v2 [cs.NI] 8 Dec 2022

## Abstract

In the realm of wireless communications in 5G, 6G and beyond, deploying unmanned aerial vehicle (UAV) has been an innovative approach to extend the coverage area due to its easy deployment. Moreover, reconfigurable intelligent surface (RIS) has also emerged as a new paradigm with the goals of enhancing the average sum-rate as well as energy efficiency. By combining these attractive features, an energy-efficient RIS-mounted multiple UAVs (aerial RISs: ARISs) assisted downlink communication system is studied. Due to the obstruction, user equipments (UEs) can have a poor line of sight to communicate with the base station (BS). To solve this, multiple ARISs are implemented to assist the communication between the BS and UEs. Then, the joint optimization problem of deployment of ARIS, ARIS reflective elements on/off states, phase shift, and power control of the multiple ARISs-assisted communication system is formulated. The problem is challenging to solve since it is mixed-integer, non-convex, and NP-hard. To overcome this, it is decomposed into three sub-problems. Afterwards, successive convex approximation (SCA), actor-critic proximal policy optimization (AC-PPO), and whale optimization algorithm (WOA) are employed to solve these sub-problems alternatively. Finally, extensive simulation results have been generated to illustrate the efficacy of our proposed algorithms.

## Index Terms

Pyae Sone Aung, Yu Min Park, Yan Kyaw Tun, and Choong Seon Hong are with the Department of Computer Science and Engineering, Kyung Hee University, Yongin-si, Gyeonggi-do 17104, Rep. of Korea, e-mail:{pyaesoneaung, yumin0906, ykyawtun7, cshong}@khu.ac.kr.

Zhu Han is with the Electrical and Computer Engineering Department, University of Houston, Houston, TX 77004, and the Department of Computer Science and Engineering, Kyung Hee University, Yongin-si, Gyeonggi-do 17104, Rep. of Korea, email{zhan2}@uh.edu.

Aerial reconfigurable intelligent surface (ARIS), deployment, reflective elements on/off, phase shift, transmit power optimization, successive convex optimization (SCA), actor-critic proximal policy optimization (AC-PPO), whale optimization algorithm (WOA).

## I. INTRODUCTION

### A. Background and Motivations

As claimed by Cisco Networking Index (CNI), the number of Internet users reached 3.9 billion in 2018 and is anticipated to surpass 5.3 billion by 2023 [1]. Rapid growth of multimedia devices such as the Internet of Things (IoT), video streaming, online gaming, Virtual Reality (VR) and Augmented Reality (AR) applications, thrives immense challenges for current communication architecture and motivates to discover new ways to enhance spectral efficiency in both academic and industrial fields. Numerous ingenious wireless technologies have been developed in the last several years, which includes deploying unmanned aerial vehicles (UAVs) and Reconfigurable Intelligent Surfaces (RIS) elements.

Recently, UAVs have achieved a great deal of interests to deploy as a communication and computing platforms due to their high mobility and ease of deployment. The emplacement of UAVs can not only save the cost of mobile infrastructure which demands a large budget but also save time for quick on-demand deployment to provide services in rural regions or disaster areas or temporary events such as concerts, stadiums where the infrastructure is difficult to come across. In some scenarios [2], [3], UAVs are implemented with a multi-access edge computing (MEC) system to deliver the computing resources near to the user equipment (UE) which saves a considerably large amount of time for uploading, computing and downloading tasks.

The newly recent technology called RIS, which is incited from the recent development of meta-surfaces, benefits the wireless communications in extending the coverage range and improving the signal quality at the receiver [4]. RIS is a man-made meta-surface implemented with low-cost passive elements that can be programmed by integrated electronic circuits to alter the incoming electromagnetic field into the desirable way [5]. Unlike the traditional collaborative communications such as decode-and-forward (DF) and amplify-and-forward (AF), RIS does not require additional power amplifier hence, is more environmental friendly and energy-efficient [6]. Taking into account of its cost efficiency and energy efficiency, RIS technology has acquired a vast attention in 5G, 6G and beyond communications. Furthermore, since RIS structures consist

of relatively small hardware components, they can be easily integrated in several communication environments such as along the surfaces of the building [7].

## B. Challenges and Research Contributions

When UAVs are considered as communication and computing platforms , there exists several challenges in UAVs' energy consumption as they are energy-constrained devices. On the other side, even though RIS can enhance the spectral efficiency, setting up RIS structures to achieve Line-of-sight (LoS) links between UE and RIS is still quite challenging issues. Taking the advantages of RIS in enhancing spectral efficiency without the requirement of any external power sources with the aid of UAVs to obtain LoS links between UE and RIS, we propose the multiple aerial RISs (ARISs)-assisted system to extend the downlink communication links from the ground base station (BS) to the UEs. In our system, we assume that there is no dominant LoS links between the BS and UEs due to obstacles. The contributions of our paper can be organized as following:

- Firstly, we propose the downlink communication system between the BS and UEs, which is assisted by the multiple ARISS to enhance the spectral efficiency for all UEs since the dominant LoS links between BS and UEs are blocked by the obstacles. We assume the BS and ARISs are deployed by the same service operator and thus the BS is responsible for ARIS deployment and controlling the on/off states and the phase shifts for the ARIS reflective elements.

- Secondly, we formulate the problem to maximize the energy efficiency of the proposed system by jointly optimizing the ARISs deployment, ARIS reflective elements on/off states, phase shift, and power control. We show that the formulated problem is a mixed integer non-linear programming (MINLP) problem and it is challenging to solve in the polynomial time.

- To address this challenge, we decompose our formulated problem into three sub-problems: 1) ARISs deployment problem, 2) joint ARIS reflective elements on/off states and phase shift problem, and 3) power control problem. Then, successive convex approximation (SCA), actor-critic proximal policy optimization (AC-PPO), and whale optimization algorithm (WOA) are proposed to solve these sub-problems, alternatively.

- Finally, a comprehensive numerical analysis is integrated to validate efficacy of the overall performance of our proposed algorithms with several benchmark schemes, such as single-

ARIS, ARIS with fixed phase shifts (ARIS-NPS), and UAV as relay (UAV-relay) scenarios. We achieve the improvement in average sum-rate by 24% and 58% compared to the single-ARIS and the ARIS-NPS scenarios, and 43% and 72% increase in energy-efficiency compared to the single-ARIS and the UAV-relay scenarios, respectively. Moreover, our proposed multiple ARISs-assisted system achieves 69% increase in average sum-rate compared to the multiple RISs-assisted system.

The rest of the paper is categorized as follows: we present the related works in Section II. Next, we present our system model and problem formulation in Section III. Afterwards, the solution approach is proposed in Section IV, and performance evaluation is performed in Section V. Finally, Section VI concludes our paper.

## II. RELATED WORKS

### A. UAV-assisted wireless networks

An overview on the literature related to UAV-assisted wireless networks are discussed in this section [8]–[14]. The major strength of UAV in enhancing coverage area, energy-efficiency, and cost-efficiency has received significant attention in recent years [8]. In [9], the authors studied to maximize the uplink communication where UAVs are served as relays. In [10], the authors studied a single UAV-assisted device-to-device (D2D) communications and analyse how the appropriate UAV's altitude can impact the rate performance and coverage area on the D2D users' density. The authors in [11] derived the channel model of the LoS probability for the air-to-ground UAV communications. There exists several works that studied upon UAV deployment [12]–[15]. The authors in [12] studied the UAVs deployment for UAV-to-ground communication in arbitrary spatial distribution for network planning to provide wireless services to the ground users and the authors in [13] studied the incorporation between UAVs in 3D cellular network. The work in [14] studied the adaptive UAV deployment for the dynamic users. The authors in [15] studied DRL-based dynamic UAV control instead of static UAV deployment. In all of the above works, UAV is considered either as aerial BS or MEC devices or relays, which results in higher energy consumption.

### B. RIS-assisted wireless networks

An overview on the literature related to RIS-assisted wireless networks are discussed in this section [5], [16]–[23]. In [5], the authors considered to develop the energy-efficient architecture

for the RIS structures in accordance with power allocation and phase shifting values of RIS elements while guaranteeing the individual data rate budget for each user. In [16], the authors proposed the energy-harvesting RIS elements implemented on the facades of the buildings in order to maximize the spectral efficiency while enabling the transmit power control and RIS configuration under the indeterminate wireless channel condition. The authors in [17] aimed to distinguish the principal relationship between the total sum-rate of multiple users and the required number of RIS reflective elements in wireless communications. They observed the capacity of the system could no longer efficiently rise as the number of RIS elements reaches the upper bound limit. They also investigated how the number of phase shifts can effect the performance on the achievable data rate. The authors in [18] investigated the practical case study between phase shift and finite-sized RIS to maximize the downlink multi-user system. In [19], the authors studied about RIS elements to eliminate interference between multiple D2D uplink communication network. There has also been several studies on RIS-assisted in the vehicular networks. In [20], the authors investigated the secrecy outage probability upon vehicular-to-vehicular (V2V) and vehicular-to-infrastructure (V2I). The authors in [21] aimed to maximize the data rate for each vehicle where the communication links from the road site unit (RSU) is extended by the RIS technology with discrete phase shift. The authors in [22] studied deep reinforcement learning (DRL) based RIS-assisted multi-user downlink multiple input single output (MISO) system. The work in [23] considered to improve the secrecy rate of users in RIS-assisted system by constructing DRL-based QoS-aware reward function. All of the aforementioned works only considered the RIS-assisted networks, where RIS elements are either implemented on the ground level or facades of the building, which is still challenging to achieve the dominant LoS communication links between the BS-RIS-users.

## C. UAV-RIS-assisted wireless networks

An overview on the literature related to UAV-RIS-assisted wireless networks are discussed in this section [24]–[32]. The authors in [24] examined the adaptive RIS-assisted aerial-terrestrial downlink communication system between UAVs and multi-users with respect to RIS elements allocation and reflective coefficients. In [25], the authors looked into UAV-user communication with RIS assistance in order to maximize the worst-case secrecy rate by taking into account of the transmitter's power allocation, RIS's beamforming and UAV's trajectory. The authors in [26] proposed the RIS-assisted UAV communications to maximize the received signal power

Fig. 1: System model for RIS-mounted UAVs.

at the ground user by considering the passive and active beamforming and UAV's trajetory. Furthermore, in [27], the authors minimized the energy consumption problem for both orthogonal multiple access (OMA) and non-orthogonal multiple access (NOMA) cases by jointly considering the trajectory for the UAV and passive beamforming of the RIS elements. There also exists several works on ARIS-assisted system [28], [29]. In [30], the authors considered ARIS-assisted system to satisfy the constraints of ultra-reliable low latency communication (URLLC). The authors in [31] studied the several UAVs-RISs-assisted total transmit power minimization for the heterogeneous networks. They did not consider the energy efficiency of the system. The authors in [32] considered to maximize energy efficiency for a single ARIS-assisted downlink communication for single user. They did not consider for the multiple ARISs-assisted scenario. In this paper, we propose the multiple ARISs in order to maximize the average energy efficiency for the downlink communication between the BS and the UEs.

## III. SYSTEM MODEL

Our system model includes a BS $B$ with multiple antennas, a set $\mathcal{N}$ of $N$ ARISs in which each RIS is implemented on each UAV and each ARIS $n \in \mathcal{N}$ contains an array of $\mathcal{I}_n = [1_n, 2_n, \ldots, I_n]$, reflective elements and a set $\mathcal{K}$ of $K$ UEs with single antenna as shown in Fig. 1. The coordinates of the BS is denoted by $\boldsymbol{q}_B = (x_B, y_B, z_B)$, where $z_B$ is the height of the BS. Similarly, the positions of each UE $k$ and each ARIS $n$ can be represented as $\boldsymbol{q}_k = (x_k, y_k, 0)$

and $\boldsymbol{q}_n = (x_n, y_n, z_n)$ respectively, and $z_n$ is the height where the RIS-mounted UAV is hovering. The time horizon of the system can be divided into a discrete set of $\mathcal{T} = [1, 2, ..., t, ..., T]$.

Since ARIS has limited energy, apart from hovering, the ARIS reflective elements need to be turned off when there is no connection in order to reserve excessive energy. Authors in [33] and [24] prove that turning off the the whole RIS or some surface area of RIS can preserve energy. In our work, we define $\boldsymbol{\Delta} \in \mathbb{R}^{|\mathcal{N}| \times |\mathcal{I}_n|}$ as the on/off states matrix for all reflective elements $|\mathcal{I}_n|$ for each ARIS $n$ to decide whether to turn on or off. Therefore, the on/off states of the reflective element $i_n$ in each ARIS $n$ are controlled by the decision variable $\delta_{i_n}$ as follows:

$$\delta_{i_n}[t] = \begin{cases} 1, & \text{if reflective element } i \text{ of ARIS } n \text{ is switched on at time } t, \\ 0, & \text{otherwise.} \end{cases} \tag{1}$$

### A. Communication Model

We adopt both direct and indirect communication links between the BS and the UEs. For the direct link, we assume there is no dominant propagation along the LoS signal between the BS and UEs. Therefore, we adopt the Rayleigh fading model and the channel gain for the BS-UE link at time $t$ can be obtained as follows:

$$\mathbf{H}_{B,k}[t] = \sqrt{\kappa d_{B,k}^{-\alpha}[t]} \tilde{h}, \tag{2}$$

where $\kappa$ is the channel gain at the reference distance 1 m, $\alpha \geq 2$ is the path loss exponent, $|d_{B,k}[t]| = ||\boldsymbol{q}_B[t] - \boldsymbol{q}_k[t]||$ is the Euclidean distance between the BS and UE $k$ at time $t$, and $\tilde{h}$ is the complex Gaussian random scattering component with zero mean and unit variance.

For the indirect communication, there exists two links: BS-ARIS link and ARIS-UE link, respectively. For the BS-ARIS link, we assume there is only LoS signal between the BS and ARIS, and thus the channel fading here is assumed to experience the Rician channel fading with only LoS components. Therefore, the channel gain between the BS and ARIS $n$ at time $t$ can be defined as:

$$\mathbf{h}_{B,n}[t] = \sqrt{\kappa d_{B,n}^{-\alpha}[t]} \sqrt{\frac{\hat{R}}{1 + \hat{R}}} \mathbf{h}_{B,n}^{\text{LoS}}[t], \tag{3}$$

where $\hat{R}$ is the Rician factor, and $|d_{B,n}[t]| = ||\boldsymbol{q}_B[t] - \boldsymbol{q}_n[t]||$ is the distance between the BS to ARIS $n$ at time $t$. $\mathbf{h}_{B,n}^{\text{LoS}}[t]$ is the deterministic LoS component between the BS and ARIS $n$ in correspondence with the azimuth angle-of-arrival (AoA) of the link at time $t$ [24]. For the ARIS-UE link, there are both LoS and non-line-of-sight (NLoS) propagation between ARISs and

UEs. Consequently, the Rician fading model is adopted and the channel gain for the ARIS-UE link at time $t$ can be obtained as follows:

$$\mathbf{h}_{n,k}[t] = \sqrt{\kappa d_{n,k}^{-\alpha}[t]}\sqrt{\frac{\hat{R}}{1+\hat{R}}}\mathbf{h}_{n,k}^{\text{LoS}}[t] + \sqrt{\frac{1}{1+\hat{R}}}\mathbf{h}_{n,k}^{\text{NLOS}}, \tag{4}$$

where $|d_{n,k}[t]| = ||\boldsymbol{q}_n[t] - \boldsymbol{q}_k[t]||$ is the distance between ARIS $n$ and UE $k$ at time $t$. $\mathbf{h}_{n,k}^{\text{LoS}}[t]$ is the deterministic LoS component between ARIS $n$ and UE $k$ corresponding with the azimuth angle-of-departure (AoD) of the link and $\mathbf{h}_{n,k}^{\text{NLOS}}$ is the non-LoS component which follows the identically and independently distributed circularly-symmetric complex Gaussian distribution.

Furthermore, at time $t$, the incident signals are reflected by each reflective element $i$ of ARIS $n$ from the feasible range of phase shift values specified by

$$\theta_{i_n}[t] = e^{\left(\frac{2\pi\phi}{2^b}\right)}, \tag{5}$$

where $\phi$ is the phase shift index, and $b$ is the phase shift resolution in bits [34]. Therefore, a vector of $\boldsymbol{\theta}_{i_n}[t] = [\theta_{1_n}[t], \theta_{2_n}[t], \ldots, \theta_{I_n}[t]]$ represents the phase shift values of ARIS $n$. Following that, the reflection coefficient matrix can be denoted by

$$\boldsymbol{\Theta}_n[t] = \text{diag}(\beta_{1_n}e^{j\theta_{1_n}[t]}, \beta_{2_n}e^{j\theta_{2_n}[t]}, \ldots, \beta_{I_n}e^{j\theta_{I_n}[t]}), \tag{6}$$

where $\beta_{i_n} \in [0,1]$ denotes the amplitude reflection coefficient of the $i$-th reflective element of the $n$-th ARIS, and $j$ is the imaginary unit of a complex number. Therefore, the received signal at UE $k$ can be achieved as follows:

$$y_k[t] = \left(\mathbf{H}_{B,k}[t] + \sum_{n=1}^{N}\sum_{i=1}^{I_n}\delta_{i_n}[t]\mathbf{h}_{n,k}[t]\boldsymbol{\Theta}_n[t]\mathbf{h}_{B,n}[t]\right)\boldsymbol{x} + \omega_k, \tag{7}$$

where $\boldsymbol{x} = \sum_{k=1}^{K}\boldsymbol{g_k}[t]s_k$ is the transmitted signal from the BS with beamforming vector $\boldsymbol{g_k}[t]$ at time $t$, and the unit-power complex based information symbol $s_k$ for UE $k$, while $\omega_k \sim \mathcal{CN}(0, \sigma^2)$ denotes the additive white Gaussian noise (AWGN) at UE $k$. Based on (7), the signal-to-interference-plus-noise ratio (SINR) received at UE $k$ can be obtained as

$$\gamma_k[t] = \frac{\left|\left(\mathbf{H}_{B,k}[t] + \sum_{n=1}^{N}\sum_{i=1}^{I_n}\delta_{i_n}[t]\mathbf{h}_{n,k}[t]\boldsymbol{\Theta}_n[t]\mathbf{h}_{B,n}[t]\right)\boldsymbol{g_k}[t]\right|^2}{\sum_{l=1,l\neq k}^{K}|(\mathbf{H}_{B,k}[t] + \sum_{n=1}^{N}\sum_{i=1}^{I_n}\delta_{i_n}[t]\mathbf{h}_{n,k}[t]\boldsymbol{\Theta}_n[t]\mathbf{h}_{B,n}[t])\boldsymbol{g_l}[t]|^2 + \sigma^2}, \tag{8}$$

Afterwards, based on (8), the achievable data rate of UE $k$ can be formulated as follows:

$$r_k[t] = W\log_2(1 + \gamma_k[t]), \tag{9}$$

where $W$ is the transmission bandwidth available for each UE. Therefore, the sum-rate of all users can be described as follows:

$$R[t] = \sum_{k=1}^{K}r_k[t]. \tag{10}$$

## B. Power Consumption Model

In our scenario, we need to take account of the power consumption of ARIS hovering. We assume that ARISs are considered to be hovering at the designated altitude and thus, rotary-wing UAV is adopted. Therefore, the power consumption for the hovering of the rotary-wing UAV, $P_{\text{UAV}}$ can be obtained as follows [35]:

$$P_{\text{UAV}} = \frac{\nu}{8}\varphi\Lambda\eta v_a^3\varrho + (1+\iota)\frac{\tilde{w}^{3/2}}{\sqrt{2\varphi\eta}}, \tag{11}$$

which contains two terms: power required to rotate the rotor blades, and power required to endure the induced drag generated by the lift. The symbols $\nu$, $\varphi$, $\Lambda$, $\eta$, $v_a$, and $\varrho$ represent the coefficient of the profile drag, density of the air, rotor solidity, disc area of the rotor, blade angular velocity, and radius of the rotor, respectively. Moreover, $\iota$ and $\tilde{w}$ denote the incremental correction factor, and weight of the aircraft, respectively.

Furthermore, in this work, the BS controls the phase shifts of the ARIS reflective elements. Hence, the total power of the considered multiple ARISs-assisted downlink system includes: 1) transmit power of the BS, 2) circuit power of the each UE $k$, 3) circuit power consumption of ARIS and 4) hovering power of the rotary-wing UAV [33], and is defined as

$$P[t] = \sum_{k=1}^{K}(\zeta\boldsymbol{g_k}[t]^H\boldsymbol{g_k}[t] + P_k^{\text{cir}}) + \sum_{n=1}^{N}\sum_{i=1}^{I_n}\delta_{i_n}[t]I_nP_{\text{ARIS}} + P_{\text{UAV}}, \tag{12}$$

where $\zeta = 1/\mu$ with $\mu$ being the transmit power amplifier efficiency, $P_k^{\text{cir}}$ is the circuit power of each user $k$, and $P_{\text{ARIS}}$ is the power consumption for each ARIS. The transmit signal power of the BS has the constraint as follows:

$$\text{tr}(\boldsymbol{g}[t]^H\boldsymbol{g}[t]) \leq P_{\text{max}}, \forall t \in \mathcal{T}, \tag{13}$$

where $\text{tr}(\mathbf{S})$ means the trace of square matrix $\mathbf{S}$, $\boldsymbol{g} = [\boldsymbol{g_1}; \ldots; \boldsymbol{g_k}]$ and $P_{\text{max}}$ is the maximum transmission power available at the BS.

## C. Problem Formulation

The main objective of this work is to maximize energy efficiency of the system, i.e., to maximize the average sum-rate $R[t]$ for the UEs under the constraint of the power consumption $P[t]$ of both ARISs and the BS. To accomplish this, we need to jointly optimize the deployment of ARIS, ARIS reflective element on/off states, and phase shift, and power control of the BS. Prior to problem formulation, we define the required constraints as follows:

Each UE $k$ is necessary to fulfill the demand for the specified data rate at time $t$, which is defined as:

$$r_k[t] \geq r_k^{\min}[t], \forall k \in \mathcal{K}, \forall n \in \mathcal{N}, \forall i \in \mathcal{I}_n, \forall t \in \mathcal{T}, \tag{14}$$

where $r_k^{\min}[t]$ is the minimum data rate requirement for each UE $k$ at time $t$. The accessible phase shift value of $i$-th reflective element $n$-th ARIS at time $t$ should be between 0 to $2\pi$ as follows:

$$0 \leq \theta_{i_n}[t] < 2\pi, \forall n \in \mathcal{N}, \forall i \in \mathcal{I}_n, \forall t \in \mathcal{T}. \tag{15}$$

A safe distance between two adjacent ARISs is necessary to ensure that the coverage area of each ARIS does not overlap with that of other. Thereby, it can avoid the interference between different ARIS. We denote $d_{\min}[t]$ as the threshold distance between two adjacent ARISs at time $t$, and can be defined as follows:

$$||\boldsymbol{q}_i[t] - \boldsymbol{q}_j[t]||^2 \geq d_{\min}[t], \forall i, j \in \mathcal{N}, i \neq j, \forall t \in \mathcal{T}. \tag{16}$$

Furthermore, each reflective element $i$ of ARIS $n$ can only be either turned on or off at one time slot and can be given as follows:

$$\delta_{i_n}[t] \in \{0, 1\}, \forall n \in \mathcal{N}, \forall i \in \mathcal{I}_n, \forall t \in \mathcal{T}. \tag{17}$$

Given the above mentioned network characteristics, our optimization problem $\mathcal{E}$ can be mathematically formulated as follows:

$$\mathbf{P:} \max_{\boldsymbol{q}, \boldsymbol{\Delta}, \boldsymbol{\Theta}, \boldsymbol{g}} \mathcal{E}(\boldsymbol{q}, \boldsymbol{\Delta}, \boldsymbol{\Theta}, \boldsymbol{g}) \tag{18a}$$

$$\text{s.t.}$$

$$r_k[t] \geq r_k^{\min}[t], \forall k \in \mathcal{K}, \forall n \in \mathcal{N}, \forall i \in \mathcal{I}_n, \forall t \in \mathcal{T}, \tag{18b}$$

$$0 \leq \theta_{i_n}[t] < 2\pi, \forall n \in \mathcal{N}, \forall i \in \mathcal{I}_n, \forall t \in \mathcal{T}, \tag{18c}$$

$$||\boldsymbol{q}_i[t] - \boldsymbol{q}_j[t]||^2 \geq d_{\min}[t], \forall i, j \in \mathcal{N}, i \neq j, \forall t \in \mathcal{T}, \tag{18d}$$

$$\text{tr}(\boldsymbol{g}[t]^H \boldsymbol{g}[t]) \leq P_{\max}, \forall t \in \mathcal{T}, \tag{18e}$$

$$\delta_{i_n}[t] \in \{0, 1\}, \forall n \in \mathcal{N}, \forall i \in \mathcal{I}_n, \forall t \in \mathcal{T}. \tag{18f}$$

The objective function in (18a) is shown in (19). The problem $\mathbf{P}$ is a mixed integer non-linear programming (MINLP) problem which is non-convex. Therefore, it is challenging to solve the whole problem in polynomial time. Furthermore, there are couplings in both the objective

Fig. 2: Flow diagram of joint ARIS deployment, ARIS reflective elements on/off states, phase shift and power control problem.

$$\mathcal{E}(\boldsymbol{q},\boldsymbol{\Delta},\boldsymbol{\Theta},\boldsymbol{g}) = \frac{1}{T}\sum_{t=0}^{T}\frac{R[t]}{P[t]} = \frac{\sum_{k=1}^{K} W \log_2\left(1 + \frac{\left|\left(\mathbf{H}_{B,k}[t]+\sum_{n=1}^{N}\sum_{i=1}^{I_n}\delta_{i_n}[t]\mathbf{h}_{n,k}[t]\boldsymbol{\Theta}_n[t]\mathbf{h}_{B,n}[t]\right)\boldsymbol{g_k}[t]\right|^2}{\sum_{l=1,l\neq k}^{K}|(\mathbf{H}_{B,k}[t]+\sum_{n=1}^{N}\sum_{i=1}^{I_n}\delta_{i_n}[t]\mathbf{h}_{n,k}[t]\boldsymbol{\Theta}_n[t]\mathbf{h}_{B,n}[t])\boldsymbol{g_l}[t]|^2+\sigma^2}\right)}{\sum_{k=1}^{K}(\zeta\boldsymbol{g_k}[t]^H\boldsymbol{g_k}[t]+P_k^{\mathrm{cir}})+\sum_{n=1}^{N}\sum_{i=1}^{I_n}\delta_{i_n}[t]I_nP_{\mathrm{ARIS}}+P_{\mathrm{UAV}}} \tag{19}$$

function and constraints between the ARIS deployment, ARIS reflective elements on/off states and phase shift, and power control of the BS. Therefore, problem **P** is quite implausible to solve and there is no effective solution approach to deal with these difficulties.

Thus, we first decompose our optimization problem **P** into three sub-problems, **P1:** ARIS deployment problem, **P2:** ARIS reflective elements on/off states and phase shift problem, and **P3:** power control problem. Then, we solve the sub-problems iteratively until we reach the convergence and the detailed figure of our proposed solution technique is shown in Fig. 2.

$$\dot{\mathcal{E}}(\boldsymbol{q}) = \sum_{k=1}^{K} W \log_2 \left( 1 + \frac{\left( H_{B,k}[t] + \sum_{n=1}^{N} \sum_{i=1}^{I_n} \kappa \delta_{i_n}[t] g_k[t] \mathbf{h}_{ab}^T[t] \mathbf{H}'[t] \mathbf{h}_{ab}[t] \right)}{\sum_{l=1,l\neq k}^{K} (H_{B,k}[t] + \sum_{n=1}^{N} \sum_{i=1}^{I_n} \kappa \delta_{i_n}[t] g_l[t] \mathbf{h}_{ab}^T[t] \mathbf{H}'[t] \mathbf{h}_{ab}[t] + \sigma^2)} \right) \tag{22}$$

## IV. SOLUTION APPROACH

### A. ARIS Deployment Problem

For the given ARIS reflective elements on/off states $\boldsymbol{\Delta}$, phase shift values $\boldsymbol{\Theta}$ and power control $\boldsymbol{g}$, the sub-problem **P1** can be represented as follows:

$$\textbf{P1:} \max_{\boldsymbol{q}} \quad \mathcal{E}(\boldsymbol{q}) \tag{20a}$$

$$\text{s.t.} \quad r_k[t] \geq r_k^{\min}[t], \forall k \in \mathcal{K}, \forall n \in \mathcal{N}, \forall i \in \mathcal{I}_n, \forall t \in \mathcal{T}, \tag{20b}$$

$$||\boldsymbol{q}_i[t] - \boldsymbol{q}_j[t]||^2 \geq d_{\min}[t], \forall i, j \in \mathcal{N}, i \neq j, \forall t \in \mathcal{T}. \tag{20c}$$

The objective function of sub-problem **P1** remains non-concave since $\mathbf{h}_{n,k}[t]$ and $\mathbf{h}_{B,n}[t]$ are complex and non-linear with respect to ARIS deployment $\boldsymbol{q}_n$. To handle this, we use the approximation algorithm for $\mathbf{h}_{n,k}[t]$ and $\mathbf{h}_{B,n}[t]$. Then, we rewrite our sub-problem **P1** as follows:

$$\max_{\boldsymbol{q}} \quad \dot{\mathcal{E}}(\boldsymbol{q}) \tag{21a}$$

$$\text{s.t.} \quad r_k[t] \geq r_k^{\min}[t], \forall k \in \mathcal{K}, \forall n \in \mathcal{N}, \forall i \in \mathcal{I}_n, \forall t \in \mathcal{T}, \tag{21b}$$

$$||\boldsymbol{q}_i[t] - \boldsymbol{q}_j[t]||^2 \geq d_{\min}[t], \forall i, j \in \mathcal{N}, i \neq j, \forall t \in \mathcal{T}, \tag{21c}$$

where $\dot{\mathcal{E}}(\boldsymbol{q_n})$ is shown in (22), and

$$\mathbf{h}_{ab}[t] = \left[ \sqrt{(d_{n,k}[t])^{-\alpha}}, \sqrt{(d_{B,n}[t])^{-\alpha}} \right]^T,$$

$$\mathbf{H}'[t] = \left[ H_{B,k}^H, (\mathbf{h}_{n,k}^{(\hat{i}-1)}[t])^H \boldsymbol{\Theta}_n[t] \mathbf{h}_{B,n}^{(\hat{i}-1)}[t] \right]$$

$$\left[ H_{B,k}^H, (\mathbf{h}_{n,k}^{(\hat{i}-1)}[t])^H \boldsymbol{\Theta}_n[t] \mathbf{h}_{B,n}^{(\hat{i}-1)}[t] \right]^H.$$

$$\ddot{\mathcal{E}}(\ddot{\boldsymbol{r}}) = \sum_{k=1}^{K} W \log_2 \left( 1 + \frac{\left( H_{B,k}[t] + \sum_{n=1}^{N} \sum_{i=1}^{I_n} \kappa \delta_{i_n}[t] g_k[t] \ddot{r}[t] \right)}{\sum_{l=1, l \neq k}^{K} (H_{B,k}[t] + \sum_{n=1}^{N} \sum_{i=1}^{I_n} \kappa \delta_{i_n}[t] g_l[t] \ddot{r}[t] + \sigma^2)} \right) \tag{24}$$

Next, we introduce the slack variables $\boldsymbol{a} = \{a[t]\}_{t=1}^{T}$, $\boldsymbol{b} = \{b[t]\}_{t=1}^{T}$, and $\ddot{\boldsymbol{r}} = \{\ddot{r}[t]\}_{t=1}^{T}$, and the problem (21) is transformed into the following problem as

$$\max_{\boldsymbol{q}, \boldsymbol{a}, \boldsymbol{b}, \ddot{\boldsymbol{r}}} \quad \ddot{\mathcal{E}}(\ddot{\boldsymbol{r}}) \tag{23a}$$

$$\text{s.t.} \quad 0 < a[t] \leq \sqrt{(d_{B,n}[t])^{-\alpha}}, \forall t \in \mathcal{T}, \tag{23b}$$

$$0 < b[t] \leq \sqrt{(d_{n,k}[t])^{-\alpha}}, \forall t \in \mathcal{T}, \tag{23c}$$

$$\tilde{\boldsymbol{h}}_{ab}^{T}[t] \boldsymbol{H}'[t] \tilde{\boldsymbol{h}}_{ab}[t] \geq \ddot{r}[t], \forall t \in \mathcal{T}, \tag{23d}$$

$$(20b), (20c). \tag{23e}$$

where $\ddot{\mathcal{E}}(\ddot{\boldsymbol{r}})$ is given in (24), and $\tilde{\boldsymbol{h}}_{ab} = [a[t], b[t]]^{T}$. In order to simplify the derivations, we expand (23b) and (23c) as follows [25]:

$$x_B^2 + x_n[t]^2 + y_B^2 + y_n[t]^2 - 2x_B x_n[t] - 2y_B y_n[t] + (z_B - z_n[t])^2 - (a[t])^{-\frac{4}{\alpha}} \leq 0, \tag{25}$$

$$x_n[t]^2 + x_k[t]^2 + y_n[t]^2 + y_k[t]^2 - 2x_n[t] x_k[t] - 2y_n[t] y_k[t] + (z_n[t] - z_k[t])^2 - (b[t])^{-\frac{4}{\alpha}} \leq 0. \tag{26}$$

Still it is discovered that (25) and (26) are in non-convex feasible regions. Therefore, we apply the SCA method to solve this non-convexity. The SCA approach is advantageous because it allows for the replacement of the original non-convex function with simpler surrogates to achieve a suboptimal solution [36]. Firstly, to obtain the global upper bound for the concave function, we first utilize the first-order Taylor expansion to find the linear approximation of the function. To do so, firstly, (23a) can be transformed into the difference of two concave functions as follows [37]:

$$\ddot{\mathcal{E}}(\ddot{\boldsymbol{r}}) \approx \hat{h}(\ddot{\boldsymbol{r}}) - \hat{l}(\ddot{\boldsymbol{r}}), \tag{27}$$

where

$$\hat{h}(\ddot{\boldsymbol{r}}) = \sum_{k=1}^{K} \log_2 \left( h_{B,k}[t] + \sum_{n=1}^{N} \sum_{i=1}^{I_n} \kappa \delta_{i_n}[t] g_k[t] \ddot{r}[t] + \sigma^2 \right), \tag{28}$$

and

$$\hat{l}(\ddot{\boldsymbol{r}}) = \sum_{l=1, l \neq k}^{K} \log_2 \left( h_{B,k}[t] + \sum_{n=1}^{N} \sum_{i=1}^{I_n} \kappa \delta_{i_n}[t] g_l[t] \ddot{r}[t] + \sigma^2 \right). \tag{29}$$

Both the functions $\hat{h}(\ddot{\boldsymbol{r}})$ and $\hat{l}(\ddot{\boldsymbol{r}})$ are convex. However the difference between them is neither convex nor concave, as represented in (27). Then, we find the feasible solution $\ddot{\boldsymbol{r}}'$ to problem (23) by computing the concave lower bound, i.e. the surrogate function of the non-concave objective, specified in (27). By implementing the first-order Taylor expansion to replace the $\hat{l}(\ddot{\boldsymbol{r}})$, we can construct its lower bound as follows:

$$\ddot{\mathcal{E}}(\ddot{\boldsymbol{r}}, \ddot{\boldsymbol{r}}') = \hat{h}(\ddot{\boldsymbol{r}}) - \hat{\tilde{l}}((\ddot{\boldsymbol{r}}, \ddot{\boldsymbol{r}}')), \tag{30}$$

where

$$\hat{\tilde{l}}((\ddot{\boldsymbol{r}}, \ddot{\boldsymbol{r}}')) \triangleq \hat{l}(\ddot{\boldsymbol{r}}') - \nabla \hat{l}(\ddot{\boldsymbol{r}}')(\ddot{\boldsymbol{r}} - \ddot{\boldsymbol{r}}'), \tag{31}$$

where $\nabla \hat{l}(\ddot{\boldsymbol{r}}')$ is the gradient of the $\hat{l}(\ddot{\boldsymbol{r}})$ at the given point $\ddot{\boldsymbol{r}}'$, and $\hat{\tilde{l}}((\ddot{\boldsymbol{r}}, \ddot{\boldsymbol{r}}'))$ represents the first-order Taylor's approximation of $\hat{l}(\ddot{\boldsymbol{r}})$ near $\ddot{\boldsymbol{r}}'$ in the feasible area of the solution space. The gradient for ARIS $n$ can be expressed as follows:

$$\nabla_n \hat{l}(\ddot{\boldsymbol{r}}') = \frac{\partial \hat{l}(\ddot{\boldsymbol{r}}')}{\partial \ddot{r}'} = \frac{1}{\ln 2} \sum_{l=1, l \neq k}^{K} \frac{\sum_{n=1}^{N} \sum_{i=1}^{I_n} \kappa \delta_{i_n}[t] g_l[t]}{H_{B,k}[t] + \sum_{n=1}^{N} \sum_{i=1}^{I_n} \kappa \delta_{i_n}[t] g_l[t] \ddot{r}[t] + \sigma^2}. \tag{32}$$

The surrogate function given in (31) is concave. Next, the upper bound of function $\hat{l}(\ddot{\boldsymbol{r}})$ may also be found using the first-order Taylor's expansion.

**Lemma 1.** *The first-order Taylor approximation provides the global upper bound of a concave function or the global lowest bound of a convex function.*

*Proof.* Initially, we define the first-order Taylor series as follows:

$$f(x_0) + f'(x_0)(x - x_0). \tag{33}$$

Afterwards, we have

$$\hat{l}(\ddot{\boldsymbol{r}}) \leq \hat{l}(\ddot{\boldsymbol{r}}') + \nabla \hat{l}(\ddot{\boldsymbol{r}}')(\ddot{\boldsymbol{r}} - \ddot{\boldsymbol{r}}'). \tag{34}$$

Therefore, we can derive the observations by examining (27), (30), and (34) as follows:

$$\begin{aligned} \ddot{\mathcal{E}}(\ddot{\boldsymbol{r}}) &= \hat{h}(\ddot{\boldsymbol{r}}) - \hat{l}(\ddot{\boldsymbol{r}}) \\ &\geq \hat{h}(\ddot{\boldsymbol{r}}) - \left\{ \hat{l}(\ddot{\boldsymbol{r}}') + \nabla \hat{l}(\ddot{\boldsymbol{r}}')(\ddot{\boldsymbol{r}} - \ddot{\boldsymbol{r}}') \right\} \\ &\geq \hat{h}(\ddot{\boldsymbol{r}}) - \hat{l}(\ddot{\boldsymbol{r}}') - \nabla \hat{l}(\ddot{\boldsymbol{r}}')(\ddot{\boldsymbol{r}} - \ddot{\boldsymbol{r}}') \\ &= \ddot{\mathcal{E}}(\ddot{\boldsymbol{r}}, \ddot{\boldsymbol{r}}'), \end{aligned} \tag{35}$$

where (35) denotes that the surrogate function provides the lower bound of the original function. As a result, at point $\ddot{\boldsymbol{r}}'$, i.e., $\ddot{\mathcal{E}}(\ddot{\boldsymbol{r}}, \ddot{\boldsymbol{r}}')|_{\ddot{r}=\ddot{r}'} = \ddot{\mathcal{E}}(\ddot{\boldsymbol{r}}')$, the two functions are tangent to each other.

Thereby, our objective function of sub-problem (23) has the lower bound function as obtained in (35). □

Consequently, we replace our objective function in problem (23) which is non-convex, by its surrogates as presented in (30). Furthermore, we take the first-order Taylor expansions of $(a[t])^{-\frac{4}{\alpha}}, (b[t])^{-\frac{4}{\alpha}}$, and $\tilde{\mathbf{h}}_{ab}^T[t]\mathbf{H}'[t]\tilde{\mathbf{h}}_{ab}[t]$ at the given feasible points $\boldsymbol{a_0} = \{a_0[t]\}_{t=1}^T$, $\boldsymbol{b_0} = \{b_0[t]\}_{t=1}^T$, and $\tilde{\mathbf{H}}_{\mathbf{0}ab} = \{\tilde{\mathbf{h}}_{\mathbf{0}ab}[t]\}_{t=1}^T$ are expressed as follows:

$$(a[t])^{-\frac{4}{\alpha}} \geq (a_0[t])^{-\frac{4}{\alpha}} - \frac{4}{\alpha}(a_0[t])^{-\frac{4}{\alpha}-1}(a[t] - a_0[t]), \tag{36}$$

$$(b[t])^{-\frac{4}{\alpha}} \geq (b_0[t])^{-\frac{4}{\alpha}} - \frac{4}{\alpha}(b_0[t])^{-\frac{4}{\alpha}-1}(b[t] - b_0[t]), \tag{37}$$

$$\tilde{\mathbf{h}}_{ab}^T[t]\mathbf{H}'[t]\tilde{\mathbf{h}}_{ab}[t] \geq -\tilde{\mathbf{h}}_{\mathbf{0}ab}^T[t]\mathbf{H}'[t]\tilde{\mathbf{h}}_{\mathbf{0}ab}[t] + 2\Re\left[\tilde{\mathbf{h}}_{\mathbf{0}ab}^T[t]\mathbf{H}'[t]\tilde{\mathbf{h}}_{ab}[t]\right]. \tag{38}$$

By combining (25) and (36), (26) and (37), we get

$$x_B^2 + x_n[t]^2 + y_B^2 + y_n[t]^2 - 2x_B x_n[t] - 2y_B y_n[t]+ \tag{39}$$
$$(z_B - z_n[t])^2 - (1 + \frac{4}{\alpha})(a_0[t])^{-\frac{4}{\alpha}} + \frac{4}{\alpha}(a_0[t])^{-\frac{4}{\alpha}-1}a[t] \leq 0,$$

$$x_n[t]^2 + x_k[t]^2 + y_n[t]^2 + y_k[t]^2 - 2x_n[t]x_k[t] - 2y_n[t]y_k[t]+ \tag{40}$$
$$(z_n[t] - z_k[t])^2 - (1 + \frac{4}{\alpha})(b_0[t])^{-\frac{4}{\alpha}} + \frac{4}{\alpha}(b_0[t])^{-\frac{4}{\alpha}-1}b[t] \leq 0.$$

Similarly, we apply the first-order Taylor expansion to convert $||\boldsymbol{q}_i[t] - \boldsymbol{q}_j[t]||^2$ in constraint (20c) to a linear function since it is a convex function with respect to $\boldsymbol{q}_i$ and $\boldsymbol{q}_j$. This can be expressed as follows:

$$||\boldsymbol{q}_i[t] - \boldsymbol{q}_j[t]||^2 \geq 2(\boldsymbol{q}_i[t-1] - \boldsymbol{q}_j[t-1])^T(\boldsymbol{q}_i[t] - \boldsymbol{q}_j[t]) - ||\boldsymbol{q}_i[t-1] - \boldsymbol{q}_j[t-1]||^2. \tag{41}$$

Afterwards, we can denote the above equation as follows:

$$G_0[t-1](\boldsymbol{q}_i[t] - \boldsymbol{q}_j[t]) \triangleq 2(\boldsymbol{q}_i[t-1] - \boldsymbol{q}_j[t-1])^T(\boldsymbol{q}_i[t] - \boldsymbol{q}_j[t]) - ||\boldsymbol{q}_i[t-1] - \boldsymbol{q}_j[t-1]||^2. \tag{42}$$

Finally, we can substitute (42) into (20c), and problem (23) can be rewritten as follows:

$$\min_{\boldsymbol{q}, \boldsymbol{a}, \boldsymbol{b}, \ddot{\boldsymbol{r}}} \quad -\ddot{\mathcal{E}}(\ddot{\boldsymbol{r}}, \ddot{\boldsymbol{r}}') \tag{43a}$$

$$\text{s.t.} \quad \ddot{r}[t] + \tilde{\boldsymbol{h}}_{\mathbf{0}ab}^T[t]\boldsymbol{H}'[t]\tilde{\boldsymbol{h}}_{\mathbf{0}ab}[t] - 2\Re\left[\tilde{\boldsymbol{h}}_{\mathbf{0}ab}^T[t]\boldsymbol{H}'[t]\tilde{\boldsymbol{h}}_{ab}[t]\right] \leq 0, \forall t \in \mathcal{T}, \tag{43b}$$

$$G_0[t-1](\boldsymbol{q}_i[t] - \boldsymbol{q}_j[t]) \geq d_{\min}[t], \forall i, j \in \mathcal{N}, i \neq j, \forall t \in \mathcal{T}, \tag{43c}$$

$$(20b), (39), (40). \tag{43d}$$

---

**Algorithm 1** SCA algorithm for ARIS deployment

---

**Input:** Initial feasible points $\{\boldsymbol{q}^0, \boldsymbol{a}^0, \boldsymbol{b}^0\}$, $r_k^{\min}[t]$, $d_{\min}[t]$, iteration index $\hat{i} = 0$, $\hat{i}_{max}$, stopping

　　criterion $\varepsilon_1$.

　1: **repeat**

　2:　　Set $\hat{i} \leftarrow \hat{i} + 1$.

　3:　　Update $\boldsymbol{q}^{\hat{i}}, \boldsymbol{a}^{\hat{i}}, \boldsymbol{b}^{\hat{i}}$ with given $\boldsymbol{q}^{\hat{i}-1}, \boldsymbol{a}^{\hat{i}-1}, \boldsymbol{b}^{\hat{i}-1}$.

　4:　　Acquire $\ddot{\mathcal{E}}(\ddot{\boldsymbol{r}}, \ddot{\boldsymbol{r}}') = \hat{h}(\ddot{\boldsymbol{r}}) - \hat{\hat{l}}((\ddot{\boldsymbol{r}}, \ddot{\boldsymbol{r}}'))$ based on (30).

　5:　　Solve (43) to obtain $\ddot{\boldsymbol{r}}^{\hat{i}}$.

　6: **until** $|\ddot{\mathcal{E}}(\ddot{\boldsymbol{r}}^{\hat{i}}) - \ddot{\mathcal{E}}(\ddot{\boldsymbol{r}}^{\hat{i}-1})| \leq \varepsilon_1$ or $\hat{i} > \hat{i}_{max}$.

**Output:** Optimal ARIS deployment $\boldsymbol{q}^*$.

---

Problem (43) becomes a convex optimization problem, which we can solve by using CVXPY solver in python programming. The overall algorithm of the SCA method is shown in Algorithm 1.

### B. Joint ARIS Reflective Elements On/off States and Phase Shift Problem

For the given ARIS deployment $\boldsymbol{q}$ and power control $\boldsymbol{g}$, the sub-problem **P2** can be represented as follows:

$$\textbf{P2:} \max_{\boldsymbol{\Delta}, \boldsymbol{\Theta}} \quad \mathcal{E}(\boldsymbol{\Delta}, \boldsymbol{\Theta}) \tag{44a}$$

$$\text{s.t.} \quad r_k[t] \geq r_k^{\min}[t], \forall k \in \mathcal{K}, \forall n \in \mathcal{N}, \forall i \in \mathcal{I}_n, \forall t \in \mathcal{T}, \tag{44b}$$

$$0 \leq \theta_{i_n}[t] < 2\pi, \forall n \in \mathcal{N}, \forall i \in \mathcal{I}_n, \forall t \in \mathcal{T}, \tag{44c}$$

$$\delta_{i_n}[t] \in \{0, 1\}, \forall n \in \mathcal{N}, \forall i \in \mathcal{I}_n, \forall t \in \mathcal{T}. \tag{44d}$$

This problem is still mixed-integer, non-convex, and quite challenging to solve in polynomial time, since the information of the environment is unknown. Moreover, the real-time ARIS reflective elements on/off states requires extensive computation and hardware cost, and conventional optimization methods cannot be applied. The exhaustive search method can be used to find the optimal solution, however it is impractical for large-scale networks. Due to these reasons, we propose DRL approach to solve sub-problem **P2**. The reason we do not apply DRL for the whole optimization problem is that the action spaces combined for all ARIS deployment, ARIS reflective elements on/off states, phase shift, and power control matrices will be too large and demands

high computational cost. Here, we implement Actor-Critic Proximal Policy Optimization (AC-PPO) [38] as it always provides an improved policy by using data that are currently accessible by the agent and thereby ensuring data efficiency and reliable performance. It could also be utilised in the environments where action spaces are discrete or continuous. Typically, since DRL is interpreted as Markov Decision Process (MDP), we first need to define state space $\mathcal{S}$, action space $\mathcal{A}$ and reward $\tilde{R}$.

*1) State Space:* For each state at time $t$, $s_t \in \mathcal{S}$ can be expressed as the tuples of the users' locations and ARISs' locations, the channel gain of the direct link, the channel gain of the ARIS-UE link and BS-ARIS link, and power control at time $t$, respectively, and can be represented by $s_t = \{q_k[t], q_n[t], \mathbf{H}_{B,k}[t], \mathbf{h}_{n,k}[t], \mathbf{h}_{B,n}[t], g_k[t], \forall k \in \mathcal{K}, \forall n \in \mathcal{N}, \forall i \in \mathcal{I}_n\}$.

*2) Action space:* The action at time $t$, $a_t \in \mathcal{A}$ contains the combination of the ARIS reflective elements on/off states variable $\delta_{i_n}[t]$, and phase shift values $\theta_{i_n}[t]$ at time $t$, and can be denoted as $a_t = \{\delta_{i_n}[t], \theta_{i_n}[t], \forall n \in \mathcal{N}, \forall i \in \mathcal{I}_n\}$

*3) Reward:* Since the goal of our system is to maximize the energy efficiency, our reward function is defined as

$$\tilde{R}_t(s_t|a_t) = \begin{cases} -1, & \text{if } \sum_{k=1}^{K} r_k[t] < r_k^{\min}[t], \\ \mathcal{E}(\mathbf{\Delta}, \mathbf{\Theta}), & \text{otherwise.} \end{cases} \tag{45}$$

As shown in Fig 3, in our AC-PPO algorithm, the states information, $s_t$ from the environment is obtained by the agent at the BS, and the agent observes and monitors the status of the location of the users and ARISs, channel gain of the links, and power control for each user. The agent includes the actor model and the critic model [39]. The actor model has the stochastic policy model $\pi_\psi(a_t|s_t)$ with its own parameter $\psi$ and learns to take which action under the observation of the input states. The policy $\pi_\psi(a_t|s_t)$ takes the observed states $s_t$ from the environment as an input and suggests actions $a_t$ to take as an output, and calculates the immediate reward $\tilde{R}_t(s_t|a_t)$ depending on the action taken. The reward then provides as feedback to the agent, and the new state information $s(t+1)$ is obtained. Taking into account of the requirements for the users, under given policy $\pi_\psi(a_t|s_t)$ and reward function $\tilde{R}_t(s_t|a_t)$, the cumulative discounted reward function at time $t$ can be denoted as follows:

$$V^{\pi_\psi}(s_t) = \hat{\mathbb{E}}_t \left[ \sum_{t'=t}^{T-1} \xi^{t'-t} \tilde{R}_{t'}(s_{t'}|a_{t'}) \right], \forall s_t \in \mathcal{S}, \tag{46}$$

where $0 < \xi < 1$ is the discount factor to prevent the total reward from reaching to infinity.

Fig. 3: AC-PPO Algorithm for joint ARIS reflective elements on/off states and phase shift.

Moreover, the critic model contains the advantage function, $\hat{A}_t$ which is the estimate of the relative value of the selected action in the current state is defined as [40]:

$$\hat{A}_t = V^{\pi_\psi}(s_t) - b(s_t), \forall s \in \mathcal{S}, \tag{47}$$

where $b(s_t)$ is the baseline estimate value function which provides the estimate of the discounted return starting from the current state $s_t$.

The surrogate objective function of AC-PPO is to find the policy that maximizes the total rewards from the environment and can be expressed as follows [38]:

$$L^{CLIP}(\psi) = \hat{\mathbb{E}}_t \left[ \min \left( r_t(\psi)\hat{A}_t, \text{clip}(r_t(\psi), 1 - \epsilon, 1 + \epsilon)\hat{A}_t \right) \right], \tag{48}$$

where

$$r_t(\psi) = \frac{\pi_\psi(a_t|s_t)}{\pi_{\psi_{old}}(a_t|s_t)},$$

means the probability ratio. Given the states and actions, $r_t(\psi) > 1$ is the action is more plausible currently than it was in the old version of the policy, and $0 < r_t(\psi) < 1$ if it is less plausible, and $\epsilon$ is the clipping parameter. The clipping part of the objective function ensures that the PPO does not always favor actions with positive advantage and/or consistently avoid actions with negative advantage. The overall algorithm of the AC-PPO algorithm is described in Algorithm 2.

---

**Algorithm 2** AC-PPO algorithm for ARIS reflective elements on/off states and phase shift

---

**Input:** Network states $s_t$, learning rate, discount factor $\xi$, clipping parameter $\epsilon$;

1: **Initialization** Base policy $\pi_\psi(a_t|s_t)$ with random parameters $\psi$ and clipping parameter $\epsilon$ and initial value function $V^{\pi_\psi}(s_t)$

2: **for** $k \in K$ **do**

3:     **for** each episode $\hat{t} \in \hat{T}$ **do**

4:         Collect the network observations: ARIS deployments $\boldsymbol{q}$ from Algorithm 1 and power control $\boldsymbol{g}$ from Algorithm 3 to achieve the initial state $s_0$

5:         **for** each $t \in T$ **do**

6:             Forward the network states $s_t \in \mathcal{S}$ to the AC-PPO algorithm

7:             Observe the input states $s_t$ and run the actor network

8:             Select action $a_t \in \mathcal{A}$ based on policy $\pi_\psi(a_t|s_t)$

9:             Obtain the reward $\tilde{R}_t(s_t|a_t)$ and $s_{t+1}$

10:            Calculate the probability ratio, $r_t$

11:            Compute $\hat{A}_t$ based on current $V^{\pi_\psi}(s_t)$ at the critic network according to (47)

12:            Compute $L^{CLIP}(\psi)$ according to (48)

13:            Update $\pi_{\psi_{old}} \leftarrow \pi_\psi$

14:         **end for**

15:     **end for**

16: **end for**

**Output:** Optimal AC-PPO network with $\boldsymbol{\Delta}^*, \boldsymbol{\Theta}^*$.

---

## C. Power Control Problem

For the fixed ARIS deployment $\boldsymbol{q}$, ARIS reflective elements on/off states $\boldsymbol{\Delta}$, and phase shift $\boldsymbol{\Theta}$, sub-problem **P3** can be represented as follows:

$$\textbf{P3:} \max_{\boldsymbol{g}} \quad \mathcal{E}(\boldsymbol{g}) \tag{49a}$$

$$\text{s.t.} \quad r_k[t] \geq r_k^{\min}[t], \forall k \in \mathcal{K}, \forall n \in \mathcal{N}, \forall i \in \mathcal{I}_n, \forall t \in \mathcal{T}, \tag{49b}$$

$$\text{tr}(\boldsymbol{g}[t]^H \boldsymbol{g}[t]) \leq P_{\max}, \forall t \in \mathcal{T}. \tag{49c}$$

Sub-problem **P3** is still a non-convex and NP-hard problem due to constraint (49b). Therefore, it is challenging to obtain the solutions in the polynomial time. Therefore, we adopt Whale Optimization Algorithm (WOA) to solve sub-problem **P3**. The WOA is a meta-heuristic algorithm which mimics the whales hunting strategy. The WOA has substantial advantages. First, unlike gradient-based algorithms, which involve computing and updating the gradients and step size throughout every iteration of the optimization process, WOA allows for such computation to be relaxed. Second, WOA is not influenced by the initial feasible solutions, which might have a significant impact on the convergence. Therefore, it has recently gained popularity among research community due to it being efficient optimizer. The WOA algorithm includes two states: 1) the exploitation state (the encircling prey method and spiral bubble-net attacking method), and 2) the exploration state (the searching prey method). The detail explanation of each state can be further described in the following subsections [41]–[43].

*1) Exploitation State:* The exploitation state of WOA includes two fundamental methods: the encircling prey method, and the spiral bubble-net attack method, which are discussed as follows:

**Encircling Prey Method**. Once the whales detect the location of their preys, they encircle them. Theoretically, the location of the prey is unknown in the search space, therefore, WOA assumes that the current best search agent is the target prey (optimum or close to optimum). The other whales (search agents) update their locations towards to the best search agent. This behaviour can be mathematically implemented as follows [41]:

$$\vec{D} = \left| \vec{C} \cdot \vec{g}^*(\hat{j}) - \vec{g}(\hat{j}) \right|, \tag{50}$$

$$\vec{g}(\hat{j} + 1) = \vec{g}^*(\hat{j}) - \vec{A} \cdot \vec{D}, \tag{51}$$

where $\vec{g}^*$ is the location of the best search agent, $\hat{j}$ is the current iteration, $|\cdot|$ is the absolute value. $\vec{C}$ and $\vec{A}$ are coefficient vectors, and are computed as follows:

$$\vec{A} = 2\vec{a} \cdot \vec{r} - \vec{a}, \tag{52}$$

$$\vec{C} = 2 \cdot \vec{r}, \tag{53}$$

where $\vec{r}$ is the random vector between $0$ to $1$, and $\vec{a}$ is the control parameter vector linearly declining from $2$ to $0$ over the iterations, both in exploitation and exploration states. The aim of (52) and (53) is to balance between exploitation and exploration. When $A \geq 1$, WOA will perform exploration, and exploitation is done when $A < 1$.

**Spiral Bubble-net Attack Method**. This method combines both shrinking encircling mechanism and spiral movement mechanism of whales. Its purpose is to update the new location to fall between the current agent's location and the best search agent. To mimic the helical shape movement of the whales, the equation can be expressed as

$$\vec{D'} = \left| \vec{g}^*(\hat{j}) - \vec{g}(\hat{j}) \right|, \tag{54}$$

$$\vec{g}(\hat{j}+1) = \vec{D'} \cdot e^{bj} \cdot \cos(2\pi l) + \vec{g}^*(\hat{j}), \tag{55}$$

where $\vec{D'}$ indicates the distance between the current search agent and the target prey. Moreover, $b$ is the constant for defining the shape of the logarithmic spiral, and $l$ is the random number between $-1$ and $1$. Here, coefficient vector $\vec{A}$ is updated by setting random values in $[-1, 1]$.

Conventionally, once the whales locate the prey, they approach it using either shrinking encircling method or spiral bubble-net method synchronously. To imitate this synchronous behaviour, we set the $50\%$ probability to choose between these two methods to update the location of the whales for the optimization. Mathematically, it can be modeled as follows:

$$\vec{g}(\hat{j}+1) = \begin{cases} \vec{g}^*(\hat{j}) - \vec{A} \cdot \vec{D}, & \text{if } p < 0.5, \\ \vec{D'} \cdot e^{bj} \cdot \cos(2\pi l) + \vec{g}^*(\hat{j}), & \text{if } p \geq 0.5. \end{cases} \tag{56}$$

where $p = [0, 1]$ is the random number to represent the probability to choose between two mechanisms. When $p < 0.5$, WOA chooses the shrinking encircling mechanism, and if $p \geq 0.5$, WOA chooses the sprial movement mechanism.

*2) Exploration State:* : The exploration state of WOA includes the searching for prey method. This state is necessary to prevent the solution from being trapped at the local optimum, and failing to achieve the global optimum.

**Searching for Prey Method**. This approach is similar to encircling prey method in exploitation state, but instead of claiming the location of best search agent, and here, a random location is

---

**Algorithm 3** WOA for power control

---

**Input:** Current power control $\boldsymbol{g}$, given $\boldsymbol{q}$, $\boldsymbol{\Delta}$, and $\boldsymbol{\Theta}$;

1: **Initialization** At iteration $\hat{j} = 1$, initialize the total number of whale population $g_u$, where $u = \{1, \ldots, U\}$, and maximum number of iteration $\hat{j}_{\mathrm{max}}$.

2: According to (59), calculate the fitness of the search agents $g_u$ and identify the best search agent $\vec{g}^*(0)$.

3: **repeat**

4:     **for** $u \leftarrow 1$ to $U$ (the number of whales) **do**

5:         Update $a, A, C, l$ and $p$.

6:         **if** $p < 0.5$ **then**

7:             **if** $|A| < 1$ **then**

8:                 Update $\vec{D}$ by (50) and $\vec{g}$ by (51).

9:             **else**

10:                 Select a random $\vec{g}_{\mathrm{rand}}$ and update $\vec{D}$ by (57).

11:                 Update the location $\vec{g}$ by (58).

12:             **end if**

13:         **else**

14:             Update $\vec{D}$ by (54) and $\vec{g}$ by (55).

15:         **end if**

16:     **end for**

17:     Calculate the fitness of each search agent by (59).

18:     Update the location of the best search agent $\vec{g}^*(\hat{j})$.

19:     Update $\hat{j} \leftarrow \hat{j} + 1$.

20: **until** $\hat{j} > \hat{j}_{\mathrm{max}}$

**Output:** Optimal power control $\boldsymbol{g}^*$.

---

selected to update the locations of other search agents. It can be mathematically represented as follows:

$$\vec{D} = \left| \vec{C} \cdot \vec{g}_{\mathrm{rand}}(\hat{j}) - \vec{g}(\hat{j}) \right|, \tag{57}$$

$$\vec{g}(\hat{j} + 1) = \vec{g}_{\mathrm{rand}}(\hat{j}) - \vec{A} \cdot \vec{D}, \tag{58}$$

---

**Algorithm 4** Proposed joint ARIS deployment, ARIS reflective elements on/off states, phase shift, power control optimization algorithm

---

1: **Initialization:** At $\tau = 0$, initialize the variables, $\boldsymbol{q}(0), \boldsymbol{\Delta}(0), \boldsymbol{\Theta}(0), \boldsymbol{g}(0)$;

2: **repeat**

3:   By applying **Algorithm 1**, solve problem **P1** for given $\boldsymbol{\Delta}(\tau), \boldsymbol{\Theta}(\tau), \boldsymbol{g}(\tau)$ to obtain $\boldsymbol{q}(\tau + 1)$.

4:   By applying **Algorithm 2**, solve problem **P2** for given $\boldsymbol{q}(\tau + 1), \boldsymbol{g}(\tau)$ to obtain $\boldsymbol{\Delta}(\tau + 1), \boldsymbol{\Theta}(\tau + 1)$.

5:   By applying **Algorithm 3**, solve problem **P3** for given $\boldsymbol{q}(\tau + 1), \boldsymbol{\Delta}(\tau + 1), \boldsymbol{\Theta}(\tau + 1)$ to obtain $\boldsymbol{g}(\tau + 1)$.

6:   Update $\tau \leftarrow \tau + 1$.

7: **until** objective value (18) reaches convergence.

---

where $\vec{g}_{\mathrm{rand}}(\hat{j})$ is the location of the search agent randomly selected from the search space.

Since WOA algorithm is designed only for unconstrained optimization, we apply the penalty method to our sub-problem **P3** in order to deal with the minimum achievable date rate constraint (49b) in the problem [43]. In our scenario, UEs are considered as a search agent, and the power control of the BS $\boldsymbol{g}$ represents the location of the search agents. At each iteration $\hat{j}$, the power control $\boldsymbol{g}$ can be updated by either the encircling prey method, spiral bubble-net attack method, or searching for prey method. The fitness function of our problem which chooses the optimal search agent can be expressed as follows:

$$\text{Fitness}(\boldsymbol{g}) = -\frac{R(\boldsymbol{g})}{P(\boldsymbol{g})} + \varpi \sum_{k=1}^{K} F_k(f_k(\boldsymbol{g})) f_k^2(\boldsymbol{g}), \tag{59}$$

where $f_k(\boldsymbol{g}) = r_k^{\min}[t] - r_k[t]$ is the inequality function, and $\varpi$ is the penalty factor coefficient. Since our sub-problem **P3** is the maximization problem, we add the negative sign ahead of the objective function to convert into a minimization problem. The index function $F_k(f_k(\boldsymbol{g})) = 1$ if $f_k(\boldsymbol{g}) < 0$, and $F_k(f_k(\boldsymbol{g})) = 0$ if $f_k(\boldsymbol{g}) \geq 0$. The pseudo-code of our WOA based power control can be described as in Algorithm 3.

### D. Overall Algorithm Complexity Analysis

The overall iterative algorithm for solving our optimization problem (18) is described in Algorithm 4 with the aforementioned proposed solutions to three sub-problems. According to the results in [25], [37] and [43], the complexity of our solutions can be obtained by each algorithm for each sub-problem. For ARIS deployment sub-problem, the SCA is adopted as in Algorithm 1. Since there are $K$ users, the computational complexity of the SCA method is

obtained as $\mathcal{O}_1\left(K^{3.5}\log(1/\varepsilon_1)\right)$ where $\varepsilon_1$ is the variable to control the accuracy of the SCA algorithm. For the AC-PPO algorithm for ARIS reflective elements on/off states and phase shift as in Algorithm 2, the computational complexity is $\mathcal{O}_2\left(a^2K\right)$, where $a \in \mathcal{A}$ is the total number of actions taken by the agent. With WOA for power control as in Algorithm 3, the computational complexity is $\mathcal{O}_3\left(\hat{J}U(m+K)\right)$, where $\hat{J}$ is the number of iterations for WOA, $U = 30$ denotes the number of whale populations, and $m$ represents the number of inequality constraints in sub-problem **P3**. Henceforth, the overall computational complexity for solving (18) can be acquired as $\mathcal{O}\left(\hat{\tau}K^{3.5}\log(1/\varepsilon_1) + \hat{\tau}a^2K + \hat{\tau}\hat{J}U(m+K)\right)$, where $\hat{\tau}$ denotes the number of iterations for Algorithm 4.

## V. Performance Evaluation

In this section, we evaluate our proposed technique of energy-efficient multiple ARISs-assisted downlink communication system via numerical analysis. The network design comprises of 12 UEs uniformly distributed within 100 m × 100 m square region and the BS with 15 multiple antennas located at the center of the coverage area. There are 4 ARISs to support communication, and each ARIS is integrated with 10 reflective elements. The ARISs can hover at a maximum altitude of 100 m. The simulation parameters can be observed in Table I. To evaluate our proposed algorithm, we compare our method with four benchmark schemes, which are explained as follows:

- *Single-ARIS:* In this scheme, we implement a single ARIS instead of using multiple ARISs to support the downlink communications from BS. The optimization problem is then solved by using our proposed algorithms.
- *ARIS (NPS):* In this approach, we deploy 4 ARISs with the fixed phase shifts. The ARIS deployment problem is solved by SCA, the ARIS reflective elements on/off states problem is solved by AC-PPO, and transmit power allocation problem is solved by WOA alternatively.
- *Random:* In this design, we randomly deploy the ARISs, fixed the reflective elements ON/OFF states, and fixed the transmit power of the BS.
- *UAV-Relay:* In this method, ARIS is not used. Instead, 4 UAVs are deployed as relays and the incident signal is linearly processed, and re-transmit them toward the required destination. The optimization problem is then solved by using our proposed algorithms.

Fig. 4 compares the average sum-rate of users with our proposed DRL-based algorithm towards the above-mentioned benchmark schemes. In all circumstances, the average sum-rate increases

TABLE I: Simulation parameters

| Parameter | Value |
|---|---|
| Number of reflective elements on ARIS $n$, $I_n$ | 10 |
| Bandwidth $W$ | 2 MHz |
| Noise power $\sigma^2$ | -174 dBm |
| Path loss exponent $\alpha$ | 4 |
| Channel gain at reference distance $\kappa$ | -40 dBm |
| Rician factor $\hat{R}$ | 10 |
| Circuit power of each RIS element $P_{\text{RIS}}$ | 10 dBm [33] |
| BS power amplifier efficiency $\mu$ | 0.8 [33] |
| Circuit power of each user $P_k^{\text{cir}}$ | 10 dBm [33] |
| Clipping parameter $\epsilon$ | 0.2 |
| Learning rate | 0.0002 |
| Discount factor $\xi$ | 0.9 |
| Mini batch size | 64 |
| Number of episodes | 1,000 |
| Number of time steps | 300,000 |



Fig. 4: Performance comparison of sum-rate for different transmit power.



Fig. 5: Performance comparison of energy efficiency for different transmit power.

as the maximum transmit power rises. Our proposed algorithm outperforms ARIS (NPS) by $24\%$ and single-ARIS by $58\%$, respectively. This demonstrates how the multiple ARISs can achieve better outcomes than a single ARIS since it can provide several paths between the BS and UEs. Our algorithm outperforms most benchmark schemes in average sum-rate except for the UAV-relay scenario. UAV-relay provides the highest performance since it processes and re-transmits

Fig. 6: Performance comparison of cumulative rewards for different $P_{\max}$.



Fig. 7: Performance comparison of cumulative rewards for different learning rate.



Fig. 8: CDF of sum-rate with different number of ARIS reflective elements for RIS and proposed system.



Fig. 9: CDF of energy-efficiency with different number of ARIS.

the incident signal using a dedicated power source. As a consequence, it consumes more energy which can be observed in Fig. 5.

Fig. 5 depicts the comparison of the energy efficiency under different algorithms. The smooth data is demonstrated by the solid curved line, which represents the Savitzky-Golay filter. In all scenarios, energy efficiency increases faster until the transmit power of the BS reaches to $10$ dBm. Since then, the energy efficiency hasn't improved much as the function does not increase monotonically with respect to to transmit power. Our proposed algorithm achieves $72\%$ increase compared to the UAV-relay scenario and $43\%$ increase compared to the single-ARIS scenario.

Fig. 10: CDF of energy-efficiency with different number of UEs.

Next, we evaluate the convergence of our proposed AC-PPO algorithm with different values of $P_{\max}$ ranging from 0 dBm to 40 dBm. As shown in Fig. 6, it can be observed that in all cases, the convergence of our cumulative rewards increases with respect to increase in transmit power. We can see a significant difference in the performance when $P_{\max}$ is low, and the performance difference becomes lesser as $P_{\max}$ becomes higher. This suggests that SINR has significant impact on the overall performance of the cumulative rewards.

Following that, we examine how various learning rates affect on our cumulative rewards, ranging from the set of $\{0.02, 0.002, 0.0002, 0.00002\}$. As seen in Fig. 7, a higher learning rate does not enable our cumulative reward to converge faster but provides less performance. Although it takes longer to converge, the learning rate of 0.0002 delivers better performance than 0.002 and 0.00002. In this case, we chose a learning rate of 0.0002 since it produces the highest cumulative rewards for our proposed method.

Furthermore, we compare the spectral efficiency of our proposed multiple ARISs-assisted system to that of multiple RISs-assisted systems. In this approach, we employ 4 RISs on the ground level rather than mounted on the UAVs. Fig. 8 demonstrates the cumulative distribution function (CDF) values of average sum-rates. As seen in Fig. 8, our proposed system achieved $69\%$ performance increase compared to RISs-assisted system. This is because our proposed system takes into account the deployment of UAVs, which provides improved LOS communications between the BS and UEs. Concurrently, we experiment the performance of spectral efficiency with different numbers of reflective elements. In all scenarios, the results show that as the number of reflective elements increases, so do the UEs' average sum-rates. This indicates that the spectral

efficiency will be improved by increasing the number of reflective parts.

Finally, we examine the energy efficiency with various ARIS numbers and different numbers of UEs, respectively. As shown in Fig. 9, we can observe that the energy efficiency improves as the number of ARIS components increases. Moreover, when the number is low, the energy efficiency improvement is significantly more compared to a larger number of ARIS. This indicates that for the small cell network with 12 UEs, we do not need to install a large amount of ARISs. Next, as shown in Fig. 10, we can observe that the energy efficiency almost linearly increases with increasing number of UEs between 6 to 12. We can perceive that more ARISs and more UEs help improve the energy efficiency of the multiple ARISs-assisted downlink communication system.

## VI. CONCLUSION

In this paper, we have studied an energy-efficient multiple ARISs-assisted downlink communication system. To maximize energy efficiency, we formulated a joint ARIS deployment, ARIS reflective elements on/off states, phase shift, and power control problem. As the formulated problem is MINLP and NP hard, we decompose our problem into three sub-problems: ARIS deployment problem, joint reflective elements ON/OFF states and phase shift problem, and power control problem. We then proposed SCA approach, AC-PPO method and WOA to solve our sub-problems, alternatively. Through extensive numerical analysis, we have proved that by integrating multiple ARISs in the downlink communication system, it can significantly outperform several benchmark schemes; especially in spectral efficiency compared to a multiple RISs-assisted scenario and energy efficiency compared to a single ARIS-assisted scenario.

## REFERENCES

[1] U. Cisco, "Cisco annual internet report (2018–2023) white paper," Mar. 2020.

[2] Y. K. Tun, Y. M. Park, N. H. Tran, W. Saad, S. R. Pandey, and C. S. Hong, "Energy-efficient resource management in uav-assisted mobile edge computing," *IEEE Communications Letters*, vol. 25, no. 1, pp. 249–253, Sep. 2020.

[3] X. Hu, K.-K. Wong, K. Yang, and Z. Zheng, "UAV-assisted relaying and edge computing: Scheduling and trajectory optimization," *IEEE Transactions on Wireless Communications*, vol. 18, no. 10, pp. 4738–4752, Jul. 2019.

[4] M. A. ElMossallamy, H. Zhang, L. Song, K. G. Seddik, Z. Han, and G. Y. Li, "Reconfigurable intelligent surfaces for wireless communications: Principles, challenges, and opportunities," *IEEE Transactions on Cognitive Communications and Networking*, vol. 6, no. 3, pp. 990–1002, May. 2020.

[5] C. Huang, A. Zappone, G. C. Alexandropoulos, M. Debbah, and C. Yuen, "Reconfigurable intelligent surfaces for energy efficiency in wireless communication," *IEEE Transactions on Wireless Communications*, vol. 18, no. 8, pp. 4157–4170, Jun. 2019.

[6] Y. Liu, X. Liu, X. Mu, T. Hou, J. Xu, M. Di Renzo, and N. Al-Dhahir, "Reconfigurable intelligent surfaces: Principles and opportunities," *IEEE Communications Surveys & Tutorials*, vol. 23, no. 3, pp. 1546–1577, May. 2021.

[7] L. Li, H. Ruan, C. Liu, Y. Li, Y. Shuang, A. Alù, C.-W. Qiu, and T. J. Cui, "Machine-learning reprogrammable metasurface imager," *Nature communications*, vol. 10, no. 1, pp. 1–8, Mar. 2019.

[8] I. Bucaille, S. Hethuin, A. Munari, R. Hermenier, T. Rasheed, and S. Allsopp, "Rapidly deployable network for tactical applications: Aerial base station with opportunistic links for unattended and temporary events absolute example," in *Proc. IEEE Military Communications Conference (MILCOM)*, San Diego, USA, Nov. 2013.

[9] P. Zhan, K. Yu, and A. L. Swindlehurst, "Wireless relay communications with unmanned aerial vehicles: Performance and optimization," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 47, no. 3, pp. 2068–2085, Jul. 2011.

[10] M. Mozaffari, W. Saad, M. Bennis, and M. Debbah, "Unmanned aerial vehicle with underlaid device-to-device communications: Performance and tradeoffs," *IEEE Transactions on Wireless Communications*, vol. 15, no. 6, pp. 3949–3963, Feb. 2016.

[11] Q. Feng, E. K. Tameh, A. R. Nix, and J. McGeehan, "WLCp2-06: Modelling the likelihood of line-of-sight for air-to-ground radio propagation in urban environments," in *Proc. IEEE Global Communications Conference (Globecom)*, California, USA, Nov. 2006.

[12] M. Mozaffari, W. Saad, M. Bennis, and M. Debbah, "Wireless communication using unmanned aerial vehicles (UAVs): Optimal transport theory for hover time optimization," *IEEE Transactions on Wireless Communications*, vol. 16, no. 12, pp. 8052–8066, Sep. 2017.

[13] M. Mozaffari, A. T. Z. Kasgari, W. Saad, M. Bennis, and M. Debbah, "Beyond 5G with UAVs: Foundations of a 3D wireless cellular network," *IEEE Transactions on Wireless Communications*, vol. 18, no. 1, pp. 357–372, Nov. 2018.

[14] Z. Wang, L. Duan, and R. Zhang, "Adaptive deployment for UAV-aided communication networks," *IEEE transactions on wireless communications*, vol. 18, no. 9, pp. 4531–4543, Jul. 2019.

[15] C. H. Liu, Z. Chen, J. Tang, J. Xu, and C. Piao, "Energy-efficient uav control for effective and fair communication coverage: A deep reinforcement learning approach," *IEEE Journal on Selected Areas in Communications*, vol. 36, no. 9, pp. 2059–2070, Aug. 2018.

[16] G. Lee, M. Jung, A. T. Z. Kasgari, W. Saad, and M. Bennis, "Deep reinforcement learning for energy-efficient networking with reconfigurable intelligent surfaces," in *Proc. IEEE International Conference on Communications (ICC)*, Virtual, Jun. 2020.

[17] H. Zhang, B. Di, L. Song, and Z. Han, "Reconfigurable intelligent surfaces assisted communications with limited phase shifts: How many phase shifts are enough?" *IEEE Transactions on Vehicular Technology*, vol. 69, no. 4, pp. 4498–4502, Feb. 2020.

[18] B. Di, H. Zhang, L. Song, Y. Li, Z. Han, and H. V. Poor, "Hybrid beamforming for reconfigurable intelligent surface based multi-user communications: Achievable rates with limited discrete phase shifts," *IEEE Journal on Selected Areas in Communications*, vol. 38, no. 8, pp. 1809–1822, Jun. 2020.

[19] Y. Chen, B. Ai, H. Zhang, Y. Niu, L. Song, Z. Han, and H. V. Poor, "Reconfigurable intelligent surface assisted device-to-device communications," *IEEE Transactions on Wireless Communications*, vol. 20, no. 5, pp. 2792–2804, Dec. 2020.

[20] Y. Ai, F. A. P. de Figueiredo, L. Kong, M. Cheffena, S. Chatzinotas, and B. Ottersten, "Secure vehicular communications through reconfigurable intelligent surfaces," *IEEE Transactions on Vehicular Technology*, vol. 70, no. 7, pp. 7272–7276, Jun. 2021.

[21] A. Al-Hilo, M. Shokry, M. Elhattab, C. Assi, and S. Sharafeddine, "Reconfigurable intelligent surface enabled vehicular communication: Joint user scheduling and passive beamforming," *IEEE Transactions on Vehicular Technology*, vol. 71, no. 3, pp. 2333–2345, Jan. 2022.

[22] C. Huang, R. Mo, and C. Yuen, "Reconfigurable intelligent surface assisted multiuser MISO systems exploiting deep reinforcement learning," *IEEE Journal on Selected Areas in Communications*, vol. 38, no. 8, pp. 1839–1850, Jun. 2020.

[23] H. Yang, Z. Xiong, J. Zhao, D. Niyato, L. Xiao, and Q. Wu, "Deep reinforcement learning-based intelligent reflecting surface for secure wireless communications," *IEEE Transactions on Wireless Communications*, vol. 20, no. 1, pp. 375–388, Sep. 2020.

[24] X. Cao, B. Yang, C. Huang, C. Yuen, M. Di Renzo, D. Niyato, and Z. Han, "Reconfigurable intelligent surface-assisted aerial-terrestrial communications via multi-task learning," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 10, pp. 3035–3050, Jun. 2021.

[25] S. Li, B. Duo, M. Di Renzo, M. Tao, and X. Yuan, "Robust secure uav communications with the aid of reconfigurable intelligent surfaces," *IEEE Transactions on Wireless Communications*, vol. 20, no. 10, pp. 6402–6417, Apr. 2021.

[26] L. Ge, P. Dong, H. Zhang, J.-B. Wang, and X. You, "Joint beamforming and trajectory optimization for intelligent reflecting surfaces-assisted UAV communications," *IEEE Access*, vol. 8, pp. 78 702–78 712, Apr. 2020.

[27] X. Liu, Y. Liu, and Y. Chen, "Machine learning empowered trajectory and passive beamforming design in UAV-RIS wireless networks," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 7, pp. 2042–2055, Jul. 2020.

[28] B. Shang, R. Shafin, and L. Liu, "UAV swarm-enabled aerial reconfigurable intelligent surface (SARIS)," *IEEE Wireless Communications*, vol. 28, no. 5, pp. 156–163, Oct. 2021.

[29] B. Shang, H. V. Poor, and L. Liu, "Aerial reconfigurable intelligent surfaces meet mobile edge computing," *IEEE Wireless Communications (Early Access)*, pp. 1–9, May. 2022.

[30] Y. Li, C. Yin, T. Do-Duy, A. Masaracchia, and T. Q. Duong, "Aerial reconfigurable intelligent surface-enabled URLLC UAV systems," *IEEE Access*, vol. 9, pp. 140 248–140 257, Oct. 2021.

[31] A. Khalili, E. M. Monfared, S. Zargari, M. R. Javan, N. M. Yamchi, and E. A. Jorswieck, "Resource management for transmit power minimization in uav-assisted ris hetnets supported by dual connectivity," *IEEE Transactions on Wireless Communications*, vol. 21, no. 3, pp. 1806–1822, Aug. 2021.

[32] J. J. Quispe, T. F. Maciel, Y. C. Silva, and A. Klein, "Joint beamforming and bs selection for energy-efficient communications via Aerial-RIS," in *Proc. IEEE Global Communications Conference Workshops (GC Wkshps)*, Madrid, Spain, Dec. 2021.

[33] Z. Yang, M. Chen, W. Saad, W. Xu, M. Shikh-Bahaei, H. V. Poor, and S. Cui, "Energy-efficient wireless communications with distributed reconfigurable intelligent surfaces," *IEEE Transactions on Wireless Communications*, vol. 21, no. 1, pp. 665–679, Jul. 2021.

[34] C. Huang, G. C. Alexandropoulos, A. Zappone, M. Debbah, and C. Yuen, "Energy efficient multi-user MISO communication using low resolution large intelligent surfaces," in *Proc. IEEE Global Communications Conference Workshops (GC Wkshps)*, Abu Dhabi, UAE, Feb. 2018.

[35] S. Jung, W. J. Yun, M. Shin, J. Kim, and J.-H. Kim, "Orchestrated scheduling and multi-agent deep reinforcement learning for cloud-assisted multi-UAV charging systems," *IEEE Transactions on Vehicular Technology*, vol. 70, no. 6, pp. 5362–5377, Feb. 2021.

[36] G. Scutari and Y. Sun, "Parallel and distributed successive convex approximation methods for big-data optimization," in *Multi-agent Optimization*. Springer, Nov. 2018, vol. 2224, pp. 141–308.

[37] S. Salman Hassan, Y. M. Park, Y. Kyaw Tun, W. Saad, Z. Han, and C. S. Hong, "3TO: THz-enabled throughput and trajectory optimization of UAVs in 6G networks by proximal policy optimization deep reinforcement learning," *arXiv e-prints*, pp. arXiv–2202, Feb. 2022.

[38] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, Aug. 2017.

[39] H.-K. Lim, J.-B. Kim, J.-S. Heo, and Y.-H. Han, "Federated reinforcement learning for training control policies on multiple IoT devices," *Sensors*, vol. 20, no. 5, p. 1359, Mar. 2020.

[40] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," in *Proc. International Conference on Machine Learning (ICML)*, New York, USA, Jun. 2016.

[41] S. Mirjalili and A. Lewis, "The whale optimization algorithm," *Advances in engineering software*, vol. 95, pp. 51–67, May. 2016.

[42] M. Mafarja and S. Mirjalili, "Whale optimization approaches for wrapper feature selection," *Applied Soft Computing*, vol. 62, pp. 441–453, Jan. 2018.

[43] Q.-V. Pham, S. Mirjalili, N. Kumar, M. Alazab, and W.-J. Hwang, "Whale optimization algorithm with applications to resource allocation in wireless networks," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 4, pp. 4285–4297, Feb. 2020.