

# Non-iterative Optimization of Trajectory and Radio Resource for Aerial Network

Hyeonsu Lyu, *Student Member, IEEE*, Jonggyu Jang, *Member, IEEE*, Harim Lee, *Member, IEEE*, and Hyun Jong Yang, *Member, IEEE*

**Abstract**—We address a joint trajectory planning, user association, resource allocation, and power control problem to maximize proportional fairness in the aerial IoT network, considering practical end-to-end quality-of-service (QoS) and communication schedules. Though the problem is rather ancient, apart from the fact that the previous approaches have never considered user- and time-specific QoS, we point out a prevalent mistake in coordinate optimization approaches adopted by the majority of the literature. Coordinate optimization approaches, which repetitively optimize radio resources for a fixed trajectory and vice versa, generally converge to local optima when all variables are differentiable. However, these methods often stagnate at a non-stationary point, significantly degrading the network utility in mixed-integer problems such as joint trajectory and radio resource optimization. We detour this problem by converting the formulated problem into the Markov decision process (MDP). Exploiting the beneficial characteristics of the MDP, we design a non-iterative framework that cooperatively optimizes trajectory and radio resources without initial trajectory choice. The proposed framework can incorporate various trajectory planning algorithms such as the genetic algorithm, tree search, and reinforcement learning. Extensive comparisons with diverse baselines verify that the proposed framework significantly outperforms the state-of-the-art method, nearly achieving the global optimum. Our implementation code is available at <https://github.com/hslyu/dbspf>.

**Index Terms**—Trajectory-planning, user association, resource allocation, power control, quality-of-service, Markov decision process.

## I. INTRODUCTION

Aerial networks have been considered a key enabler of the sixth-generation (6G) wireless communication. ITU-R’s suggestion of the key 6G usage scenarios highlights the role of aerial networks in achieving massive communications and ubiquitous connectivity [1]. This entails utilizing aerial vehicles for on-demand services and extensive coverage [2].

These use cases of aerial networks align closely with current research focus on internet-of-things (IoT) networks, as documented in the 3GPP white papers [3], [4]. One of the most critical agendas for realizing aerial IoT networks

Hyeonsu Lyu, Jonggyu Jang and Hyun Jong Yang are with Department of Electrical Engineering, Pohang University of Science and Technology (POSTECH), Pohang 37673, Korea, (e-mail: {hslyu4, jgjang, hyun-yang}@postech.ac.kr). Harim Lee is with the School of Electronic Engineering, Kumoh National Institute of Technology, Gumi, Gyeongbuk 39248, Korea, (e-mail: hrlee@kumoh.ac.kr).

Hyeonsu Lyu and Jonggyu Jang equally contributed to this work. Hyun Jong Yang is the corresponding author.

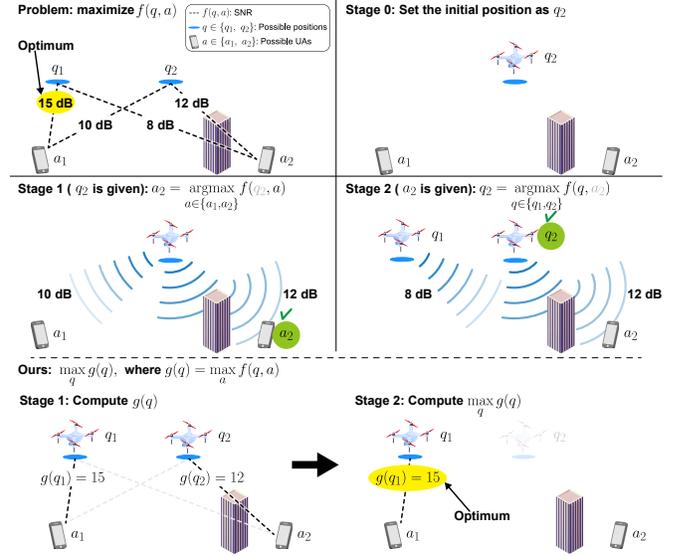


Fig. 1. An illustrative failure scenario in coordinate optimization. Stages 0-2 demonstrate how coordinate optimization approaches converge to the non-stationary point. Once the initial  $q$  is chosen (Stage 0), the following  $a$  is accordingly determined (Stage 1). Then, the next iteration fails to find the optimal  $q$  (Stage 2). Meanwhile, our approach adopts a hierarchical optimization that sweeps the entire variable space.

involves communication scheduling to compensate for the constrained memory and energy capacity [5], [6].

There is a lack of research concerning practical end-to-end requirements, including service scheduling and quality-of-service (QoS) requirements despite the proliferation of research on aerial networks. This is mainly because optimizing aerial networks becomes significantly intractable when such practical requirements are considered. Communication scheduling entangles positional variables and physical-layer variables, such as user association (UA); frequency resource allocation (RA); and power control (PC) variables, along with multiple time steps.

Most research detours the variable entanglement by adopting coordinate optimization approaches which reciprocally optimize TP variables for fixed radio resource variables and vice versa [7]. In the optimization-theoretic viewpoint, these approaches are well-known to converge to local optima in differentiable problems, but may halt at the non-stationary point<sup>1</sup> when optimizing mixed-integer problems, even if the objective is convex [8].

<sup>1</sup>A point is stationary if changes in any direction make utility degradation (e.g. local maxima). A point is non-stationary if the point is not stationary.

We have found an interesting, but catastrophic interaction between trajectory-planning (TP) and radio resource management (RRM), characterized by “*non-optimal initialization can trap the network variables in the non-stationary point.*” Figure 1 illustrates the curse of initialization that could deteriorate the network utility unless the initial variable choice is optimal, which is rarely the case. The discovery suggests that a surprising number of research works may inadvertently optimize the utility of aerial networks in a possibly incorrect manner [9]–[24].

We ask ourselves “*How can we optimize aerial networks while avoiding the curse of initialization?*” and find the answer in the non-iterative approach. As in Fig. 1, we first convert radio resources as a function of trajectory. Then, we find the optimal trajectory by integer programming<sup>2</sup>. The optimal RRM is, in turn, determined by the trajectory, resulting in the optimal solution.

TABLE I  
SUMMARY OF THE RELATED WORKS

Ref.*	3D Traj.*	UA	RA	PC	QoS	Fairness	TCN*
[25]		✓	✓				✓
[26], [27]		✓					✓
[28], [29]	✓						
[12]	✓			✓			
[13]				✓			
[14], [30]		✓			✓		✓
[31]	✓	✓				✓	
[32]		✓					
[33]	✓			✓	✓		
[10], [11]	✓	✓	✓	✓	✓	✓	
Ours	✓	✓	✓	✓	✓	✓	✓

\* Ref.: Reference, Traj.: Trajectory, TCN: Time-critical networks.

While addressing the aforementioned research question, we additionally solve an open problem for the aerial IoT network. As highlighted in Table I, there is a lack of research that studies proportional fairness (PF) maximization with practical end-to-end requirements. Thus, we consider a joint TP; UA, RA, and PC problem for a single UAV base station (UAV-BS) network where IoT devices punctually request a downlink stream above a certain QoS rate. Appendix A provides an in-depth explanation of why the network scenario is considered, referring to adjacent research works.

We hierarchically optimize the PF maximization problem by subordinating the RRM to the TP. Then, we can formulate the proposed problem as a Markov decision process (MDP) where various decision-making schemes can be applied. For the RRM, we employ the Lagrangian method with Karush-Kuhn-Tucker (KKT) conditions, offering a low-complexity algorithm to accelerate the TP optimization.

<sup>2</sup>We assume a discretized trajectory for the non-convex, non-differentiable problem. For a continuously differentiable trajectory, convex relaxations and optimizations can be applied.

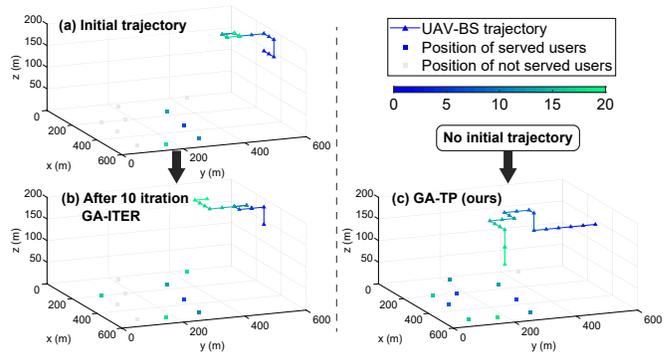


Fig. 2. Motivated illustration of our contribution. Graphs (a) and (b) represent iterative optimization (GA-ITER), and graph (c) represents the proposed method (GA-TP). The two methods share the same environment configurations and algorithms except for the optimization order. While GA-ITER optimizes radio resources for a given initial trajectory, GA-TP adaptively optimizes radio resources according to the changing trajectory. Note that the initial trajectory in (a) restricts the updated trajectory in (b), resulting in a large difference between the result trajectories in (b) and (c). Detail descriptions are provided in the Sec. V.

The salient contributions of the paper can be summarized as follows.

- **Resolve the curse of initialization:** This research is, to the authors’ knowledge, the first to highlight that iterative optimization might be inherently unhelpful in aerial networks. We address the problem through the MDP reformulation, which provides a straightforward, powerful non-iterative framework by separately optimizing TP and RRM problems. Figure 2 briefly compares the result of the iterative optimization and the proposed method. In Sec. V, we show that the proposed method outperforms various comparison schemes and validate the robustness of the proposed method.
- **Propose a novel temporal decoupling method:** We suggest a novel method to separate the PF problem into unit-time sub-problems. The PF entangles all controllable variables in the logarithm over the entire time steps. Then, applying well-known RRM techniques, specialized to optimize unit-time network snapshots, becomes challenging. We devise multiple mathematical tricks to convert the logarithm of sum-rate (PF formulation) into a summation of logarithms. Standing on the temporal decoupling, we show that the proposed RRM scheme tightly achieves the global optimum.
- **Introduce a generalized water-filling algorithm:** We redesign the water-filling algorithm that can be applied to practical scenarios where QoS requirements exist. The conventional water-filling algorithm cannot find the UA and RA solution. The proposed algorithm automatically finds the optimal UA and RA combination when both the minimum resource requirements and resource budget co-exist. Numerical experiments show that the proposed method achieves the global maximum found by the genetic algorithm with the complexity of  $\mathcal{O}(I^2)$  for  $I$  users.

The remainder of the paper is organized as follows. Section II introduces the system model of the IoT networks served

by a UAV-BS and its mathematical representations. Section III-A formulates the PF maximization problem for TP, UA, RA, and PC, and Sec. III-B decomposes the problem into the TP and RRM problems. Section IV presents in-depth explanations for the TP and RRM optimization. Section V provides numerical evaluations of the proposed method. Section VI concludes the paper.

## II. SYSTEM MODEL

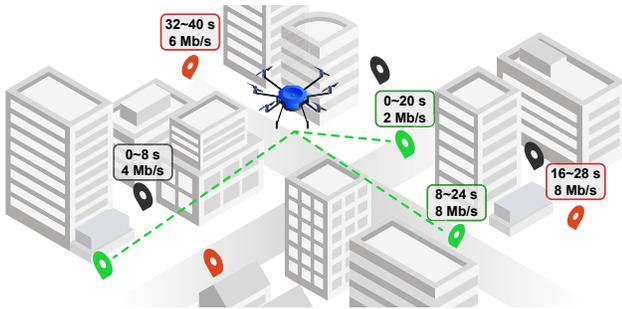


Fig. 3. Illustration of the system model at time  $t = 20$ . The dark gray marks represent users whose request periods have already expired. The green marks with dashed links represent users getting serviced at the current time. The orange marks represent users whose request periods do not expire, but that are not currently being served. Users who have not expired may not be served if their request period does not start. The UAV-BS selects service users by jointly optimizing TP, UA, RA, and PC.

### A. Scenario Description

We consider an IoT network where a single UAV-BS serves  $I$  IoT devices for  $T$  service time slots, as illustrated in Fig. 3. In the network, the UAV-BS selectively provides downlink service to users due to the limitations in the available bandwidth  $B$  and power  $P$ . Each user has a downlink request period and QoS rate to satisfy its service requirements. The users are indexed by the index set  $\mathcal{I} = \{1, \dots, I\}$ .

We assume that the flight map and service timeline are discretized. The time slots are indexed by  $\mathcal{T} = \{1, \dots, T\}$ , each of which has the duration of  $\Delta T$  (s). The UAV-BS covers a square-shaped ground area with the size of  $w \times w$  m<sup>2</sup>, and the flight map is represented by a set  $\mathcal{Q}$  of three-dimensional grids discretized by an interval  $\Delta Q$  (m) as follows:

$$\mathcal{Q} = \{\Delta Q[i, j, k]^T | [i, j, k]^T \in \mathbb{N}^3, h_{\min} \leq k\Delta Q \leq h_{\max}, 0 \leq i\Delta Q \leq w, 0 \leq j\Delta Q \leq w\}, \quad (1)$$

where  $h_{\min}$  and  $h_{\max}$  are the minimum and maximum altitude of the UAV-BS, respectively.

When the  $t$ -th time slot begins, the UAV-BS must be located on a certain grid point of the flight map. We denote  $\mathbf{q}^{(t)}$  as the position of the UAV-BS at the beginning of the  $t$ -th time slot, where  $\mathbf{q}^{(t)}$  is defined as

$$\mathbf{q}^{(t)} = [x^{(t)}, y^{(t)}, h^{(t)}]^T \in \mathcal{Q}. \quad (2)$$

The initial position of the UAV-BS is denoted as  $\mathbf{q}^{(0)}$ . The UAV-BS cannot exceed its maximum velocity  $v$ . In other words, the position vector  $\mathbf{q}^{(t)}$  is constrained by

$$\|\mathbf{q}^{(t)} - \mathbf{q}^{(t-1)}\|_2 \leq v\Delta T, \quad \forall t \in \mathcal{T}. \quad (3)$$

Combining (2) and (3), all possible positions  $\mathbf{q}^{(t)}$  at time slot  $t$  are determined by the set  $S(\mathbf{q}^{(t-1)}) = \{\mathbf{q} \in \mathcal{Q} | \|\mathbf{q} - \mathbf{q}^{(t-1)}\|_2 \leq v\Delta T\}$  as follows:

$$\mathbf{q}^{(t)} \in S(\mathbf{q}^{(t-1)}). \quad (4)$$

Users are stationary during the total service time. The  $i$ -th user is located at  $\mathbf{q}_i = [x_i, y_i, 0]^T$ . To prolong the device lifetime, each user requests a downlink service only if the current flight time slot  $t$  is within the short-term period<sup>3</sup>  $[s_i, s_i + T_i)$ . We define a binary indicator  $d_i^{(t)}$  to specify user indices requesting services at time slot  $t$  as follows:

$$d_i^{(t)} = \begin{cases} 1, & \text{if } s_i \leq t < s_i + T_i \\ 0, & \text{otherwise.} \end{cases} \quad (5)$$

For each time slot  $t$ , the UAV-BS determines a set of users to provide downlinks. Accordingly, we define the UA variables as

$$\alpha_i^{(t)} = \begin{cases} 1, & \text{if the UAV-BS serves user } i \text{ at time slot } t \\ 0, & \text{otherwise.} \end{cases} \quad (6)$$

We denote  $\beta_i^{(t)}$  as the amount of frequency resources that the UAV-BS allocates to the  $i$ -th user at time slot  $t$ . The summation of the allocated bandwidth proportions must not exceed the available bandwidth  $B$ , so we have the following constraint:

$$\sum_{i \in \mathbf{A}^{(t)}} \beta_i^{(t)} \leq B, \quad \forall t \in \mathcal{T}, \quad (7)$$

where  $\mathbf{A}^{(t)} = \{i \in \mathcal{I} | \alpha_i^{(t)} = 1\}$  denotes the set of associated users at time slot  $t$ .

The UAV-BS regulates its transmission power by adjusting the power spectral density (PSD). The variable  $\rho_i^{(t)}$  is defined as the assigned PSD to the  $i$ -th user at time slot  $t$ . Denoting the maximum transmission power at each time slot as  $P$ , the transmission power is constrained as follows:

$$\sum_{i \in \mathbf{A}^{(t)}} \rho_i^{(t)} \beta_i^{(t)} \leq P, \quad \forall t \in \mathcal{T}. \quad (8)$$

### B. Pathloss and Data Rate Model

The propagation channel of the  $i$ -th user follows a probabilistic LoS channel [34]. The probability that the  $i$ -th user has an LoS wireless link is defined as

$$\mathbb{P}_i^{(t)} = \left(1 + a \cdot \exp(-b(\theta_i^{(t)} - a))\right)^{-1}, \quad (9)$$

<sup>3</sup>If  $[s_i, T_i) = [0, T)$  for  $i \in \mathcal{I}$ , the proposed problem has a similar formulation with the problem in [11] except for the objective function. However, the difference in the objective function causes a significant performance gap, which is shown in Sec. V.

where  $\theta_i^{(t)} = \arcsin\left(\frac{h^{(t)}}{\|\mathbf{q}^{(t)} - \mathbf{q}_i\|_2}\right)$  is the elevation angle between the UAV-BS and the  $i$ -th user at time slot  $t$ . Variables  $a$  and  $b$  in (9) are environmental parameters that are determined according to the environmental characteristics [34], [35]. For example, possible pairs of  $(a, b)$  could be (4.76, 0.37) for a rural environment; and (9.64, 0.06) for a dense urban environment. We remark that the  $i$ -th user has an NLoS link at time slot  $t$  with a probability of  $1 - P_i^{(t)}$ .

The average pathloss for the  $i$ -th user at time slot  $t$  is defined as

$$\xi_i^{(t)} = \text{PL}(\mathbf{q}^{(t)}, \mathbf{q}_i) + \mathbb{P}_i^{(t)} \eta_{\text{LoS}} + (1 - \mathbb{P}_i^{(t)}) \eta_{\text{NLoS}}, \quad (10)$$

where pathloss  $\text{PL}(\mathbf{q}^{(t)}, \mathbf{q}_i) = 20 \log(4\pi f \|\mathbf{q}^{(t)} - \mathbf{q}_i\|_2 / c)$ ;  $f$  is the carrier frequency;  $c$  is the speed of light; constant  $\eta_{\text{LoS}}$  and  $\eta_{\text{NLoS}}$  are the expected value of excessive pathloss for the LoS and NLoS links [34], [36]. The first term of (10) is the free-space pathloss and the sum of the remaining terms represents expected excessive pathloss [34].

The data rate of the  $i$ -th user at time slot  $t$  is defined as

$$R_i^{(t)} = d_i^{(t)} \beta_i^{(t)} \log_2 \left( 1 + \frac{\rho_i^{(t)} 10^{-\xi_i^{(t)}/10}}{N_0} \right), \quad (11)$$

where constant  $N_0$  denotes additive white Gaussian noise PSD.

To ensure the served users' QoS requirements, the UAV-BS should serve user  $i$  with a link that has a data rate higher than  $r_i$  if the  $i$ -th user receives downlink data at time slot  $t$ . Therefore, we have

$$R_i^{(t)} \geq \alpha_i^{(t)} r_i. \quad (12)$$

We define  $R_i$  as the sum-rate of the  $i$ -th user across the UAV-BS service timeline, denoted as

$$R_i = \sum_{t \in \mathcal{T}} R_i^{(t)}. \quad (13)$$

Then, the PF of the served users is written by  $\sum_{i \in \mathbf{A}} \log R_i$ . We aim to maximize the PF that specializes in prioritized service scheduling. PF accounts for both the total network throughput and the number of served users as the logarithm operation prevents the network throughput from being concentrated on a few number of users. However, sum-rate maximization prioritizes a user with the greatest spectral efficiency, not considering the fairness of service. Then the number of served users inevitably decreases when we target to maximize the sum-rate.

### III. PROBLEM FORMULATION

#### A. Original Full-Time Problem

By gathering the previously defined constraints, the joint problem of TP, UA, RA, and PC is formulated to maximize

the PF as

$$\mathcal{P}\mathbf{1} : \max_{\mathbf{Q}, \mathbf{A}, \mathbf{B}, \mathbf{P}} \sum_{i \in \mathbf{A}} \log R_i \quad (14a)$$

$$\text{s.t. } \mathbf{q}^{(t)} \in S(\mathbf{q}^{(t-1)}), \quad (14b)$$

$$\beta_i^{(t)} \geq 0, \forall i \in \mathcal{I}, \quad (14c)$$

$$\sum_{i \in \mathbf{A}^{(t)}} \beta_i^{(t)} \leq B, \quad (14d)$$

$$\rho_i^{(t)} \geq 0, \forall i \in \mathcal{I}, \quad (14e)$$

$$\sum_{i \in \mathbf{A}^{(t)}} \rho_i^{(t)} \beta_i^{(t)} \leq P, \quad (14f)$$

$$R_i^{(t)} \geq \alpha_i^{(t)} r_i, \forall i \in \mathcal{I}, \quad (14g)$$

$$\forall t \in \mathcal{T} \text{ for (14b)-(14g)}.$$

For brevity of the notations, we define augmented vectors and matrices<sup>4</sup> of variables as follows:  $\mathbf{Q} = [\mathbf{q}^{(t)}]_{t \in \mathcal{T}}$ ,  $\mathbf{B} = [\beta_i^{(t)}]_{i \in \mathcal{I}, t \in \mathcal{T}}$  and  $\mathbf{P} = [\rho_i^{(t)}]_{i \in \mathcal{I}, t \in \mathcal{T}}$ . Also,  $\mathbf{A} = \bigcup_{t \in \mathcal{T}} \mathbf{A}^{(t)}$  represents a set of users that have been served at least once during the UAV-BS service timeline.

Constraints (14b) to (14g) are described as follows: (14b) correspond to (4); (14c) and (14e) represent the non-negativity of the allocated bandwidth and PSD, respectively; (14d), (14f), and (14g) are equivalent with (7), (8), and (12), respectively.

Problem  $\mathcal{P}\mathbf{1}$  is a non-convex mixed-integer problem that is generally known to be NP-hard [37], because the association variables  $\alpha_i^{(t)}, \forall i, t$  are binary integers, and since the objective function is non-convex with respect to the position vector  $\mathbf{q}^{(t)}$ . Then, solving the problem within a feasible complexity becomes challenging if numerous time slots are involved in Problem  $\mathcal{P}\mathbf{1}$  because all the variables are coupled over the entire service timeline  $\mathcal{T}$ . Therefore, we decompose Problem  $\mathcal{P}\mathbf{1}$  into  $T$  sub-problems. Each sub-problem has a Markov (or memoryless) property [38], making themselves independent of each other.

#### B. On the Separation of Control and RRM

We first decompose Problem  $\mathcal{P}\mathbf{1}$  into sub-problems where each sub-problem corresponds to a single time slot. All controllable variables across time slots  $t \in \mathcal{T}$  are closely coupled in the objective function of Problem  $\mathcal{P}\mathbf{1}$ . However, we can separate the variables and constraints related to the position from those related to the radio resources once Problem  $\mathcal{P}\mathbf{1}$  is temporally decomposed.

The key idea of the separation is converting summation into multiplication. For example,  $\sum_{i=1}^n i$  is equivalent with

$$\sum_{i=1}^n i = 1 \cdot \frac{1+2}{1} \cdot \frac{1+2+3}{1+2} \cdots \frac{1+\cdots+n}{1+\cdots+n-1}. \quad (15)$$

Using the trick, we convert  $\log R_i = \log\left(\sum_{t \in \mathcal{T}} R_i^{(t)}\right)$  into  $\log R_i = \sum_{t \in \mathcal{T}} \log\left(1 + R_i^{(t)} / \sum_{k=1}^{t-1} R_i^{(k)}\right)$ . Then, we can

<sup>4</sup>We define a vector-building operator  $[e^{(t)}]_{t \in \mathcal{T}}$  to represent  $[e^{(1)}, e^{(2)}, \dots, e^{(T)}]^T$ . Then we define a matrix-building operator  $[e_i^{(t)}]_{i \in \mathcal{I}, t \in \mathcal{T}}$  as  $[[e_1^{(t)}]_{t \in \mathcal{T}}, [e_2^{(t)}]_{t \in \mathcal{T}}, \dots, [e_I^{(t)}]_{t \in \mathcal{T}}]$ .

define a lower-bound of Problem  $\mathcal{P}1$  which consists of  $T$  sub-problems.

**Proposition 1** (Lower bound of Problem  $\mathcal{P}1$ ). *The following inequality holds:*

$$\max_{\mathbf{Q}, \mathbf{A}, \mathbf{B}, \mathbf{P}} \sum_{i \in \mathcal{A}} \log R_i \quad (16)$$

$$= \max_{\mathbf{Q}} \max_{\mathbf{A}, \mathbf{B}, \mathbf{P}} \sum_{t \in \mathcal{T}} \sum_{i \in \mathcal{A}} \log \left( 1 + \frac{R_i^{(t)}}{\sum_{k=0}^{t-1} R_i^{(k)}} \right) \quad (17)$$

$$\geq \max_{\mathbf{Q}} \sum_{t \in \mathcal{T}} \max_{\mathbf{A}^{(t)}, \mathbf{B}^{(t)}, \mathbf{P}^{(t)}} \sum_{i \in \mathcal{A}^{(t)}} \log \left( 1 + \frac{R_i^{(t)}}{\sum_{k=0}^{t-1} R_i^{(k)}} \right), \quad (18)$$

where  $R_i^{(0)} = 1$  for all  $i \in \mathcal{I}$ ,  $\mathbf{A}^{(t)} = \{i \in \mathcal{I} | \alpha_i^{(t)} = 1\}$ ,  $\mathbf{B}^{(t)} = [\beta_i^{(t)}]_{i \in \mathcal{I}}$ , and  $\mathbf{P}^{(t)} = [\rho_i^{(t)}]_{i \in \mathcal{I}}$ .

*Proof:* The proof is shown in Appendix B.  $\blacksquare$

From Proposition 1, we relax Problem  $\mathcal{P}1$  to maximize the lower bound of the problem; that is, Eq. (18) can be re-written as

$$\mathcal{P}2 : \max_{\mathbf{q}^{(t)} \in \mathcal{S}(\mathbf{q}^{(t-1)})} \sum_{t=1}^T f(\mathbf{q}^{(t)}), \quad (19)$$

where

$$f(\mathbf{q}^{(t)}) = \max_{\mathbf{A}^{(t)}, \mathbf{B}^{(t)}, \mathbf{P}^{(t)}} \sum_{i \in \mathcal{A}^{(t)}} \log \left( 1 + \frac{R_i^{(t)}}{\sum_{k=0}^{t-1} R_i^{(k)}} \right) \quad (20a)$$

$$\text{s.t.} \quad (14c) - (14g). \quad (20b)$$

Here, Problem  $\mathcal{P}2$  only involves with the positional variable  $\mathbf{q}^{(t)}$  and Eq. (20) associates with the radio resource variables  $\mathbf{A}^{(t)}$ ,  $\mathbf{B}^{(t)}$ , and  $\mathbf{P}^{(t)}$ . We remark that computing  $f(\mathbf{q}^{(t)})$  requires  $R_i^{(k)}$  for  $k = 1, \dots, t-1$ . Thus,  $f(\mathbf{q}^{(t)})$  need to be sequentially computed  $f(\mathbf{q}^{(t)})$  from  $t = 1$  to  $t = T$ .

The computing process can be considered a Markov decision process (MDP) in that sequential actions  $\mathbf{q}^{(1)}, \dots, \mathbf{q}^{(T)}$  determine a reward as  $\sum_{t \in \mathcal{T}} f(\mathbf{q}^{(t)})$ . Once the optimal RRM scheme  $f(\mathbf{q}^{(t)})$  is given, the control problem  $\mathcal{P}2$  can be considered as finding the best solution of the MDP, as depicted in Fig. 4.

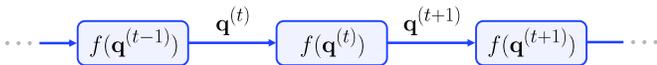


Fig. 4. Illustration of the objective in problem  $\mathcal{P}2$  in the viewpoint of MDP. For each time slot  $t$ , state is given as all variables before  $t$ ; action is  $\mathbf{q}^{(t)}$ ; and reward is  $f(\mathbf{q}^{(t)})$ .

#### IV. OPTIMIZATION OF RRM AND CONTROL

The MDP reward  $f(\mathbf{q}^{(t)})$  is ideally designed if we can find the optimal UA, RA, and PC. Through the optimization-theoretic approach, we develop a fast and accurate RRM that makes various decision-making schemes affordable in terms of the computing resource budget. Then, the global optimum of Problem  $\mathcal{P}2$  can be obtained by finding the

optimal trajectory  $\mathbf{q}^{(1)}, \dots, \mathbf{q}^{(T)}$ . We aim to find the solution of Problem  $\mathcal{P}2$  by applying various decision-making schemes to find the optimal trajectory.

##### A. Finding the Optimal Radio Resource Management

We first optimize the integer variable  $\mathbf{A}^{(t)}$  by suggesting a generalized water-filling algorithm that can jointly optimize  $\mathbf{A}^{(t)}$  and  $\mathbf{B}^{(t)}$ . Then, the RA and PC variables,  $\mathbf{B}^{(t)}$  and  $\mathbf{P}^{(t)}$ , are iteratively optimized by using the Lagrangian method and the KKT conditions.

1) *Initial User Association:* To determine optimal  $\mathbf{A}^{(t)}$ , we first assume that the PSD variables  $\rho_i^{(t)}$  are equally assigned as  $\rho_i^{(t)} = P/B$  for all  $i \in \mathcal{I}$ . Then the problem (20) in  $f(\mathbf{q})$  can be simplified as

$$\max_{\mathbf{A}^{(t)}, \mathbf{B}^{(t)}} \sum_{i \in \mathcal{I}} \log \left( 1 + \frac{R_i^{(t)}}{\sum_{k=0}^{t-1} R_i^{(k)}} \right) \quad (21a)$$

$$\text{s.t.} \quad \beta_i^{(t)} \geq \alpha_i^{(t)} r_i / e_i^{(t)}, \forall i, \quad (21b)$$

$$\sum_{i \in \mathcal{A}^{(t)}} \beta_i^{(t)} \leq B, \quad (21c)$$

where  $e_i^{(t)} = \log_2 \left( 1 + \frac{\rho_i^{(t)} 10^{-\xi_i^{(t)}/10}}{N_0} \right)$  is the spectral efficiency of the  $i$ -th user at time slot  $t$ . This can be considered as a water-filling problem with non-zero resource requirements  $\alpha_i^{(t)} r_i / e_i^{(t)}$ . Therefore, the optimal  $\mathbf{B}^{(t)}$  is also provided as a water-filling solution.

We can obtain the optimal  $\mathbf{B}^{(t)}$  by using the KKT conditions. For a given  $\mathbf{A}^{(t)}$  that does not violate constraint (21b) and (21c), the optimal  $\beta_i^{(t)}$  is

$$\beta_i^{(t)} = \begin{cases} \max \left( \frac{r_i}{e_i^{(t)}}, \frac{1}{\lambda^{(t)}} - \frac{1}{w_i^{(t)}} \right) & , \text{ if } \alpha_i^{(t)} d_i^{(t)} = 1, \\ 0 & , \text{ otherwise,} \end{cases} \quad (22)$$

where  $\lambda^{(t)}$  can be uniquely determined to satisfy the constraint (21c) with equality condition<sup>5</sup> and  $1/w_i^{(t)} = \sum_{k=0}^{t-1} R_i^{(k)} / e_i^{(t)}$ . Physically, the solution implies that the UAV-BS allocates a little frequency resource to the user who has a high total received data and low spectral efficiency. The derivation of (22) is provided in Appendix C. Figure 5 depicts our intuition on the lower-bounded water-filling solution.

Now, the objective (21a) can be maximized by finding the optimal  $\mathbf{A}^{(t)}$  as the optimal  $\mathbf{B}^{(t)}$  can be obtained for arbitrary  $\mathbf{A}^{(t)}$ . We adopt an incremental algorithm that finds a local-optimal solution (Alg. 1). The key idea of the algorithm is to sequentially find a better UA combination than the current UA combination  $\mathbf{A}_{\text{current}}^{(t)}$ .

Let  $\mathbf{A}_{\text{current}}^{(t)}$  be a UA set and  $\mathbf{A}_{\text{current}+i}^{(t)}$  contains one more user  $i \notin \mathbf{A}_{\text{current}}^{(t)}$ , denoted as

$$\mathbf{A}_{\text{current}+i}^{(t)} = \mathbf{A}_{\text{current}}^{(t)} \cup \{i\}. \quad (23)$$

If  $\mathbf{A}_{\text{current}+i}^{(t)}$  satisfies constraint (21b) and (21c),  $\mathbf{A}_{\text{current}+i}^{(t)}$  could be a maximizer of Problem (21). Two constraints (21b)

<sup>5</sup>We use the bisection method to find the solution for  $\lambda$ .

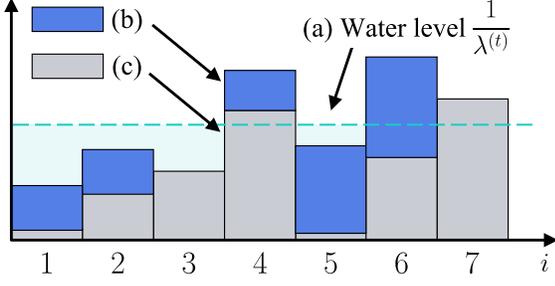


Fig. 5. Visualization of the water-filling solution  $\beta_i^{(t)}$  in (22) for 7 users. (a) The dashed line implies the water level. (b) The blue column indicates the minimum RA requirement  $r_i/e_i^{(t)}$ . (c) The grey column represents the given ground heights  $1/w_i^{(t)} = \sum_{k=0}^{t-1} R_i^{(k)}/e_i^{(t)}$ , which implies that the UAV-BS allocates a low bandwidth proportion to a user if the user has a high total received data and low spectral efficiency.

---

#### Algorithm 1: User Association Algorithm

---

```

1: Input  $\mathbf{q}^{(t)}, \mathbf{P}^{(t)}$ 
2: Output  $\mathbf{A}^{(t)}, \mathbf{B}^{(t)}$ 
3:  $\mathbf{A}_{\text{current}}^{(t)} \leftarrow \{\}$ .
4: while  $\mathbf{A}_{\text{current}}^{(t)} \neq \mathbf{A}_{\text{next}}^{(t)}$  do
5:    $\mathcal{I}_{\text{feasible}} \leftarrow \{i \in \mathcal{I} \mid \sum_{i \in \mathbf{A}_{\text{current}}^{(t)}} r_i/e_i^{(t)} \leq B\}$ 
6:    $\mathbf{A}_{\text{next}}^{(t)} \leftarrow$  The best UA set among  $\mathbf{A}_{\text{current}}^{(t)}$  and
      $\mathbf{A}_{\text{current}+i}^{(t)}$ ,  $i \in \mathcal{I}_{\text{feasible}}$ .
7: end
8:  $\mathbf{A}^{(t)} \leftarrow \mathbf{A}_{\text{next}}^{(t)}$ 
9: Compute  $\mathbf{B}^{(t)}$  for  $\mathbf{A}^{(t)}$  by using (22).

```

---

and (21c) can be simplified, so we can define a set of feasible additive users as

$$\mathcal{I}_{\text{feasible}} = \{i \in \mathcal{I} \mid \sum_{i \in \mathbf{A}_{\text{current}}^{(t)}} r_i/e_i^{(t)} \leq B\}. \quad (24)$$

Then, we can define a new set  $\mathbf{A}_{\text{next}}^{(t)}$  that makes the objective (21a) be the greatest among the set  $\mathbf{A}_{\text{current}}^{(t)}$  and  $\mathbf{A}_{\text{current}+i}^{(t)}$  for  $i \in \mathcal{I}_{\text{feasible}}$ . If  $\mathbf{A}_{\text{next}}^{(t)}$  is equal to  $\mathbf{A}_{\text{current}}^{(t)}$ , set  $\mathbf{A}_{\text{next}}^{(t)}$  is a solution. Otherwise,  $\mathbf{A}_{\text{current}}^{(t)}$  is updated to  $\mathbf{A}_{\text{next}}^{(t)}$  and the above procedure is repeated until finding the solution.

By iteratively updating the UA set from the empty set  $\mathbf{A}_{\text{current}}^{(t)} = \{\}$ , the proposed algorithm can find the local-optimal UA and RA variables.

2) *Resource Allocation and Power Control*: The optimal  $\mathbf{B}^{(t)}$  and  $\mathbf{P}^{(t)}$  for the given  $\mathbf{A}^{(t)}$  is optimized by iteratively optimizing the following two problems.

Once the initial  $\mathbf{A}^{(t)}$  and  $\mathbf{B}^{(t)}$  is determined by Alg. 1,  $\mathbf{P}^{(t)}$  can be accordingly determined by solving

$$\max_{\mathbf{P}^{(t)}} \sum_{i \in \mathbf{A}^{(t)}} \log \left( 1 + \frac{R_i^{(t)}}{\sum_{k=0}^{t-1} R_i^{(k)}} \right) \quad (25a)$$

$$\text{s.t. (14f), (14e), and (14g).} \quad (25b)$$

---

#### Algorithm 2: Radio Resource Management $f(\cdot)$

---

```

1: Input  $\mathbf{q}^{(t)}$ 
2: Output  $f(\mathbf{q}^{(t)}), \mathbf{A}^{(t)}, \mathbf{B}^{(t)}, \mathbf{P}^{(t)}$ 
3: Initialize  $f \leftarrow 0, f_{\text{prev}} \leftarrow \infty, \rho_i^{(t)} = P/B, \forall i \in \mathcal{I}$ 
4: Function  $f(\mathbf{q}^{(t)})$ :
5:   Update  $\mathbf{A}^{(t)}$  and  $\mathbf{B}^{(t)}$  by using Alg. 1.
6:   while  $\|f - f_{\text{prev}}\| > \epsilon$  do
7:     Update  $\mathbf{B}^{(t)}$  by using Alg. 3.
8:     Update  $\mathbf{P}^{(t)}$  by using (61).
9:      $f_{\text{prev}} \leftarrow f, f \leftarrow \sum_{i \in \mathbf{A}^{(t)}} \log \left( 1 + \frac{R_i^{(t)}}{\sum_{k=0}^{t-1} R_i^{(k)}} \right)$ 
10:  end
11: return  $f$ 

```

---

Similarly, fixing  $\mathbf{A}^{(t)}$  and  $\mathbf{P}^{(t)}$  gives

$$\max_{\mathbf{B}^{(t)}} \sum_{i \in \mathbf{A}^{(t)}} \log \left( 1 + \frac{R_i^{(t)}}{\sum_{k=0}^{t-1} R_i^{(k)}} \right) \quad (26a)$$

$$\text{s.t. (14c), (14d), (14g), and (14f).} \quad (26b)$$

Using the convexity of these two problems,

The Lagrangian dual method and KKT conditions are applied to minimize the computing time, and we find an  $\mathcal{O}(I^2)$  algorithm (Alg. 3 in Appendix D) and closed-form solution (Eq. (61) in Appendix E) for Problem (25) and (26), respectively. Appendices D and E describe in-depth optimization processes.

3) *Overall RRM Algorithm*: Combining the UA, RA, and PC schemes altogether, the RRM scheme to compute  $f(\mathbf{q}^{(t)})$  can be summarized as Alg. 2. We first initialize the PC variables as  $\rho_i^{(t)} = P/B$  for all  $i \in \mathcal{I}$ . Then, in Line 5, we obtain initial UA and RA variables for given  $\mathbf{P}^{(t)}$  by Alg. 1. After the UA variables are initialized, the RA and PC variables are iteratively optimized in Lines 6 to 9, until the objective function converges. Figure 6a visualizes these processes as algorithm flowcharts.

4) *Computational Complexity*: The complexity of Alg. 2 can be computed by adding the complexity of Alg. 1 and the complexity of the iteration in Alg. 2.

The complexity of Alg. 1 is  $\mathcal{O}(I^2)$  as the maximum number of comparisons (Line 6 in Alg. 1) is  $\frac{I(I+1)}{2}$ . The RA scheme in Alg. 3 utilizes gradient descent with the complexity of  $\mathcal{O}(\frac{I}{\epsilon})$  and the closed-form solution (61) takes  $\mathcal{O}(1)$ . Then, the complexity of the Alg. 2 is  $\mathcal{O}(\frac{I^2}{\epsilon^2})$  since the iterations from Line 6 to Line 9 in Alg. 2 takes  $\mathcal{O}(\frac{1}{\epsilon})$ .

#### B. Decision-Making Algorithms for Trajectory-Planning

A consecutive update of  $\mathbf{q}^{(1)}, \dots, \mathbf{q}^{(T)}$  provides a solution of the reformulated lower-bound problem  $\mathcal{P2}$ , as function  $f(\mathbf{q}^{(t)})$  now can be obtained by Alg. 2. The lower-bounded problem is considered a Markov decision process as illustrated in Fig. 4, where various optimization strategies exist.

We adopt a genetic algorithm (GA), limited depth-first search (DFS), and deep q-learning (DQN) as the main

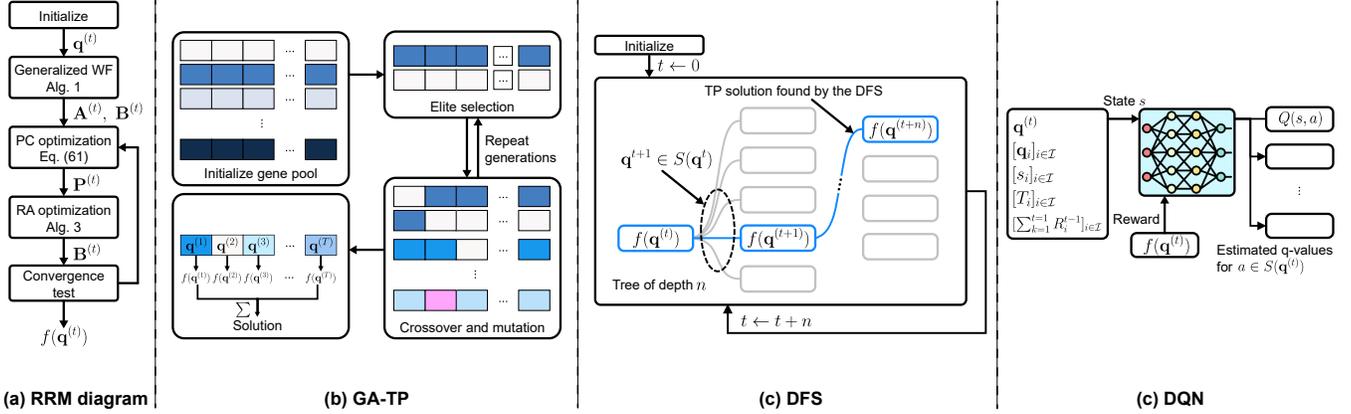


Fig. 6. Visualization of the proposed algorithms (a) Algorithm flowchart of the RRM scheme in Alg. 2. (b) Visualization of GA-TP. Each gene corresponds to a trajectory. Fitness for each gene is computed at the elite selection. (c) Visualization of DFS with a tree representation of length  $n$  sub-trajectory. (d) State, reward, and action for DQN.

strategies. A detailed algorithm design, characteristics, and summary of pros and cons can be arranged as follows:

- **GA-TP [39] (Naive upper-bound)**: A canonical GA with elitism can find the global optimum after a sufficient number of generations and populations [39]. The gene is defined as a trajectory  $\mathbf{Q} = [\mathbf{q}^{(t)}]_{t \in \mathcal{T}}$ ; and the fitness is defined as  $\sum_{t \in \mathcal{T}} f(\mathbf{q}^{(t)})$ . The genetic algorithm is configured as 10,000 generations, 10 elite counts, 50 populations, and 10% mutation probability. Figure 6b demonstrates the overall algorithm flows. In the elite selection, the fitness is computed for each gene.
- **DFS [40]**: DFS targets to find the best sub-trajectory  $q^{(t)}, \dots, q^{(t+n)}$  by the depth-first search and repeats the computation until DFS covers the entire time slots. As illustrated in Fig. 6c, DFS considers all possible actions as a tree and finds the best sub-trajectory. Then, DFS can find the global optimal solution when  $t = 0$  and  $n = T$ . However, the practical depth  $n$  is limited due to the limitations in the memory size and computing time [40]. We set the search depth  $n$  as  $\{1, 3, 5\}$ .
- **DQN [41]**: DQN is a well-known algorithm that can optimize a discrete MDP. To determine  $\mathbf{q}^{(t)}$ , we configured the state as a concatenation of  $\mathbf{q}^{(t-1)}$ ,  $[\mathbf{q}_i]_{i \in \mathcal{I}}$ ,  $[s_i]_{i \in \mathcal{I}}$ ,  $[T_i]_{i \in \mathcal{I}}$ ,  $[\sum_{k=1}^{t-1} R_i^{(k)}]_{i \in \mathcal{I}}$ ; action space as  $S(\mathbf{q}^{(t-1)})$ ; and reward  $r(\cdot)$  as  $r(\mathbf{q}^{(t)}) = f(\mathbf{q}^{(t)})$ . Fig. 6c visualizes states, rewards, and actions of DQN. Implementation details are provided in the Appendix F.

The complexity of the three TP schemes combined with the RRM scheme is computed by multiplying the complexity of the TP scheme and Alg. 2. When  $m$  generations and  $n$  populations are given, the complexity of GA is  $\mathcal{O}(mn \frac{T^2}{\epsilon^2})$ . For DFS, if the maximum number of reachable grid points  $m$  with search depth  $n$  is defined as  $m = \max_{\mathbf{q} \in \mathcal{Q}} |S(\mathbf{q})|$ , the complexity of DFS is  $\mathcal{O}(mn \frac{T^2}{\epsilon^2})$ . The complexity of DQN is  $\mathcal{O}(\frac{T^2}{\epsilon^2})$  as the feed-forward of the network takes  $\mathcal{O}(1)$ .

## V. SIMULATION RESULTS

In various performance metrics, we compare the proposed schemes with the other PF-maximizing schemes. We first

TABLE II  
PARAMETER CONFIGURATIONS

Parameter	Value
Map width $w$ (m)	600
Minimum altitude $h_{\min}$ (m)	50
Maximum altitude $h_{\max}$ (m)	200
Unit length of the grid map $\Delta Q$ (m)	40
Number of users $I$	10 to 80
Number of time slots $T$	20
Unit length of the time slot $\Delta T$ (s)	3
Vehicle velocity $v$ (m/s)	15
Carrier frequency (GHz)	2
Bandwidth $B$ (MHz)	$\{2, 5, 10\}$
Transmission power $P$ (dBm)	23
Noise spectral density (dBm/Hz)	-173.8
Data rate constraint $r_i, \forall i \in \mathcal{I}$ (Mbps)	0 to 10
Environmental parameter $a$	9.64
Environmental parameter $b$	0.06
Excessive LoS pathloss $\eta_{\text{LoS}}$ (dB)	1
Excessive NLoS pathloss $\eta_{\text{NLoS}}$ (dB)	40
Rician $K$ -factor	12

show that the proposed RRM converges to the global optimum of Problem (20), then evaluate three control schemes (GA, DFS, and DQN). Simulation parameters are chosen from existing works [36], [42], and 3GPP specifications [43], which are listed in Table II.

In the experiments, the length of the request period  $T_i, \forall i \in \mathcal{I}$  and the initial request time  $s_i, \forall i \in \mathcal{I}$ , follow the discrete uniform distributions,  $\mathcal{U}[4, 8]$  and  $\mathcal{U}[0, T]$ , respectively.

We have developed a TP simulator that implements the proposed method and several comparison schemes. The simulations are implemented under Python 3.9 on AMD Ryzen™ 9 5950X processor.

To evaluate the performance of the RRM scheme Alg. 2, we additionally implement a genetic RRM (GA-RRM) and maximum SINR scheme (MAX SINR) that provide a solution  $\mathbf{A}^{(t)}, \mathbf{B}^{(t)}$ , and  $\mathbf{P}^{(t)}$  in Problem (20).

- **GA-RRM [39] (Naive upper-bound)**: The gene is defined as a concatenation of  $\mathbf{A}^{(t)}, \mathbf{B}^{(t)}$ , and  $\mathbf{P}^{(t)}$ ; and the fitness is  $f(\mathbf{q}^{(t)})$ . The GA-RRM is configured as

10,000 generations, 10 elite counts, 100 populations, and 10% mutation probability.

- **MAX SINR [44], [45]:** For the UA scheme, a user with the greatest SINR is associated with the UAV-BS to maximize the throughput. This UA scheme is prevalent in the sum-rate maximization problems [44], [45], providing a benchmark of the heuristic UA scheme. RA and PC schemes are the same as those of Alg. 2.

Then, we evaluate the proposed TP schemes in Sec. IV-B with the following comparison schemes:

- **GA-ITER [7], [46] (Conventional iterative optimization)** Iterative optimization of the TP and RRM. The canonical GA and Alg. 2 are used for TP and RRM, respectively. The GA parameters are identically configured as in GA-TP. Properties are listed as follows:
  - *Stationary RRM:* The RRM variables in GA-ITER are fixed when optimizing TP by GA. However, the RRM in GA-TP changes for every gene.
  - *Iterations:* Starting with a random trajectory, GA-ITER optimizes radio resources  $\mathbf{A}$ ,  $\mathbf{B}$ , and  $\mathbf{P}$  for a fixed trajectory. Then, the trajectory is optimized with the fixed  $\mathbf{A}$ ,  $\mathbf{B}$ , and  $\mathbf{P}$ . Iteration stops either when there is no improvement in the objective (19) or when GA-ITER reaches 10 iterations<sup>6</sup>.
- **OFDMA [10], [11]:** We implement the closest research<sup>7</sup> that optimizes TP, UA, RA, and PC variables. The key differences are listed as follows:
  - *The objective function:* The objective function is relaxed from the summation of a logarithm to the weighted sum-rate. For a fair comparison, the PF of all schemes are computed as a summation of logarithms (14a).
  - *QoS constraint:* The proposed scheme adopts the data rate as a QoS constraint, [10] and [11] adopt the SNR threshold. For a fair comparison, the SNR thresholds  $\gamma_i^{\text{OFDMA}}$ ,  $i \in \mathcal{I}$  are chosen to satisfy  $r_i = B \log_2(1 + \gamma_i^{\text{OFDMA}})$ , where  $\gamma_i^{\text{OFDMA}}$  implies the minimal SNR requirements for achieving data rate  $r_i$ .
  - *Problem formulation:* The trajectory is obtained by sequentially optimizing the position for every time step, similar to DFS with tree depth  $n = 1$ .
- **Circular-TP (Inspired by [47], [48]):** While UA, RA, and PC are optimized by Alg. 2, the UAV-BS position  $\mathbf{q}^{(t)}$  is parametrized as

$$\mathbf{q}^{(t)} = \begin{bmatrix} 100 \cos(\theta_t) + 300 \\ 100 \sin(\theta_t) + 300 \\ 200 \end{bmatrix}, \quad (27)$$

where  $\theta_t = \theta + \frac{v\Delta T}{r}t$  for randomly selected  $\theta \in [0, 2\pi]$ . Then, the variable  $\theta$  determines the initial position, and the UAV-BS rotates the orbit with angular velocity  $\frac{v\Delta T}{r}$ .

- **Fixed-TP:** For all  $t \in \mathcal{T}$ , the UAV-BS position is fixed at  $\mathbf{q}^{(t)} = [300, 300, 200]^T$ . Same as Circular-TP, this scheme determines UA, RA, and PC by Alg. 2 for given position  $\mathbf{q}^{(t)}$ .

<sup>6</sup>On average, GA-ITER converges in less than 5 iterations

<sup>7</sup>The energy consumption and relaying model in [11] are not considered in our implementation.

## A. Optimality of the Proposed RRM Scheme

TABLE III  
REGULARIZED\*  $f(\mathbf{q}^{(t)})$ (%) AND COMPLEXITY OF RRM SCHEMES

Algorithm	Number of users				Complexity (Flops)
	5	10	20	40	
GA-RRM	100	100	100	100	$\mathcal{O}(IG)**$
Alg. 2 (Ours)	99.95	99.93	99.97	99.99	$\mathcal{O}(I^2)$
Max SINR	73.54	55.19	43.91	37.78	$\mathcal{O}(I)$

\* The PF values are regularized by those of GA-RRM.

\*\*  $G$  is the number of GA generations, where  $I \ll G$ .

Once the proposed RRM scheme accomplishes global optimality, optimal decision-making for the proposed MDP formulation (as illustrated in Fig. 4) provides a global optimum of the lower-bound problem  $\mathcal{P}2$ .

Table III shows the regularized PF of the three RRM schemes to compute  $f(\mathbf{q}^{(t)})$ . The number of users is assumed to be  $\{5, 10, 20, 40\}$ ; QoS thresholds are configured as  $r_i = 5, i \in \mathcal{I}$  (Mbps); and initial data  $R_i^{(0)}, i \in \mathcal{I}$  is randomly chosen from  $\mathcal{U}[10, 30]$  (Mb).

The proposed RRM scheme tightly achieves the global-optimal PF value, which is computed by GA-RRM. Meanwhile, the PF values of the MAX SINR scheme decrease as the number of users increases. From a PF standpoint, this is because providing service to multiple users, even at the cost of some sum-rate sacrifice, proves a superior strategy compared to prioritizing spectral efficiency.

## B. Measuring the Curse of Initialization

Figure 7 illustrates PF differences in GA-TP and GA-ITER. GA-TP and GA-ITER share the same GA parameters for TP; and Alg. 2 for RRM. However, the PF gap increases when the map width and minimum data requirements increase. These factors affect both the user pathloss and candidate users whose QoS might be satisfied, so consequently determining user association within a given trajectory. Therefore, GA-ITER is more likely to converge to the local optimum as the environment becomes rigorous, which leads to an increase in the PF gap.

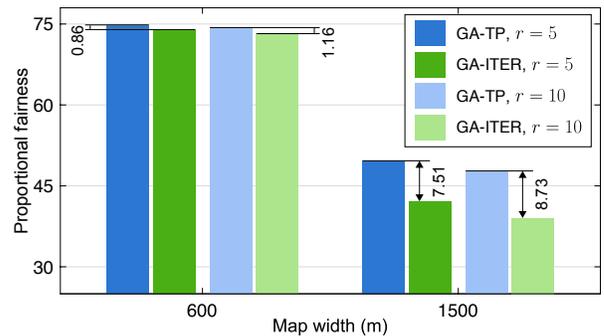


Fig. 7. PF values of GA-ITER and GA-TP with bandwidth  $B = 10$  (MHz) and minimum data rate  $r_i = r$  for  $i \in \mathcal{I}$ .

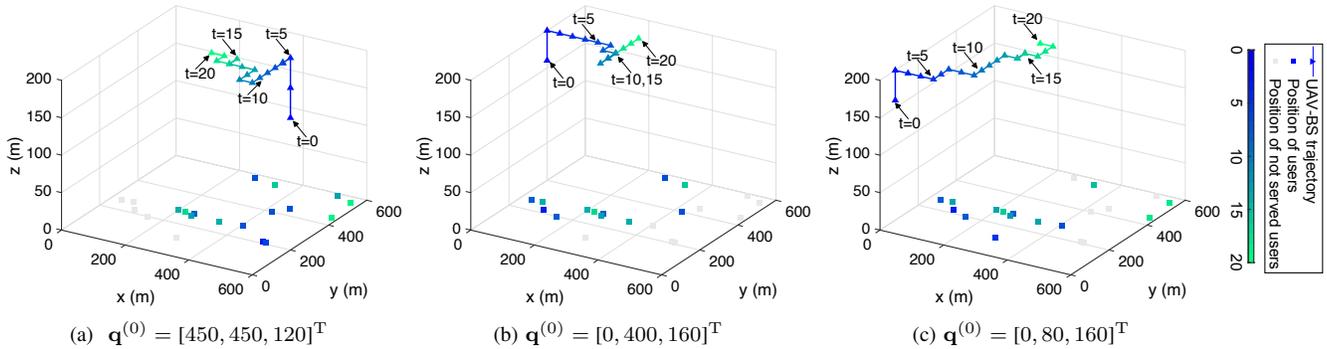


Fig. 8. Trajectories obtained from DQN with the same environment configurations, except for the initial position  $\mathbf{q}^{(0)}$ .

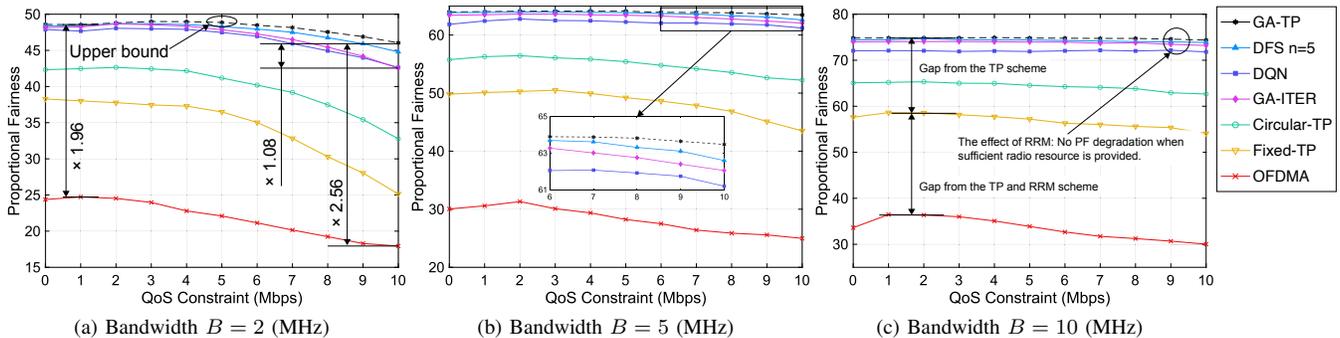


Fig. 9. Proportional fairness for 20 users with various QoS constraints.

### C. Adaptive Trajectory-Planning of MDP formulation

Figure 8 demonstrates that the proposed MDP formulation effectively resolves the dependency on the initial choice. In the early motivation presented in Fig. 1, we have argued that the initial position choice fetters subsequent variable optimization, thereby leading to a deterioration in network utility. However, Fig. 8 illustrates that the MDP formulation enables the UAV-BS to make an adaptive trajectory regardless of the initial position.

### D. Proportional Fairness for Various QoS Constraints

Figure 9a depicts the PF for 20 users under various QoS constraints when the bandwidth is 2 MHz. The figure shows that the three proposed TP schemes achieve more than 90% higher PF value than the OFDMA scheme in the worst case. This is because the OFDMA scheme maximizes the PF in the format of the weighted sum-rate. Then, the OFDMA scheme's objective becomes linear in terms of the users' sum-rate, providing all resources to the most effective user could be an optimal strategy. However, the strategy decreases the fairness of the time-critical mobile network as the users' request periods expire while a few users are selectively serviced.

We remark that the proposed schemes show a smaller decline in the PF value than the comparison schemes, even at the stringent QoS constraints. Compared with the greatest PF value of each scheme in Fig. 9a, the PF values at QoS constraints of 10 Mbps decrease 8% for GA-TP and DFS; 11% for DQN; 14% for GA-ITER 23% for Circular-TP; 34% for Fixed-TP; and 27% for the OFDMA scheme, respectively.

Figs. 9b and 9c depict the PF values with 5 and 10 MHz bandwidth, respectively. The proposed schemes show consistent PF values regardless of the QoS constraints, while the PF values of Circular-TP and Fixed-TP slightly decrease. This implies that the proposed schemes benefit from the high SNR channels because Circular-TP and Fixed-TP share the same RRM scheme with GA-TP.

### E. Percentage of Served Users for Various QoS Constraints

Figure 10 illustrates the percentage of served users (%UE) for various QoS constraints. The percentage is calculated by dividing the count of users served at least once by the total number of users. The overall tendency shows that the %UE decreases as the QoS constraint increases.

The GA-TP scheme serves at least {56, 57, 40}%p more users even in the worst case than the OFDMA scheme for bandwidth {2, 5, 10} MHz, respectively. This gap is mainly derived from the different formulations of the objective function. As mentioned in the previous section V-D, serving a few users could be an optimal UA policy in the weight sum-rate formulation. However, the objective of the proposed problem is a summation of logarithms. Then, concentrating all resources on a single user is not optimal because the growth of the logarithmic function diminishes with the increasing user sum-rate. Therefore, multiple users are highly likely associated with the UAV-BS at a single time slot in the proposed scheme.

For Circular-TP and Fixed-TP, the gap from GA-TP increases as users' QoS constraints increase. The proposed scheme shows a small gap in unconstrained cases, but the gap

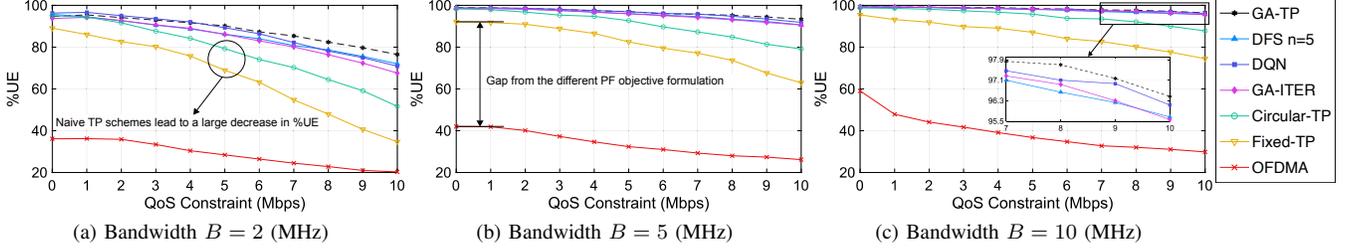


Fig. 10. Percentage of served users with various QoS constraints. The total number of users is configured to be 20.

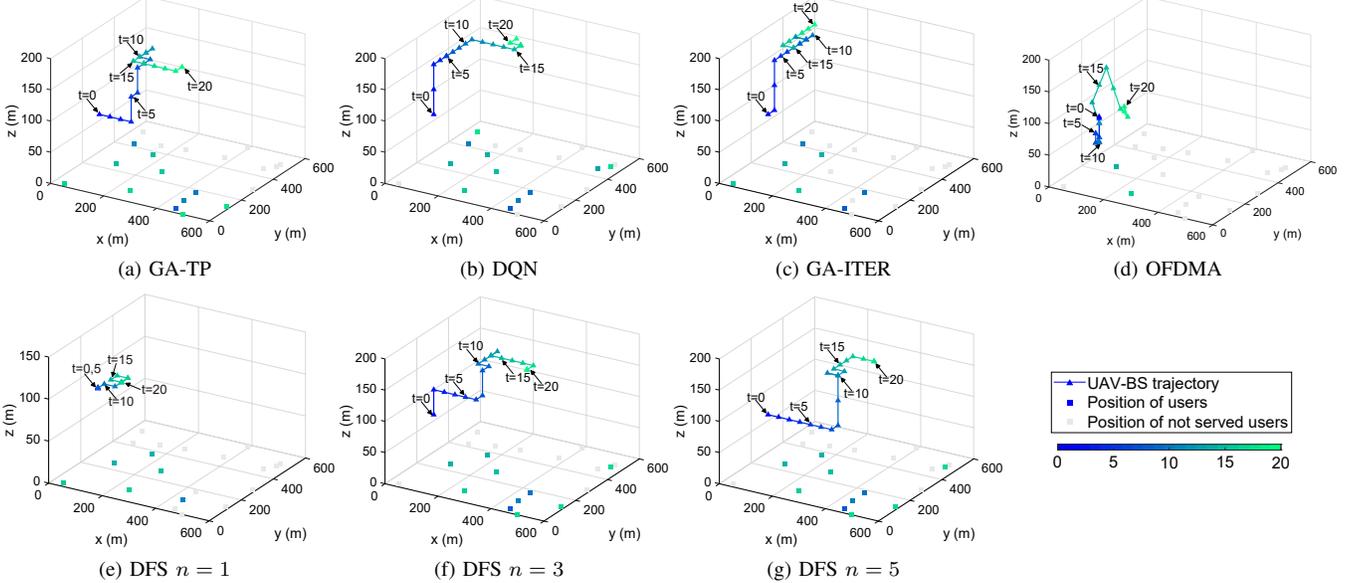


Fig. 11. Visualized trajectories for various TP schemes. The colors of the UAV-BS trajectory represent the index of time slots, and the colors in the position of users represent the first time slot where the user receives service. The trajectory of the OFDMA scheme is not latticed as we preserve the system models of [10], [11].

increases up to 20%p for Circular-TP and 37%p for Fixed-TP as QoS constraints are configured to be  $r_i = 10$ ,  $\forall i \in \mathcal{I}$ .

As depicted in Figs. 10b and 10c, the proposed schemes gradually converge to serving 100% users regardless of the lookahead variable  $n$  when the available bandwidth increases. However, the OFDMA scheme shows little enhancements due to the winner-take-all tendency in the user association.

#### F. Comparisons of the Trajectory and User Association

Figure 11 demonstrates the trajectories and user associations of various schemes. We assume that 20 users exist in the service area with QoS constraints  $r_i = 10$  Mbps for  $i \in \mathcal{I}$  with the 2 MHz bandwidth.

As we argue in the previous two sections, we can observe that the OFDMA scheme takes the winner-take-all strategy, serving only the adjacent users adjacent to the UAV-BS's position in Fig. 11d.

In Figs. 11a, 11b, and 11g, we can observe a similar UA and TP tendency as the schemes closely converge to the optimum. The number of served users for DFS accordingly increases as  $n$  increases. This is because UAV-BS can be pre-positioned near users on demand, taking into account the user's future requirements. The difference in trajectories

accumulates along the service timeline, so there is a big gap in the number of served users between  $n = 1$  and  $n = 5$  as the service time of the UAV-BS increases.

#### G. Proportional Fairness for Various Number of Users

Figure 12 presents the PF for  $\{10, 20, \dots, 80\}$  users with  $\{2, 5, 10\}$  MHz bandwidth. The QoS constraints are assumed to be  $r_i = 5$  Mbps for all  $i \in \mathcal{I}$ .

The PF of the  $\{GA-TP, DFS, DQN, GA-ITER\}$  schemes increases up to  $\{296, 295, 290, 290\}\%$  as the number of users increases (Fig. 12a). Also, Circular-TP and Fixed-TP have similar increases, but they show 18% and 40% lower PF than DFS, respectively. This implies that the proposed schemes benefit from the high SNR channels obtained from the TP scheme because Circular-TP and Fixed-TP share the same RRM scheme with the proposed schemes. Meanwhile, the PF of the OFDMA scheme shows a 41% increase.

As illustrated in Fig. 12b, the PF for 5 MHz bandwidth shows a similar trend with that for 2 MHz bandwidth (Fig. 12a). However, the PF gap between the proposed schemes and the OFDMA scheme increases as the available bandwidth increases. The GA-TP and DFS schemes show around 60% higher PF value at 10 users and 363% higher PF

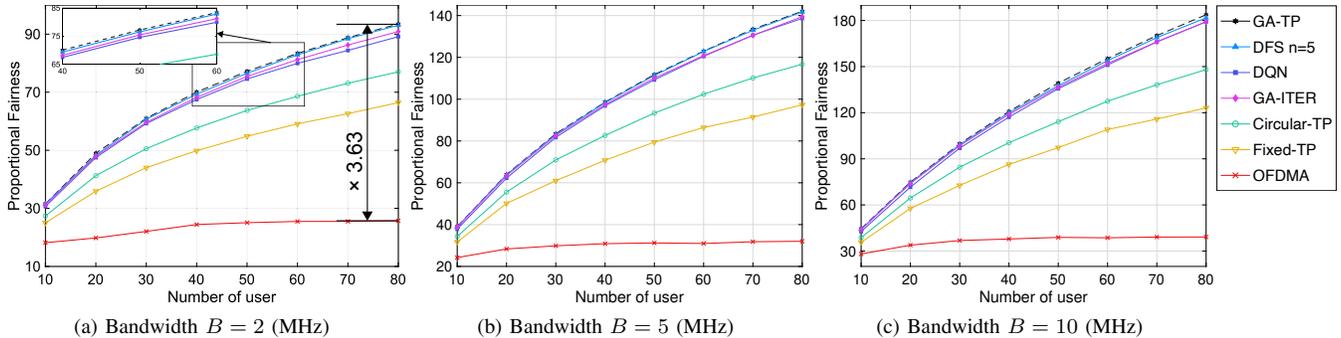


Fig. 12. Proportional fairness with a various number of users. The rate constraints  $r_i$ ,  $\forall i \in \mathcal{I}$  are configured as  $r_i = 5$ ,  $\forall i \in \mathcal{I}$ .

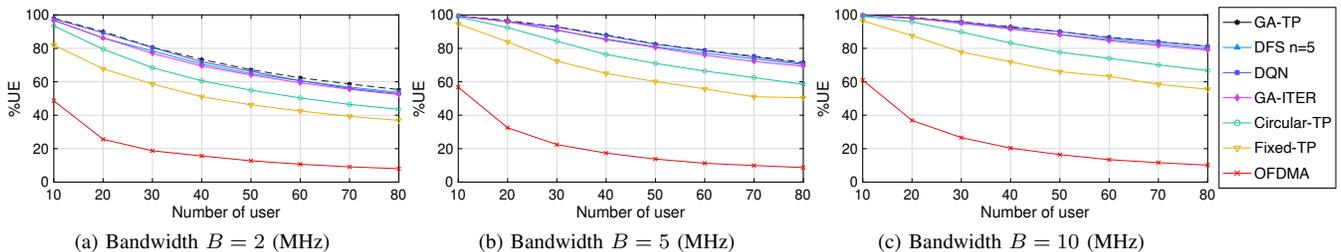


Fig. 13. Percentage of served users with a various number of users.

value at 80 users than the OFDMA scheme. Similarly, when the UAV-BS has 10 MHz bandwidth, the PF gap between the DFS and OFDMA scheme increases from 57% for 10 users to 264% for 80 users in Fig. 12c.

#### H. Percentage of Served Users for the Various Number of Users

Figure 13 depicts %UE for  $\{10, 20, \dots, 80\}$  users with  $\{2, 5, 10\}$  MHz bandwidth. The QoS constraints are configured to be  $r_i = 5$  Mbps for all  $i \in \mathcal{I}$ .

We can compute the average number of served users from Fig. 13. In Fig. 13a, the average number of served users increases from 9.8 to 44.24 for GA-TP; from 9.68 to 42.93 for DFS; from 9.76 to 42.26 for DQN; and from 9.66 to 41.85 for GA-ITER. Meanwhile, the average number of served users for the OFDMA scheme increases from 4.85 to 6.43. This implies that the increase in Fig. 12a of the OFDMA scheme is mainly caused by the channel enhancement from the high user density, not by the increase in the number of served users.

In Fig. 13b, we note that the overall percentage of served users increases when the UAV-BS utilizes 5 MHz bandwidth.  $\{\text{GA-TP, DFS, DQN}\}$  achieve about  $\{16, 17, 18\}\%$  higher %UE for 5 MHz bandwidth than for 2 MHz, but the OFDMA scheme shows almost no difference regardless of the size of available bandwidth. Figure 13c shows that the percentage of served users gradually approaches 100%, having a similar tendency with Figs. 13a and 13b.

#### I. Convergence and Computational Complexity Analysis

Figure 14 depicts the regularized PF for each iteration in the RRM scheme (Alg. 2). Zero iteration implies that the UA

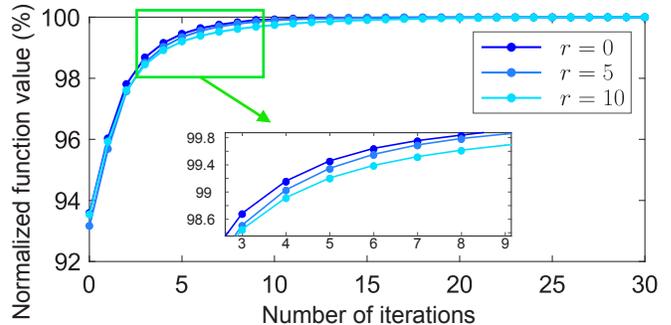


Fig. 14. The output value of Alg. 2, normalized by the convergence value. The QoS constraints are configured as  $r_i = r$ ,  $\forall i \in \mathcal{I}$ .

and RA variables are only optimized by the initial UA and RA scheme (Alg. 1) without implementing Alg. 2. Even at the zero iteration, Alg. 1 achieves 93% of the regularized PF value. Then, the PF accomplishes 99% and 99.9% after 5 and 10 iterations, respectively.

## VI. CONCLUSION

This paper proposes a separation of control and RRM in the time-critical aerial network. We find a critical drawback of conventional joint iterative optimizations, where the initial trajectory choice fetters the subsequent variable choices. However, separating TP and RRM in the practical scenario is challenging because the user requirements and fairness tightly bind network variables. We address the challenge by decomposing the proposed problem into the sum of serial sub-problems. Then, we cleverly transform the sub-problems into the MDP formulation where the TP and RRM variables are separately optimized. The proposed methods achieve the

global optimum obtained by the genetic algorithm, outperforming the state-of-the-art scheme.

The proposed approach focuses on designing a non-differentiable trajectory by integer programming. We underscore that applying the non-iterative approach to build a differentiable trajectory is another important open problem. Then continuous decision-making algorithms such as deep deterministic policy gradient [49] should be adopted.

Optimizing multi-UAV networks through the MDP formulation is another open challenge. A proliferation of research on multi-UAV networks has adopted iterative optimization, and there might be a gap in network utility due to the initial trajectories (or position) choice problem.

We anticipate that the overall network utility could be enhanced by applying the non-iterative philosophy.

## REFERENCES

- [1] ITU-R WP5D, "M.2160: framework and overall objectives of the future development of IMT for 2030 and beyond," ITU Radiocommunication Sector (ITU-R), ITU-R recommendations, Nov. 2023. [Online]. Available: <https://www.itu.int/rec/R-REC-M.2160/en>
- [2] G. Karabulut Kurt, M. G. Khoshkholgh, S. Alfattani, A. Ibrahim, T. S. J. Darwish, M. S. Alam, H. Yanikomeroglu, and A. Yongacoglu, "A vision and framework for the high altitude platform station (HAPS) networks of the future," *IEEE Commun. Surveys Tuts.*, vol. 23, no. 2, pp. 729–779, 2021.
- [3] 3GPP, "NR; radio resource control (RRC); protocol specification," 3rd Generation Partnership Project (3GPP), Technical Specification (TS) 38.831, July 2023, v17.5.0.
- [4] —, "Solution for NR to support non-terrestrial networks," 3rd Generation Partnership Project (3GPP), Technical Specification (TS) 38.821, April 2023, v16.2.0.
- [5] D. Tyrovolas, P.-V. Mekikis, S. A. Tegos, P. D. Diamantoulakis, C. K. Liaskos, and G. K. Karagiannidis, "Energy-aware design of UAV-mounted RIS networks for IoT data collection," *IEEE Trans. Commun.*, vol. 71, no. 2, pp. 1168–1178, 2023.
- [6] W. Jiang, B. Ai, M. Li, W. Wu, and X. Shen, "Average age-of-information minimization in aerial IRS-assisted data delivery," *IEEE Internet Things J.*, vol. 10, no. 17, pp. 15 133–15 146, 2023.
- [7] H. Zhou, M. Erol-Kantarci, Y. Liu, and H. V. Poor, "A survey on model-based, heuristic, and machine learning optimization approaches in RIS-aided wireless networks," *IEEE Commun. Surveys Tuts.*, pp. 1–1, 2023.
- [8] P. Tseng, "Convergence of a block coordinate descent method for nondifferentiable minimization," *J. Optim. Theory Appl.*, vol. 109, pp. 475–494, 01 2001.
- [9] S. Zhang and N. Ansari, "3D drone base station placement and resource allocation with FSO-based backhaul in hotspots," *IEEE Trans. Veh. Technol.*, vol. 69, no. 3, pp. 3322–3329, 2020.
- [10] S. Zeng, H. Zhang, and L. Song, "Trajectory optimization and resource allocation for multi-user OFDMA UAV relay networks," in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, Dec. 2019, pp. 1–6.
- [11] S. Zeng, H. Zhang, B. Di, and L. Song, "Trajectory optimization and resource allocation for OFDMA UAV relay networks," *IEEE Trans. Wireless Commun.*, vol. 20, no. 10, pp. 6634–6647, 2021.
- [12] C. Shen, T. Chang, J. Gong, Y. Zeng, and R. Zhang, "Multi-UAV interference coordination via joint trajectory and power control," *IEEE Trans. Signal Processing*, vol. 68, pp. 843–858, 2020.
- [13] O. Abbasi, H. Yanikomeroglu, A. Ebrahimi, and N. M. Yamchi, "Trajectory design and power allocation for drone-assisted NR-V2X network with dynamic NOMA/OMA," *IEEE Trans. Wireless Commun.*, vol. 19, no. 11, pp. 7153–7168, 2020.
- [14] H. Hu, K. Xiong, G. Qu, Q. Ni, P. Fan, and K. B. Letaief, "AoI-minimal trajectory planning and data collection in UAV-assisted wireless powered IoT networks," *IEEE Internet Things J.*, vol. 8, no. 2, pp. 1211–1223, 2021.
- [15] Y. Wu, W. Yang, X. Guan, and Q. Wu, "Energy-efficient trajectory design for UAV-enabled communication under malicious jamming," *IEEE Wireless Commun. Lett.*, vol. 10, no. 2, pp. 206–210, 2021.
- [16] G. Yang, R. Dai, and Y.-C. Liang, "Energy-efficient UAV backscatter communication with joint trajectory design and resource optimization," *IEEE Trans. Wireless Commun.*, vol. 20, no. 2, pp. 926–941, 2021.
- [17] Y. Gao, H. Tang, B. Li, and X. Yuan, "Joint trajectory and power design for UAV-enabled secure communications with no-fly zone constraints," *IEEE Access*, vol. 7, pp. 44 459–44 470, 2019.
- [18] M. T. Nguyen and L. B. Le, "Multi-UAV trajectory control, resource allocation, and NOMA user pairing for uplink energy minimization," *IEEE Internet Things J.*, vol. 9, no. 23, pp. 23 728–23 740, 2022.
- [19] W. Du, T. Wang, H. Zhang, Y. Dong, and Y. Li, "Joint resource allocation and trajectory optimization for completion time minimization for energy-constrained UAV communications," *IEEE Trans. Veh. Technol.*, vol. 72, no. 4, pp. 4568–4579, 2023.
- [20] W. Lu, Y. Ding, Y. Gao, S. Hu, Y. Wu, N. Zhao, and Y. Gong, "Resource and trajectory optimization for secure communications in dual unmanned aerial vehicle mobile edge computing systems," *IEEE Trans. Ind. Inform.*, vol. 18, no. 4, pp. 2704–2713, 2022.
- [21] J. Liu, Z. Xu, and Z. Wen, "Joint data transmission and trajectory optimization in UAV-enabled wireless powered mobile edge learning systems," *IEEE Trans. Veh. Technol.*, vol. 72, no. 9, pp. 11 617–11 630, 2023.
- [22] H. Yang, R. Ruby, Q.-V. Pham, and K. Wu, "Aiding a disaster spot via multi-UAV-based IoT networks: Energy and mission completion time-aware trajectory optimization," *IEEE Internet Things J.*, vol. 9, no. 8, pp. 5853–5867, 2022.
- [23] P. Yi, L. Zhu, L. Zhu, Z. Xiao, Z. Han, and X.-G. Xia, "Joint 3-D positioning and power allocation for UAV relay aided by geographic information," *IEEE Trans. Wireless Commun.*, vol. 21, no. 10, pp. 8148–8162, 2022.
- [24] G. Zhang, Q. Wu, M. Cui, and R. Zhang, "Securing uav communications via joint trajectory and power control," *IEEE Trans. Wireless Commun.*, vol. 18, no. 2, pp. 1376–1389, 2019.
- [25] M. Samir, S. Sharafeddine, C. M. Assi, T. M. Nguyen, and A. Ghrayeb, "UAV trajectory planning for data collection from time-constrained IoT devices," *IEEE Trans. Wireless Commun.*, vol. 19, no. 1, pp. 34–46, 2020.
- [26] A. Al-Hilo, M. Samir, M. Elhattab, C. Assi, and S. Sharafeddine, "RIS-assisted UAV for timely data collection in IoT networks," *IEEE Systems Journal*, pp. 1–12, 2022.
- [27] Y. Zhu and S. Wang, "Efficient aerial data collection with cooperative trajectory planning for large-scale wireless sensor networks," *IEEE Trans. Commun.*, vol. 70, no. 1, pp. 433–444, 2022.
- [28] C. Zhan, Y. Zeng, and R. Zhang, "Trajectory design for distributed estimation in UAV-enabled wireless sensor network," *IEEE Trans. Veh. Technol.*, vol. 67, no. 10, pp. 10 155–10 159, 2018.
- [29] H. Bayerlein, P. D. Kerret, and D. Gesbert, "Trajectory optimization for autonomous flying base station via reinforcement learning," in *'18 Proc. IEEE 19th Int. Workshop on Signal Processing Advances in Wireless Commun. (SPAWC)*, Aug. 2018, pp. 1–5.
- [30] P. Tong, J. Liu, X. Wang, B. Bai, and H. Dai, "Deep reinforcement learning for efficient data collection in UAV-aided internet of things," in *Proc. IEEE Int. Conf. on Commun. Workshops (ICC Workshops)*, June 2020, pp. 1–6.
- [31] C. You and R. Zhang, "3D trajectory optimization in rician fading for UAV-enabled data harvesting," *IEEE Trans. Wireless Commun.*, vol. 18, no. 6, pp. 3192–3207, 2019.
- [32] B. Zhu, E. Bedeer, H. H. Nguyen, R. Barton, and J. Henry, "UAV trajectory planning in wireless sensor networks for energy consumption minimization by deep reinforcement learning," *IEEE Trans. Veh. Technol.*, vol. 70, no. 9, pp. 9540–9554, 2021.
- [33] X. Liu, Y. Liu, Y. Chen, and L. Hanzo, "Trajectory design and power control for multi-UAV assisted wireless networks: A machine learning approach," *IEEE Trans. Veh. Technol.*, vol. 68, no. 8, pp. 7957–7969, 2019.
- [34] A. Al-Hourani, S. Kandeepan, and S. Lardner, "Optimal LAP altitude for maximum coverage," *IEEE Commun. Lett.*, vol. 3, no. 6, pp. 569–572, 2014.
- [35] J. Holis and P. Pechac, "Elevation dependent shadowing model for mobile communications via high altitude platforms in built-up areas," *IEEE Trans. Antennas Propag.*, vol. 56, no. 4, pp. 1078–1084, 2008.
- [36] A. Al-Hourani, S. Kandeepan, and A. Jamalipour, "Modeling air-to-ground path loss for low altitude platforms in urban environments," in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, Dec. 2014, pp. 2898–2904.

- [37] S. Burer and A. N. Letchford, “Non-convex mixed-integer nonlinear programming: A survey,” *Surv. Oper. Res. Manag. Sci.*, vol. 17, no. 2, pp. 97–106, 2012.
- [38] A. Papoulis and S. Pillai, *Probability, Random Variables, and Stochastic Processes*, ser. McGraw-Hill series in electrical and computer engineering. McGraw-Hill, 2002.
- [39] T. Weise, “Global optimization algorithms-theory and application,” *Self-Published Thomas Weise*, vol. 361, 2009.
- [40] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, “Elementary graph algorithms,” in *Introduction to Algorithms*, 3rd ed. The MIT Press, 2009, ch. 22, pp. 589–623.
- [41] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. A. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, pp. 529–533, 2015.
- [42] D. W. Matolak and R. Sun, “Air-ground channel characterization for unmanned aircraft systems—part III: The suburban and near-urban environments,” *IEEE Trans. Veh. Technol.*, vol. 66, no. 8, pp. 6607–6618, 2017.
- [43] 3GPP, “Typical traffic characteristics of media services on 3GPP networks,” 3rd Generation Partnership Project (3GPP), Technical Specification (TS) 26.925, 04 2022, version 17.1.0.
- [44] J. Chakareski, S. Naqvi, N. Mastrorade, J. Xu, F. Afghah, and A. Razi, “An energy efficient framework for UAV-assisted millimeter wave 5G heterogeneous cellular networks,” *IEEE Trans. Green Commun. Netw.*, vol. 3, no. 1, pp. 37–44, 2019.
- [45] A. Mahmood, T. X. Vu, S. Chatzinotas, and B. Ottersten, “Joint optimization of 3D placement and radio resource allocation for per-UAV sum rate maximization,” *IEEE Trans. Veh. Technol.*, vol. 72, no. 10, pp. 13 094–13 105, 2023.
- [46] D. G. Luenberger, Y. Ye *et al.*, *Linear and nonlinear programming*. Springer, 1984, vol. 2.
- [47] Q. Wu, Y. Zeng, and R. Zhang, “Joint trajectory and communication design for multi-UAV enabled wireless networks,” *IEEE Trans. Wireless Commun.*, vol. 17, no. 3, pp. 2109–2121, 2018.
- [48] F. Ono, H. Ochiai, and R. Miura, “A wireless relay network based on unmanned aircraft system with rate optimization,” *IEEE Trans. Wireless Commun.*, vol. 15, no. 11, pp. 7699–7708, 2016.
- [49] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, “Continuous control with deep reinforcement learning,” in *4th Int. Conf. Learn. Represent., ICLR 2016, San Juan, Puerto Rico, Conf. Track Proc.*, May 2-4 2016.
- [50] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *3rd Int. Conf. on Learn. Represent., ICLR 2015, San Diego, CA, USA, Conf. Track Proc.*, May 7-9 2015.

## APPENDIX A RELATED WORKS

Recent works that suggest TP and resource management of UAV-BSs in various scenarios are categorized into the following subsections. Key characteristics of the works are summarized in Table. I.

### A. UAV-BSs in Time-Critical Mobile Scenarios

Several research works have studied the control of UAV-BSs in time-critical wireless sensor networks [14], [25], [26], [30]. The works in [25], [26] target to maximize the number of devices that successfully transmit data within the lifetime of data. In [25], two algorithms are suggested to jointly design a 2-dimensional trajectory and bandwidth allocation of devices. Reconfigurable intelligent surfaces are additionally considered in [26] and jointly optimized with the trajectory of a UAV-BS by using deep-reinforcement learning. In [14], the authors propose an age of information minimization problem in a wireless sensor network, where sensors are powered by

a UAV-BS using energy harvesting. These works enhance the network utility for timely data, but the fairness between the users is not considered. Without considering fairness, the trajectory of a UAV-BS can be biased according to the QoS constraints and channel state of users, thereby decreasing the number of served IoT users.

### B. UAV-BSs for Wireless Sensor Networks

The work in [31] addresses a problem that maximizes the minimum average data rate in wireless sensor networks under an outage probability constraint and angle-dependent Rician fading. In [28], the authors reformulate the TP problem into an equivalent traveling salesman problem to maximize the number of visited wireless sensors. The authors of [27] suggest a TP algorithm that minimizes the required time slots for collecting data from all wireless sensors. In [32], a deep-reinforcement learning model is adopted to minimize energy consumption while collecting data from wireless sensors. These studies provide a well-designed trajectory in various mobile network scenarios, but designing optimal RA and PC still remains unsolved, which should be jointly designed with the trajectory.

### C. TP of UAV-BSs

The works in [12], [13], [29], [33] propose a TP problem maximizing the sum-rate of users in UAV-enabled networks. These works significantly enhance the sum-rate by adopting reinforcement-learning (RL) methods [29], [33] or utilizing successive convex approximation [12], [13]. However, none of these works considers the fairness between users or the RA problem.

The authors in [10], [11] propose a TP algorithm that jointly optimizes RA and PC parameters to maximize the fairness of users in the orthogonal frequency-division multiplexing access (OFDMA) network, where a single UAV-BS is utilized as a relay node. These algorithms sequentially find the position, RA, and PC parameters by optimizing these parameters for the next time slot based on the current location of the UAV-BS. The problems designed in these works are well-investigated with solid results, but future network states still need to be considered when optimizing the current movement of the UAV-BS.

## APPENDIX B PROOF OF PROPOSITION 1

The logarithm of the sum-rate of  $i$ -th user is reformulated as

$$\log R_i = \log \sum_{t \in \mathcal{T}} R_i^{(t)} \quad (28)$$

$$= \log \prod_{t=1}^T \left( \frac{\sum_{k=0}^t R_i^{(k)}}{\sum_{k=0}^{t-1} R_i^{(k)}} \right) \cdot R_i^{(0)} \quad (29)$$

$$= \sum_{t=1}^T \log \left( 1 + \frac{R_i^{(t)}}{\sum_{k=0}^{t-1} R_i^{(k)}} \right) + \log R_i^{(0)}, \quad (30)$$

where the auxiliary constant  $R_i^{(0)}$  is identical for all  $i \in \mathcal{I}$ . In what follows, we shall adopt the auxiliary constant  $R_i^{(0)}$  to consistently decouple the problem for all time slots<sup>8</sup>. For the numerical experiments,  $R_i^{(0)}$ ,  $i \in \mathcal{I}$  are chosen to be 1, but the value of  $R_i^{(0)}$  rarely affects on the outcome of the proposed algorithms. For each user  $i$ ,  $R_i^{(0)}$  determines the cumulative throughput  $\sum_{k=0}^t R_i^{(k)}$ , which affects the solution of RA and PC. A specific range of  $R_0$  can affect the short-term RRM; but does not have a significant long-term effect, as the cumulative throughput naturally grows as time goes by.

The objective in (14a) is equivalent with

$$\sum_{i \in \mathbf{A}} \log R_i = \sum_{t=1}^T \sum_{i \in \mathbf{A}^{(t)}} \log \left( 1 + \frac{R_i^{(t)}}{\sum_{k=0}^{t-1} R_i^{(k)}} \right) + \tilde{R}, \quad (31)$$

where  $\tilde{R} = \sum_{i \in \mathbf{A}} \log R_i^{(0)}$ . Then, the following inequality holds:

$$\max_{\mathbf{q}, \mathbf{A}, \mathbf{B}, \mathbf{P}} \sum_{i \in \mathbf{A}} \log R_i \quad (32)$$

$$= \max_{\mathbf{q}, \mathbf{A}, \mathbf{B}, \mathbf{P}} \sum_{t=1}^T \sum_{i \in \mathbf{A}^{(t)}} \log \left( 1 + \frac{R_i^{(t)}}{\sum_{k=0}^{t-1} R_i^{(k)}} \right) + \tilde{R} \quad (33)$$

$$\geq \sum_{t=1}^T \max_{\mathbf{q}^{(t)}, \mathbf{A}^{(t)}, \mathbf{B}^{(t)}, \mathbf{P}^{(t)}} \sum_{i \in \mathbf{A}^{(t)}} \log \left( 1 + \frac{R_i^{(t)}}{\sum_{k=0}^{t-1} R_i^{(k)}} \right). \quad (34)$$

The inequality (34) indicates that the addition of the objectives of Problem  $\mathcal{P}^{(t)}$ ,  $t \in \mathcal{T}$  is the lower bound of Problem  $\mathcal{P}1$ .

#### APPENDIX C DERIVATION OF THE INITIAL $\mathbf{B}^{(t)}$

If  $\mathbf{A}^{(t)}$  is given, Problem (21) is concave; thus can be solved by finding the KKT conditions. For Problem (21a), the Lagrangian  $\mathcal{L}_{\mathbf{B}|\mathbf{A}}(\cdot)$  is represented as

$$\begin{aligned} \mathcal{L}_{\mathbf{B}|\mathbf{A}}(\mathbf{B}^{(t)}, \boldsymbol{\nu}^{(t)}, \lambda^{(t)}) &= \sum_{i \in \mathbf{A}^{(t)}} \log(1 + w_i^{(t)} \beta_i^{(t)}) \quad (35) \\ &+ \sum_{i \in \mathbf{A}^{(t)}} \nu_i^{(t)} \left( \beta_i^{(t)} - \frac{\alpha_i^{(t)} r_i}{e_i^{(t)}} \right) - \lambda^{(t)} \left( \sum_{i \in \mathbf{A}^{(t)}} \beta_i^{(t)} - B \right), \end{aligned} \quad (36)$$

where  $w_i^{(t)} = \alpha_i^{(t)} d_i^{(t)} e_i^{(t)} / \sum_{k=0}^{t-1} R_i^{(k)}$ ;  $\boldsymbol{\nu}^{(t)} = [\nu_i^{(t)}]_{i \in \mathcal{I}}$  and  $\lambda^{(t)}$  are Lagrangian coefficient.

To meet the first-order optimality, we have

$$\nabla_{\beta_i^{(t)}} \mathcal{L}_{\mathbf{B}|\mathbf{A}}(\cdot) = \frac{w_i^{(t)}}{1 + w_i^{(t)} \beta_i^{(t)}} + \nu_i^{(t)} - \lambda^{(t)} = 0. \quad (37)$$

$$\nu_i^{(t)} \left( \beta_i^{(t)} - \frac{\alpha_i^{(t)} r_i}{e_i^{(t)}} \right) = 0. \quad (38)$$

<sup>8</sup>Otherwise, we have a different term at  $t = 1$ .

---

#### Algorithm 3: Resource Allocation

---

```

1: Input  $\mathbf{q}^{(t)}$ ,  $\mathbf{A}^{(t)}$ ,  $\mathbf{P}^{(t)}$ .
2: Output  $\mathbf{B}^{(t)}$ .
3: Initialize  $\boldsymbol{\mu}^{(t)}$ ,  $\lambda_1^{(t)}$ ,  $\lambda_2^{(t)}$ ,  $\hat{\mathcal{L}} \leftarrow 0$ ,  $\hat{\mathcal{L}}_{\text{prev}} \leftarrow \infty$ 
4: while  $\|\hat{\mathcal{L}} - \hat{\mathcal{L}}_{\text{prev}}\| > \epsilon$  do
5:   Update  $\boldsymbol{\mu}^{(t)}$ ,  $\lambda_1^{(t)}$ ,  $\lambda_2^{(t)}$  according to (45)-(47).
6:    $\hat{\mathcal{L}}_{\text{prev}} \leftarrow \hat{\mathcal{L}}$ ,  $\hat{\mathcal{L}} \leftarrow \mathcal{L}_{\beta}^*(\boldsymbol{\mu}^{(t)}, \lambda_1^{(t)}, \lambda_2^{(t)})$  using (44).
7: end
8:  $\mathbf{B}^{(t)} \leftarrow \left[ \frac{1}{\lambda_1^{(t)} + \lambda_2^{(t)} \rho_i^{(t)} - \mu_i^{(t)}} - \frac{1}{w_i^{(t)}} \right]_{i \in \mathcal{I}}$ .

```

---

By substituting  $\nu_i^{(t)}$  in (38) with (37), the following inequality holds:

$$\left( \lambda^{(t)} - \frac{w_i^{(t)}}{1 + w_i^{(t)} \beta_i^{(t)}} \right) \left( \beta_i^{(t)} - \frac{\alpha_i^{(t)} r_i}{e_i^{(t)}} \right) = 0. \quad (39)$$

For users  $i \in \mathbf{A}^{(t)}$ , combining (39) and the minimum RA requirements constraints (21b) results in

$$\beta_i^{(t)} = \max \left( \frac{r_i}{e_i^{(t)}}, \frac{1}{\lambda^{(t)}} - \frac{1}{w_i^{(t)}} \right). \quad (40)$$

Because  $\beta_i^{(t)}$  is simply determined to be zero for  $i \notin \mathbf{A}^{(t)}$ , we can conclude

$$\beta_i^{(t)} = \begin{cases} \max \left( \frac{r_i}{e_i^{(t)}}, \frac{1}{\lambda^{(t)}} - \frac{1}{w_i^{(t)}} \right) & \text{if } \alpha_i^{(t)} d_i^{(t)} = 1, \\ 0 & \text{otherwise.} \end{cases} \quad (41)$$

#### APPENDIX D OPTIMAL RESOURCE ALLOCATION

If  $\mathbf{A}^{(t)}$  and  $\mathbf{P}^{(t)}$  are fixed, the problem (20) is written as

$$\max_{\mathbf{B}^{(t)}} \sum_{i \in \mathbf{A}^{(t)}} \log \left( 1 + \frac{R_i^{(t)}}{\sum_{k=0}^{t-1} R_i^{(k)}} \right) \quad (42a)$$

$$\text{s.t.} \quad \sum_{i \in \mathbf{A}^{(t)}} \beta_i^{(t)} \leq B, \quad (42b)$$

$$\sum_{i \in \mathbf{A}^{(t)}} \rho_i^{(t)} \beta_i^{(t)} \leq P, \quad (42c)$$

$$\beta_i^{(t)} \geq \alpha_i^{(t)} r_i / e_i^{(t)}, \forall i. \quad (42d)$$

The optimal solution of the problem (42) can be obtained by using the Lagrangian dual method. Then,  $\mathbf{B}^{(t)}$  can be optimized by globally searching all feasible  $\mathbf{B}^{(t)}$  for the objective function of (42).

The Lagrangian dual problem of (42) is

$$\min_{\boldsymbol{\mu}^{(t)}, \lambda_1^{(t)}, \lambda_2^{(t)}} \mathcal{L}_{\mathbf{B}}^*(\boldsymbol{\mu}^{(t)}, \lambda_1^{(t)}, \lambda_2^{(t)}) \quad (43a)$$

$$\text{s.t.} \quad \mu_i^{(t)} \geq 0, \forall i, \quad (43b)$$

$$\lambda_1^{(t)}, \lambda_2^{(t)} \geq 0, \quad (43c)$$

where  $\boldsymbol{\mu}^{(t)} = [\mu_i^{(t)}]_{i \in \mathcal{I}}$  and

$$\begin{aligned} \mathcal{L}_{\mathbf{B}}^*(\boldsymbol{\mu}^{(t)}, \lambda_1^{(t)}, \lambda_2^{(t)}) &= - \sum_{i \in \mathbf{A}^{(t)}} \alpha_i^{(t)} - \sum_{i \in \mathbf{A}^{(t)}} \log(\lambda_1^{(t)} + \lambda_2^{(t)} \rho_i^{(t)} - \mu_i^{(t)}) \\ &\quad - \sum_{i \in \mathbf{A}^{(t)}} \mu_i^{(t)} \left( \frac{1}{w_i^{(t)}} + \frac{\alpha_i^{(t)} r_i}{e_i^{(t)}} \right) \\ &\quad + \lambda_1^{(t)} \left( \sum_{i \in \mathbf{A}^{(t)}} \frac{1}{w_i^{(t)}} + B \right) + \lambda_2^{(t)} \left( \sum_{i \in \mathbf{A}^{(t)}} \frac{\rho_i^{(t)}}{w_i^{(t)}} + P \right). \end{aligned} \quad (44)$$

The derivation of the Lagrangian dual function  $\mathcal{L}_{\mathbf{B}}^*(\boldsymbol{\mu}^{(t)}, \lambda_1^{(t)}, \lambda_2^{(t)})$  is shown in the following section Appendix D-A. The duality gap between (42) and (43) is zero, because the objective in (42a) is concave and the constraints (42b)-(42d) are linear. Then, by using the sub-gradient descent method, the Lagrangian multipliers can be updated by

$$\mu_i^{(t)} \leftarrow \mu_i^{(t)} - \gamma \left( \frac{1}{o_i^{(t)}} - \frac{1}{w_i^{(t)}} - \frac{\alpha_i^{(t)} r_i}{e_i^{(t)}} \right), \quad (45)$$

$$\lambda_1^{(t)} \leftarrow \lambda_1^{(t)} - \gamma \left( - \sum_{i \in \mathbf{A}^{(t)}} \frac{1}{o_i^{(t)}} + B + \sum_{i \in \mathbf{A}^{(t)}} \frac{1}{w_i^{(t)}} \right), \quad (46)$$

$$\lambda_2^{(t)} \leftarrow \lambda_2^{(t)} - \gamma \left( - \sum_{i \in \mathbf{A}^{(t)}} \frac{\rho_i^{(t)}}{o_i^{(t)}} + P + \sum_{i \in \mathbf{A}^{(t)}} \frac{\rho_i^{(t)}}{w_i^{(t)}} \right), \quad (47)$$

where  $\gamma$  is a learning rate and  $o_i^{(t)} = \lambda_1^{(t)} + \lambda_2^{(t)} \rho_i^{(t)} - \mu_i^{(t)}$ . The initial values of the multipliers are configured<sup>9</sup> as  $\mu_i^{(t)} = 0.5$  for all  $i \in \mathcal{I}$ ,  $\lambda_1^{(t)} = 0.5$ , and  $\lambda_2^{(t)} = 0.8$ . After that, the optimal value of  $\mathbf{B}^{(t)}$  can be obtained from

$$\beta_i^{*(t)} = \frac{1}{\lambda_1^{(t)} + \lambda_2^{(t)} \rho_i^{(t)} - \mu_i^{(t)}} - \frac{1}{w_i^{(t)}}, \quad (48)$$

where  $\mathbf{B}^{*(t)} = [\beta_i^{*(t)}]_{i \in \mathcal{I}}$ .

#### A. Derivation of the Lagrangian Dual Function $\mathcal{L}_{\mathbf{B}}^*(\boldsymbol{\mu}^{(t)}, \lambda_1^{(t)}, \lambda_2^{(t)})$ in (44)

The Lagrangian dual function (44) is derived as

$$\begin{aligned} \mathcal{L}_{\mathbf{B}}(\mathbf{B}^{(t)}, \boldsymbol{\mu}^{(t)}, \lambda_1^{(t)}, \lambda_2^{(t)}) &= \sum_{i \in \mathbf{A}^{(t)}} \log(1 + w_i^{(t)} \beta_i^{(t)}) + \sum_{i \in \mathbf{A}^{(t)}} \mu_i^{(t)} \left( \beta_i^{(t)} - \frac{\alpha_i^{(t)} r_i}{e_i^{(t)}} \right) \\ &\quad - \lambda_1^{(t)} \left( \sum_{i \in \mathbf{A}^{(t)}} \beta_i^{(t)} - B \right) - \lambda_2^{(t)} \left( \sum_{i \in \mathbf{A}^{(t)}} \rho_i^{(t)} \beta_i^{(t)} - P \right), \end{aligned} \quad (49)$$

where  $w_i^{(t)} = \alpha_i^{(t)} d_i^{(t)} e_i^{(t)} / \sum_{k=0}^{t-1} R_i^{(k)}$ .

The partial derivative of  $\mathcal{L}(\cdot)$  with respect to  $\beta_i^{(t)}$  is 0 at optimal resource allocation vector  $\mathbf{B}^{*(t)}$ , so we have

$$\nabla_{\beta_i^{(t)}} \mathcal{L}_{\mathbf{B}}(\cdot) = \frac{w_i^{(t)}}{1 + w_i^{(t)} \beta_i^{(t)}} + \mu_i^{(t)} - \lambda_1^{(t)} - \lambda_2^{(t)} \rho_i^{(t)} \quad (50)$$

$$= 0, \quad (51)$$

<sup>9</sup>We heuristically apply the bisection method to find the best combination that minimizes the convergence time.

Then, the optimal  $\beta_i^{*(t)}$  is

$$\beta_i^{*(t)} = \frac{1}{\lambda_1^{(t)} + \lambda_2^{(t)} \rho_i^{(t)} - \mu_i^{(t)}} - \frac{1}{w_i^{(t)}}, \quad (52)$$

and  $\mathbf{B}^{*(t)} = [\beta_i^{*(t)}]_{i \in \mathcal{I}}$ . Therefore, we have

$$\mathcal{L}_{\mathbf{B}}^*(\boldsymbol{\mu}^{(t)}, \lambda_1^{(t)}, \lambda_2^{(t)}) = \mathcal{L}_{\beta}(\mathbf{B}^{*(t)}, \boldsymbol{\mu}^{(t)}, \lambda_1^{(t)}, \lambda_2^{(t)}). \quad (53)$$

## APPENDIX E OPTIMAL POWER CONTROL

If  $\mathbf{A}^{(t)}$  and  $\mathbf{B}^{(t)}$  are given, the problem (20) is written as

$$\max_{\mathbf{P}^{(t)}} \sum_{i \in \mathbf{A}^{(t)}} \log \left( 1 + \frac{R_i^{(t)}}{\sum_{k=0}^{t-1} R_i^{(k)}} \right) \quad (54a)$$

$$\text{s.t. } \rho_i^{(t)} \geq 0, \forall i, \quad (54b)$$

$$\sum_{i \in \mathbf{A}^{(t)}} \rho_i^{(t)} \beta_i^{(t)} \leq P, \quad (54c)$$

$$R_i^{(t)} \geq \alpha_i^{(t)} r_i, \forall i. \quad (54d)$$

Since  $\log(1+x) \approx x$  when  $x \ll 1$ , we relax the objective function as

$$\sum_{i \in \mathbf{A}^{(t)}} \log \left( 1 + \frac{R_i^{(t)}}{\sum_{k=0}^{t-1} R_i^{(k)}} \right) \approx \sum_{i \in \mathbf{A}^{(t)}} \frac{R_i^{(t)}}{\sum_{k=0}^{t-1} R_i^{(k)}}. \quad (55)$$

We remark that the above approximation holds more accurately as  $t$  increases because the cumulative received data  $\sum_{k=0}^{t-1} R_i^{(k)}$  increases over time. Also, the constraint (54d) can be equivalently expressed with respect to the PSD of a user as follows:

$$\rho_i^{(t)} \geq \zeta_i^{(t)}, \forall i, \quad (56)$$

where  $\zeta_i^{(t)} = \frac{N_0}{10^{-\xi_i^{(t)}/10}} (2^{\alpha_i^{(t)} r_i / \beta_i^{(t)}} - 1)$ . We can replace (54b) and (54d) by (56) because  $\zeta_i^{(t)} \geq 0$ .

By applying (55) and (56), the problem (54) can be reformulated as

$$\max_{\mathbf{P}^{(t)}} \sum_{i \in \mathbf{A}^{(t)}} \frac{R_i^{(t)}}{\sum_{k=0}^{t-1} R_i^{(k)}} \quad (57a)$$

$$\text{s.t. } \sum_{i \in \mathbf{A}^{(t)}} \rho_i^{(t)} \beta_i^{(t)} \leq P, \quad (57b)$$

$$\rho_i^{(t)} \geq \zeta_i^{(t)}, \forall i. \quad (57c)$$

The objective function and constraints of (57) are concave and linear, respectively. Therefore, the zero-duality gap is guaranteed between Problem (57) and its dual problem.

The Lagrangian  $\mathcal{L}_{\mathbf{P}}$  of (57) is

$$\begin{aligned} \mathcal{L}_{\mathbf{P}}(\mathbf{P}^{(t)}, \boldsymbol{\nu}^{(t)}, \lambda^{(t)}) &= \sum_{i \in \mathbf{A}^{(t)}} \tau_i^{(t)} \log_2(1 + \omega_i^{(t)} \rho_i^{(t)}) \\ &\quad + \sum_{i \in \mathbf{A}^{(t)}} \nu_i^{(t)} (\rho_i^{(t)} - \zeta_i^{(t)}) - \lambda^{(t)} \left( \sum_{i \in \mathbf{A}^{(t)}} \rho_i^{(t)} \beta_i^{(t)} - P \right), \end{aligned} \quad (58)$$

where  $\boldsymbol{\nu}^{(t)} = [\nu_i^{(t)}]_{i \in \mathcal{I}}$ ,  $\tau_i^{(t)} = \alpha_i^{(t)} d_i^{(t)} / \sum_{k=0}^{t-1} R_i^{(k)}$  and  $\omega_i^{(t)} = 10^{\xi_i^{(t)}/10} / n_0$ .

We need to find a point at which derivative of  $\mathcal{L}_{\mathbf{P}}(\cdot)$  is zero:

$$\nabla_{\rho_i^{(t)}} \mathcal{L}_{\mathbf{P}}(\cdot) = \frac{\tau_i^{(t)} \omega_i^{(t)}}{1 + \omega_i^{(t)} \rho_i^{(t)}} + \nu_i^{(t)} - \lambda^{(t)} \beta_i^{(t)} = 0. \quad (59)$$

In addition, the pair of  $\mathbf{P}^{(t)}$ ,  $\boldsymbol{\nu}^{(t)}$  and  $\lambda^{(t)}$  should satisfy

$$\nu_i^{(t)} (\rho_i^{(t)} - \zeta_i^{(t)}) = 0 \quad (60)$$

to meet the complementary slackness of a KKT solution.

Similar to Appendix C,  $\rho_i^{(t)}$  is determined as

$$\rho_i^{(t)} = \begin{cases} \max \left( \zeta_i^{(t)}, \frac{\tau_i^{(t)}}{\beta_i^{(t)} \lambda^{(t)}} - \frac{1}{\omega_i^{(t)}} \right) & \text{if } \alpha_i^{(t)} d_i^{(t)} = 1, \\ 0 & \text{otherwise.} \end{cases} \quad (61)$$

TABLE IV  
DQN PARAMETER CONFIGURATIONS

Parameter	Value
Learning rate (LR)	$1 \cdot 10^{-4}$
LR scheduler step size	500
LR decay factor	0.9
Discount factor	0.99
Exploration prob. $\epsilon$	1
$\epsilon$ -decay per update	0.999
Soft-update weight $\tau$	$10^{-3}$
Soft-update period	1
Experience-replay buffer size	$10^5$
Experience-replay batch size	512

#### APPENDIX F

##### DQN ARCHITECTURE AND TRAINING PROCEDURE

We design a 16-layer fully connected network to produce the results in Sec. V. All layers, except for the input and output layers, have input and output vectors of size  $\mathbb{R}^{1024}$ . The input vector size is a cardinality of a state vector, determined as  $11 + 8I$  by the number of users  $I$ . The output vector size is  $\mathbb{R}^7$  for our numerical experiment, determined by the cardinality of  $S(\cdot)$ . A detailed configuration is provided on the implementation source code.

The model is trained using the adaptive moment estimation (Adam) optimizer [50], step-decaying learning rate scheduler (e.g. StepLR in PyTorch), mean squared error (MSE) loss. All hyper-parameters are listed in Table IV. The DQN models are trained during 5,000 episodes, which is equivalent to 200,000 steps. Figure 15 illustrates the training progress with the DQN environment configured as 10 MHz bandwidth; 20 users; and  $r_i = 5$  for all  $i \in \mathcal{I}$ .

In all cases, the DQN scheme shows comparable results with the GA-TP and DFS but fails to exceed the PF of the DFS scheme. The main reason is that the DFS scheme benefits from the analytical optimal reward of the sub-trajectory, but the DQN scheme should estimate the exact reward of each position  $\mathbf{q}^{(t)}$  for all  $t \in \mathcal{T}$ .

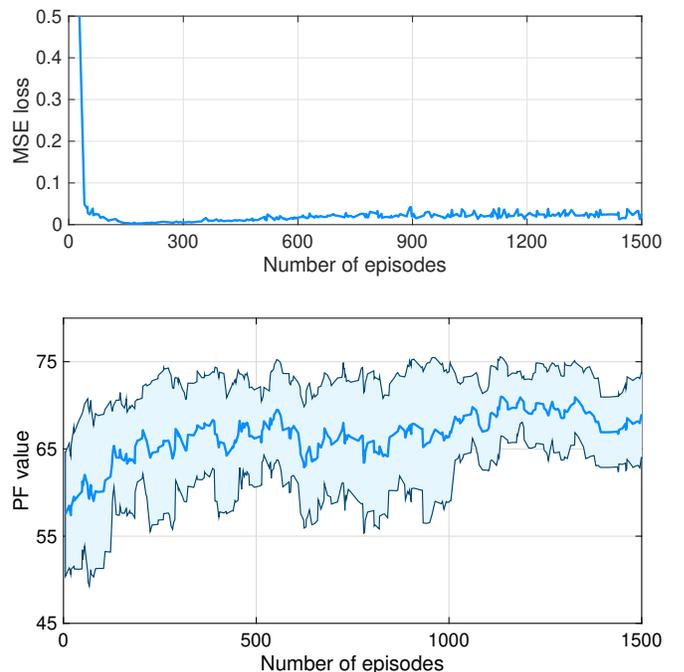


Fig. 15. Loss and PF value change for the first 1500 training episodes. The PF value is moving-averaged with the window size 13. The colored area represents the variance.