

# Multi-UAV Trajectory and Power Optimization for Cached UAV Wireless Networks With Energy and Content Recharging-Demand Driven Deep Learning Approach

Shuqi Chai<sup>ID</sup>, Associate Member, IEEE, and Vincent K. N. Lau, Fellow, IEEE

**Abstract**—In this paper, we propose a novel joint trajectory and communication scheduling scheme for multiple unmanned aerial vehicles (UAVs) enabled wireless caching networks. To exploit the favorable propagation of air-to-ground channels, we consider an ultra dense UAVs enabled content-centric wireless transmission network, where massive UAVs are deployed to transmit cached contents to a group of random distributed ground users. We formulate the problem as an infinite horizon ergodic stochastic differential game (SDG) for optimizing the users' quality-of-experience (QoE). In particular, stochastic dynamics of channel states, UAVs' mobility, energy queues and content request queues are modeled in this game. To deal with the state coupling between the UAVs, we consider a limiting problem for large number of UAV based on mean field analysis. A reduced-complexity decentralized solution can be obtained through mean-field equilibrium analysis. To further reduce the solution complexity on each UAV, we propose a model-specific deep neural network (DNN) to learn the optimal control solution in an online manner. The DNN is not arbitrarily generated but tailored to the structural properties of the value function and stationary distribution based on the homotopy perturbation method analysis. Finally, simulation results are provided to show that the proposed solution can achieve significant gain over the existing baselines.

**Index Terms**—UAV caching networks, multi-UAV trajectory design, radio resource control, mean-field game, deep learning.

## I. INTRODUCTION

IT IS widely anticipated that the capacity in 5G wireless networks will continue to increase to meet the requirements of high-data-rate demand applications, such as multimedia streaming. Currently, there is a surge of interest in utilizing unmanned aerial vehicles (UAVs) as moving base stations in wireless networks to improve the network capacity and content-delivery performance. Using UAVs-assisted wireless network has several motivations. First, traditional terrestrial base stations (BSs) are comparatively inefficient to support

mobile users with popularly requested content due to the lack of line-of-sight (LOS). Hence, they have to use the lower frequency spectrum which is heavily congested. On the other hand, UAVs can assist the BSs with a higher data transmission rate by exploiting the inherit LOS propagation, and this can allow the usage of the high frequency spectrum to serve hotspot demands. Furthermore, the agility and controllability of UAV-aided wireless networks can achieve larger wireless coverage, making them ideal to support reliable transmission for special events, such as large-scale ad-hoc activities or sports events.

Many research works have analyzed UAV-aided wireless networks in recent years. In [1]–[4], the authors studied UAVs' transmission technology and discussed various open research problems of UAVs-aided networks. The authors in [5]–[10] optimized UAVs' aerial positions to assist the ground static networks by deploying UAVs as the moving relays. In addition, many recent works have deployed UAVs as aerial mobile BSs. The works [10]–[14] analyzed the UAV-based aerial BS area coverage problem by optimizing the UAV placement. In [15], [16], the authors proposed a novel framework by jointly optimizing the UAV trajectory and radio resources to achieve UAV-aided network throughput maximization and secure communication. However, the scenario with multiple UAVs and multiple users is of more practical and appealing in the content-centric network design. In [17], the authors considered the multi-UAV trajectory and deployment optimization problem, and both the cases of static and mobile ground terminals were studied. The authors in [18] exploited multiple UAVs to form a C-RAN to assist ground BSs and designed the UAVs' travel path to optimize the spectral efficiency. In [19], the authors expanded their previous work to the multi-UAV case by jointly optimizing UAVs' trajectory, power and communication scheduling. However, these multi-UAV control framework are not cross-layer design, in the sense that only physical layer metrics are considered. In addition, the computation complexity of the multi-UAV optimization algorithms used in above works are very high. Thus the implementation of such algorithm is quite challenging for UAV online control.

The deployment of UAVs in wireless networks may induce a tremendous backhaul burden, and the BS to UAV link may be the performance bottleneck. Cache-enabled UAVs

Manuscript received October 28, 2020; revised February 24, 2021; accepted April 12, 2021. Date of publication June 14, 2021; date of current version September 16, 2021. This work was supported by the Research Grants Council, Hong Kong, under Project 16204018. (Corresponding author: Shuqi Chai.)

The authors are with the Department of Electronic and Computer Engineering, The Hong Kong University of Science and Technology, Hong Kong (e-mail: schai@ust.hk; eeknlau@ust.hk).

Color versions of one or more figures in this article are available at <https://doi.org/10.1109/JSAC.2021.3088694>.

Digital Object Identifier 10.1109/JSAC.2021.3088694

content-centric networking [20]–[22]<sup>1</sup> has thus emerged as a promising solution to alleviate the bottlenecks. In such networking, UAVs are proactively loaded with cached content and deployed to serve the users according to the instantaneous user demands, with the transmitted packets coming from the local cache instead of the BS-UAV link. There are a few research works on UAV caching communications. In [23], a deep learning-based UAV caching design was proposed. However, only optimized UAV transmission and cache content scheduling was considered, but not the trajectory design. In [24], the authors studied UAV-assisted caching in small-cell networks, while the authors in [21] proposed a UAV-based content caching and position design scheme. However, the above multi-UAV control designs suffer from exponential complexity w.r.t. the number of UAVs. Furthermore, in these works, the solutions have ignored the limited battery life as well as the limited storage capacity of the UAV. In [25], the authors proposed a cache-enabled UAV transmission and trajectory scheme taking into account the UAV endurance issue. In our recent paper [26], we proposed an online reinforcement learning algorithm to control UAV trajectory and radio resource allocation for UAV-assisted wireless networks. The solution considers a novel energy and content recharging scheme, which can effectively address the finite battery life and finite storage capacity issue of the UAV. However, only the single UAV case is considered in [25] and [26].

In this paper, we extend the work in [26] to a multi-UAV system. Specifically, we consider demand-driven dynamic control of the trajectories, the decision to return to the power and content charging station or not, and transmission power control of the multi-UAV system subject to the constraints of limited battery and limited cache storage capacity. The extension to multiple UAVs is highly non-trivial due to the strong coupling among the states of the UAVs induced by the anti-collision constraint. The following summarizes the key challenges and our contributions.

- **Multi-UAV Trajectory Optimization with Content and Energy Charging:** In most existing works, the finite battery life of the UAV as well as the finite storage have been ignored. In this paper, we propose a novel energy and content recharging scheme which enables continuous operation of the cache-assisted UAV network as well as adaptation to the drift of content popularity in the requests distribution of the user. As such, the decision space of the UAV system has to include decisions as to whether a particular UAV has to continue to serve the next demands or to fly back to the charging station to replenish the energy and cache content in addition to the trajectory control and transmission power control. The proposed solution is demand-driven in the sense that it is adaptive to both the channel state information (CSI) between the UAVs and the users as well as the request queue state information (QSI) of the demands from various users.
- **Complex Dynamic State Coupling in the Multi-UAV System:** One of the key challenges in the extension to

a multi-UAV system is the complexity issue induced by complex dynamic state coupling. Specifically, the states of the multiple UAVs are tightly coupled together due to the *inter-UAV collision avoidance constraint*. The inter-UAV collision avoidance constraint is to prevent UAVs from crashing into each other due to the wind disturbance or mis-coordinated trajectory control. Such complex state coupling induce exponential complexity in the brute-force optimization approach.

- **Reduced-Complexity Mean-Field Game (MFG) Solution:** To reduce the complex coupling issues, we first consider an ergodic stochastic optimization formulation with the state evolution of the multi-UAV system modeled as an  $N$ -dimensional continuous time stochastic differential equations (SDEs) [27], [28]. To address the inter-UAV collision issue, we define a new anti-collision loss function in the optimization objective that penalizes the collision avoidance violation. To deal with the coupling induced by the inter-collision, we exploit mean-field analysis [28]–[30] to construct a *mean-field equilibrium* which captures a limiting behavior with the number of UAVs tending to infinity. Using *propagation of chaos* in statistical mechanics, we show that the original ergodic control problem for the multi-UAV is asymptotically equivalent to a reduced-complexity mean-field game for a large number of UAVs. Such mean-field approximation enables decomposition of the original coupled HJB equations into the MFG equations.
- **Deep-learning Design and Online Training:** Solving the derived MFG equations using classical approach is still computational intractable. There are several research works on the learning based UAV trajectory and transmission control. In [31], the authors design a CNN based deep supervised learning to enable UAV-aided edge caching. In [32], the authors propose a hierarchical deep reinforcement learning based UAV trajectory and resource control. In contrast, we exploit problem-specific structures to design an online deep learning solution to learn the solution. The motivation of using the DNN method is on the fast inferencing to compute the control policy given a state realization. To enhance the expressibility of the DNN, we tailored the datapath of the DNN according to the structural properties of the optimal control and the value functions based on the *homotopy perturbation method* [33], [34]. As such, the DNN can approximate the value function more efficiently with less weights. Based on this, we develop an *online training algorithm* which does not require *labeled data*. In other words, the online training algorithm autonomously adjust the weights of the DNN to converge to the best approximator of the solution of the MFG equations [35]–[37]. We have also obtained the error bound of the DNN and showed that the complexity is independent of the number of UAVs.

The paper is organized as follows. In Section II, we introduce the architecture of the UAV-assisted wireless caching system and communication model, as well as the UAV control and energy consumption model. In Section III, we formulate the UAV trajectory and radio resource control

<sup>1</sup>Over 70% of the new capacity demands for 5G wireless networks come from high quality video streaming applications. As such, these capacity demands are mostly content-centric. These demands have huge spatial and temporal correlations between content requests by users in the network.

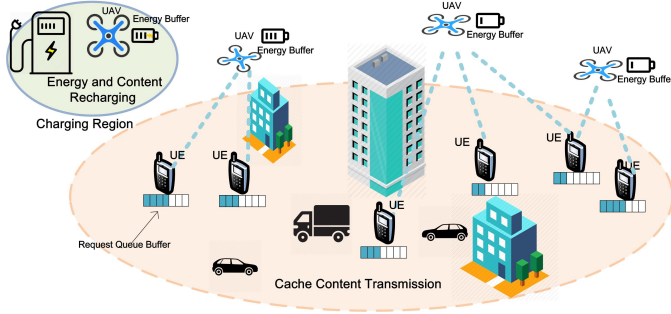


Fig. 1. Architecture of the multi-UAV wireless caching network with energy and content recharging.

problem as a stochastic differential game. In Section IV, using the mean-field game approach, we derive a reduced-complexity solution based on mean-field equilibrium analysis. In Section V, we propose a novel model-specific deep neural network (DNN) to learn the partial differential equations (PDEs) numerically by exploiting an HPM method. Section VI illustrates the numerical simulation results and presents discussion. Finally, we summarize our main results in Section VII.

## II. SYSTEM MODEL

In this section, we introduce the architecture of the multi-UAV wireless caching network, with the basic stochastic dynamic models of the UAV trajectory, user content requests queues and UAV energy/content recharging discussed in detail.

### A. UAV-Assisted Wireless Caching Network Architecture

As illustrated in Fig.1, we consider a cache-enabled wireless network where  $N$  UAVs are deployed to provide data transmission to  $M$  users by using unicast. Assume that there are a total of  $NJ$  content files in the content server, and each UAV will be responsible for the content files set  $|\mathcal{J}^n| = J$  independently. Denote the multi-UAV service groups as  $\mathcal{N} = \{1, \dots, N\}$ , and each UAV  $n$  can only be cached with a single content file  $j \in \mathcal{J}^n$  to serve the users and is subject to a finite energy and storage buffer<sup>2</sup>. Specifically, let  $\mathcal{M} = \{1, \dots, M\}$  denote the set of users served by all UAVs where  $M$  denotes the number of users. To support continuous operation of the UAV network, we assume there is a charging station at position  $\mathbf{X}^{ch} = [x^{ch}, y^{ch}, z^{ch}]^T$ . The charging station can support energy and content recharging for UAVs.

### B. UAV Communication Model

The 3D coordinates of the  $n$ -th UAV are given by  $\mathbf{X}_n = [x_n, y_n, z_n]^T$ . Let  $\bar{\mathbf{X}}_m = [\bar{x}_m, \bar{y}_m, \bar{z}_m]^T$  denote the position of user  $m$ . The distance between the UAV  $n$  and user  $m$  can be written as

$$D_{nm} = \|\mathbf{X}_n - \bar{\mathbf{X}}_m\|_2. \quad (1)$$

Based on this, the large-scale path gain from UAV  $n$  to user  $m$  at time  $t$  can be modeled as<sup>3</sup>

$$L_{nm}(t) = \rho_0 D_{nm}^{-2}(t) = \frac{\rho_0}{\|\mathbf{X}_n(t) - \bar{\mathbf{X}}_m\|_2^2}, \quad (2)$$

<sup>2</sup>The solution in this paper can be extended easily to include caching multiple files. In Sec VI, we have compare the performance with the extension.

<sup>3</sup>We assume that all users positions are static for simplicity in this paper.

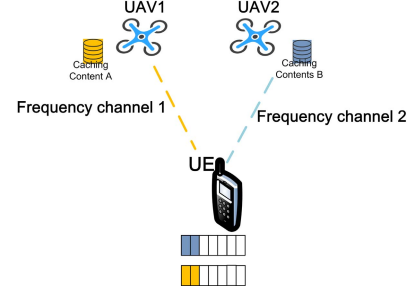


Fig. 2. Illustration of UAV transmission model with the case of 2 UAV and single user.

where  $\rho_0$  denotes the reference path gain at distance  $d_0 = 1\text{m}$ .

The instantaneous channel power gain between UAV  $n$  and UE  $m$  can be modeled as

$$|H_{nm}(t)|^2 = L_{nm}(t)|h_{nm}(t)|^2, \quad (3)$$

where  $L_{nm}(t)$  is the large-scale path gain as defined in (2), and  $h_{nm}(t)$  accounts for the short-term channel with dynamics given by the following stochastic differential equation:

$$dh_{nm}(t) = -\frac{1}{2}a_H h_{nm}(t)dt + \sqrt{a_H}dW_{nm}^h(t), \quad (4)$$

where  $a_H > 0$  determines the temporal correlation (depending on the mobility of the user) and  $W_{nm}^h(t)$  is the standard independent Wiener process for the short-term fading process [38]. According to Ito's formula, the channel power gain process  $G_{nm}(t) = |h_{nm}|^2$  will have the dynamic

$$dG_{nm}(t) = (a_H - a_H G_{nm}(t))dt + 2\sqrt{a_H G_{nm}(t)}dW_{nm}^h(t), \quad (5)$$

We assume each UAV has a dedicated channel and hence, can serve one user at a time. However, the users are equipped with multiple radios and they are allowed to receive data streams from different UAVs simultaneously. Figure 2 illustrates an example.

As each UAV has a dedicated channel to serve users, and thus there's no frequency interference over different UAV-users links. As a result, the received signal at the  $m$ -th user from the  $n$ -th UAV can be expressed as  $Y_{nm} = \sqrt{P_n L_{nm}(t)}h_{nm}(t)X_n + Z_{nm}$ , where  $X_n$  denotes the transmitted symbol, and  $P_n$  is the transmit power of the  $n$ -th UAV, and  $Z_{nm} \sim \mathcal{CN}(0, 1)$  is the channel noise from UAV  $n$  to user  $m$ . Let  $u_{nm} \in \{0, 1\}$  denote the user selection index for the  $m$ -th user at the  $n$ -th UAV with the scheduling constraint  $\sum_{m=1}^M u_{nm} = 1$ . The achievable data rate (bits/sec) from the  $n$ -th UAV to the  $m$ -th user at time  $t$  is given by

$$R_{nm}(t) = u_{nm}(t)W \log_2(1 + |H_{nm}(t)|^2 P_n(t)), \quad (6)$$

where  $W$  is the bandwidth allocated to the UAV.

We assume that the  $n$ -th UAV can only be preloaded with one content file  $j \in \mathcal{J}^n$  in the cache. Denote the caching state as  $\mathbf{s}_n = \{s_n^j, \forall j \in \mathcal{J}^n\}_{n \in \mathcal{N}}$ , where  $s_n^j \in \{0, 1\}$  is the  $j$ -th content file caching placement state at UAV  $n$ . Specifically,  $s_n^j = 1$  if the  $j$ -th content file is cached in UAV  $n$ , and  $s_n^j = 0$ , otherwise. The cache states at each UAV should satisfy  $\sum_{j=1}^J s_n^j = 1, \forall n \in \mathcal{N}$ .



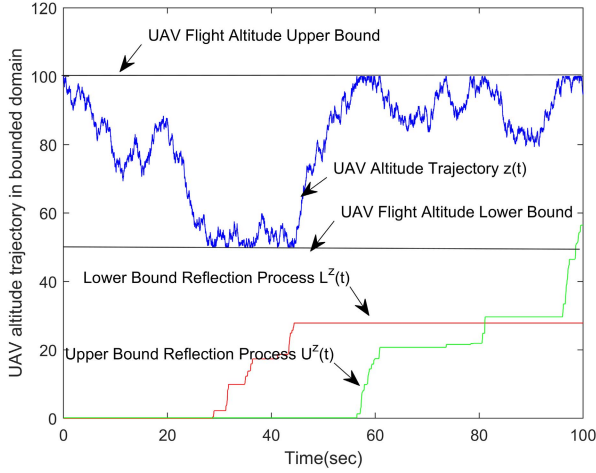


Fig. 3. Example of UAV altitude trajectory in bounded region and reflection processes with upper bound at 100m, lower bound at 50 m.

### C. UAV Trajectory Control Model

Given the UAV flight trajectory control  $\mathbf{v}_n = [\mathbf{v}_n^x, \mathbf{v}_n^y, \mathbf{v}_n^z]^T$ , the  $n$ th UAV trajectory follows the dynamic:

$$d\mathbf{X}_n(t) = \mathbf{v}_n(t)dt + \boldsymbol{\sigma}_X d\mathbf{W}_n^X(t) + d\mathbf{L}_n^X(t) - d\mathbf{U}_n^X(t), \quad (7)$$

where  $\boldsymbol{\sigma}_X = \text{diag}(\sigma_X^x, \sigma_X^y, \sigma_X^z)$  denotes the volatility of the process due to the wind disturbance, and  $\mathbf{W}_n^X(t)$  is a three dimensional independent standard Brownian motion process.  $\mathbf{L}_n^X(t) = [L_n^x(t), L_n^y(t), L_n^z(t)]^T$  and  $\mathbf{U}_n^X(t) = [U_n^x(t), U_n^y(t), U_n^z(t)]^T$  are the non-decreasing and continuous reflection processes associated with the UAVs flight space boundaries  $x(t) = L_x$ ,  $x(t) = U_x$ ,  $y(t) = L_y$ ,  $y(t) = U_y$  and  $z(t) = L_z$ ,  $z(t) = U_z$  respectively, where  $L_x, L_y, L_z$  and  $U_x, U_y, U_z$  are the boundary value. As a result,  $\mathbf{X}_n(t)$  will always be within  $[L_x, U_x] \times [L_y, U_y] \times [L_z, U_z]$ . The reflection process associated with the UAV trajectory can be uniquely determined by the following equations:

$$\begin{aligned} \int_0^t 1_{\{x_n(t) > L_x\}} dL_n^x(t) &= \int_0^t 1_{\{x_n(t) < U_x\}} dU_n^x(t) \\ &= \int_0^t 1_{\{y_n(t) > L_y\}} dL_n^y(t) = \int_0^t 1_{\{y_n(t) < U_y\}} dU_n^y(t) \\ &= \int_0^t 1_{\{z_n(t) > L_z\}} dL_n^z(t) = \int_0^t 1_{\{z_n(t) < U_z\}} dU_n^z(t) = 0 \end{aligned} \quad (8)$$

Figure 3 illustrates an example of the UAV altitude trajectory with reflective boundaries  $L_z = 60\text{m}$  and  $U_z = 100\text{m}$ .

### D. Content Request Queue Model

All users in the service group will send content request to each UAV for the specific content file, and we adopt the *request queue model* (in bits) to characterize the requests of the users. Assume that each content has the same file size. Let  $Q_{nm}^j(t)$  denote the number of bits of the  $j$ -th content ( $j \in \mathcal{J}^n$ ) requested by the  $m$ -th user in the  $n$ -th UAV.  $Q_{nm}^j(t)$  evolves according to the dynamic equation [22]:

$$dQ_{nm}^j(t) = [A_{nm}^j(t) - s_n^j(t)R_{nm}(t)]dt + dL^{Q_{nm}^j}(t), \quad (9)$$

TABLE I  
UAV ENERGY PROPULSION CONSUMPTION PARAMETERS

Parameter	Physical Meaning	Parameter	Physical Meaning
$\delta_e$	Blade drag coefficient 0.012	$s$	Rotor solidity 0.05
$\Omega_e$	Blade angular velocity 400 radians/sec	$g$	Gravity acceleration 9.8
$R_e$	Rotor radius 0.38	$A$	Rotor disc area 0.79
$\rho$	Air density 1.225	$M_{\text{UAV}}$	UAV mass 4 kg
$\kappa_p$	Induced power factor 0.1	$\mathbf{v}$	UAV velocity

where  $A_{nm}^j(t)$  is the request arrival rate (bits/sec) with  $\mathbb{E}[A_{nm}^j] = \bar{A}$ ,  $\forall n, m, j$ .  $R_{nm}(t)$  is the achievable data rate (bits/sec) of the  $n$ -th UAV to the user  $m$  as defined in (6).  $L^{Q_{nm}^j}(t)$  denotes the reflection process associated with the lower request queue boundary  $Q_{nm}^j(t) = 0$  such that the request queue  $Q_{nm}^j(t)$  will not go below zero. The reflection process  $L^{Q_{nm}^j}(t)$  can be uniquely determined by

$$\int_0^t 1_{\{Q_{nm}^j(t) < 0\}} dL^{Q_{nm}^j}(t) = 0. \quad (10)$$

Note that the request queue  $Q_{nm}^j(t)$  is a virtual queue and hence, there is no upper bound on  $Q_{nm}^j(t)$ .

### E. UAV Energy Consumption and Recharging Model

In this paper, all UAVs will be able to recharge the energy and cache content when it hovers over the charging station within the distance of  $D_0$ . Let  $E_n(t)$  denote the remaining energy (in Joule) in the battery of the  $n$ -th UAV. It follows a dynamic equation given by:

$$dE_n(t) = (-P_n(t) - E_n^p(t) + E_n^{ch}(t))dt + dL_n^E(t) - dU_n^E(t). \quad (11)$$

$E_n^p(t) = (c_1 + c_2 \|\mathbf{v}_n(t)\|_2^2)$  denotes the UAV propulsion consumption rate (Joule/sec) which depends on the UAV velocity with two terms: profile power and induced power [39]–[42], where  $c_1 = \frac{\delta_e}{8} \rho s A \Omega_e^3 R_e^3 + \kappa_p \frac{(M_{\text{UAV}} g)^{3/2}}{\sqrt{2\rho A}}$  and  $c_2 = \frac{3\delta_e}{8\Omega_e^2 R_e^2} \rho s A \Omega_e^3 R_e^3$ . The detailed parameters of the UAV propulsion consumption are listed in Table I.

Given a certain energy level  $E_n(t)$  in the battery, the  $n$ -th UAV has an option to continue its journey to serve user requests or to return to the charging station to replenish the battery. We assume there is an efficient power transfer mechanism from the charging station to the UAV such as wireless power transfer [43], [44] or laser power transfer [45]. However, the detailed mechanism is outside the scope of this paper. The UAV energy charging rate in (11) is given by:

$$E_n^{ch}(t) = a \mathbf{1}(D_n^{ch}(t) \leq D_0), \quad (12)$$

where  $a$  is the energy transfer rate (Joule/sec) from the charging station to the UAV, and  $D_n^{ch}(t)$  denotes the distance from UAV  $n$  to the charging region. In addition, we assume that the energy charging rate is larger than the consumption rate in the charging region:  $a \geq c_1 + c_2 v_c^2$ .

$L_n^E(t)$  and  $U_n^E(t)$  in (11) denote the reflection processes associated with the lower and upper energy queue boundary  $E_n(t) = 0$  and  $E_n(t) = N_E$  such that when the UAV energy buffer is empty, no energy can be consumed and when the UAV energy buffer is full, additional energy cannot be replenished. The reflection processes  $L_n^E(t)$  and  $U_n^E(t)$  together with the admissible control constraint on the available energy defined in Section III ensures the energy queue length  $E_n(t)$  lies in the domain  $[0, N_E]$ . The energy queue reflection process can be uniquely determined by

$$\int_0^t 1_{\{E_n(t) < 0\}} dL_n^E(t) = 0. \quad (13)$$

$$\int_0^t 1_{\{E_n(t) > N_E\}} dU_n^E(t) = 0. \quad (14)$$

#### F. UAV Cache Content Recharging Model

Beside the energy charging, the charging station can also update the content in the UAV cache buffer. When the  $n$ -th UAV gets to the charging region, the cache content states are determined by the following equation:

$$s_n^j(t) = \begin{cases} 1, & \text{if } j = \arg \max \sum_m Q_{nm}^j(t), \\ 0, & \text{otherwise} \end{cases}, \forall j \in \mathcal{J}^n.$$

Specifically, UAVs will be charged with the most popular content file into the caching buffer. Instantaneous popularity at time  $t$  can be characterized by the instantaneous request queue length  $\sum_m Q_{nm}^j(t)$ . The content with the longest request queue  $j^* = \arg \max \sum_m Q_{nm}^j(t)$  will be replenishment at the  $n$ -th UAV buffer.

*Remark 1:* We assume that the content recharging can be completed within a limited time using mmWave beamforming when the UAV is in close proximity with the charging station. In practice, the energy charging will take longer time compared with content charging. Furthermore, the UAV will always be triggered by insufficient energy to return to the charging station. Hence, the content charging will always be completed well before the energy charging.

### III. UAV TRAJECTORY AND TRANSMISSION CONTROL DYNAMIC OPTIMIZATION

In this section, we shall design a dynamic trajectory and radio resource control policy for the multi-UAV system by considering an infinite horizon stochastic differential game. The dynamic policy is adaptive to the instantaneous request queue length (QSI)  $Q_n(t)$ , the instantaneous energy queue length (ESI)  $E_n(t)$ , the instantaneous channel condition (CSI)  $G_n(t)$  as well as the current location of the UAV  $\mathbf{X}_n(t)$ . Adaptation to the CSI reveals **good transmission opportunities** induced by time-varying channel fading between the UAV and the user. Adaptation to the QSI reveals the **dynamic urgency of individual content flows**. Adaptation to the ESI reveals the **dynamic urgency induced by the battery life**. We first define the admissible control policy which considers multi-UAVs' safe operation conditions such as the finite energy constraint. We then formulate the problem as an  $N$ -player stochastic differential game and derive the associated optimality condition (Nash equilibrium).

#### A. UAV Trajectory and Radio Resource Control Admissible Policy

Given the system model in Section III, we now introduce the UAV control policy for content-centric caching network. Let  $\mathbf{S}_n = \{\mathbf{X}_n, \mathbf{Q}_{nm}, E_n, G_{nm}\} \in \mathcal{S}_n$  denotes the local system state at the  $n$ -th UAV and  $\mathbf{S}^N = \{\mathbf{S}_1, \dots, \mathbf{S}_N\} \in \mathcal{S}^N := \mathcal{S}_1 \times \dots \times \mathcal{S}_N$  denotes the global system state. Based on the system states, each UAV  $n$  implements an admissible control policy to determine the transmission power  $P_n^N$ , the user to serve  $\{u_{nm}^N\}_{m \in \mathcal{M}}$  and the trajectory control  $\mathbf{v}_n^N = [v_n^{Nx}, v_n^{Ny}, v_n^{Nz}]^T$ . We define the stationary admissible control policy  $\Omega^N$  as below.

*Definition 1 Admissible UAV Control Policy:* Let  $\mathcal{U}$  denote the set of admissible control policies. The admissible control policy of the  $n$ -th UAV  $\Omega_n^N \in \mathcal{U}$  is a mapping from the system state  $\mathbf{S}^N(t)$  to the trajectory and transmission actions  $\{P_n^N, \{u_{nm}^N\}_{m \in \mathcal{M}}, \mathbf{v}_n^N\}$ . In addition, the actions should satisfy the following constraints:

- $\|\mathbf{v}_n^N(t)\|_2 \leq v_c$ , where  $v_c$  denotes the UAV maximum speed in m/s;
- Under the admissible policy, the energy state and the control of the  $n$ -th UAV should satisfy

$$\begin{aligned} E_n(s) - \int_s^t (P_n(\tau) + c_1 + c_2 \|\mathbf{v}_n(\tau)\|_2^2 - E_n^{ch}(\tau)) d\tau - U_n^E(t) \\ \geq \frac{(c_1 + c_2 v_c^2)}{v_c} \left( \|\mathbf{X}(s) + \int_s^t \mathbf{v}_n(\tau) d\tau + \int_s^t \boldsymbol{\sigma}_X d\mathbf{W}_n^X(\tau) \right. \\ \left. + L_n^X(t) - U_n^X(t) - \mathbf{X}^{ch} \right\|_2 - D_0 \end{aligned} \quad (15)$$

- The control policy is a stationary unichain policy. The request queuing system  $Q_{nm}^j(t), \forall n \in \mathcal{N}, m \in \mathcal{M}, j \in \mathcal{J}^n$  under the policies  $\Omega_1^N \times \dots \times \Omega_N^N$  is stable with  $\lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T \sum_{n,m,j} \mathbb{E}^{\Omega^N} [Q_{nm}^j(t)] dt < \infty$ .

The first condition is due to the constraint on the maximum speed of a UAV. The second condition is the safety constraint, which ensures that each UAV must have enough energy to get to the charging station with the highest speed. In the third condition, we restrict our attention to stationary unichain policies, i.e., the controlled stochastic process under  $\Omega$  has a single recurrent class (and possibly some transient states) [27]. Note that such assumptions are reasonable as the optimal policy is also stationary<sup>4</sup>.

#### B. Infinite Horizon Ergodic Differential Game Formulation

We formulate the problem as an infinite horizon ergodic stochastic differential game, where the UAV control policies  $\{\Omega_1^N, \dots, \Omega_N^N\}$  are jointly designed to minimize its average power consumption cost, its average collision cost and the average request delay cost. Specifically, each UAV is a player targeted to optimize the following ergodic cost:

<sup>4</sup>In this paper, we focus on the set of the admissible control policies which give irreducible and ergodic markov process with well-defined steady state distribution.

**Problem 1: Ergodic Stochastic Differential Game for multi-UAV Trajectory and Transmission Control**

$$\min_{\Omega_n^N \in \mathcal{U}} J^n(\Omega_n^N, \Omega_{-n}^N) \quad (16)$$

$$\begin{aligned} &= \limsup_{T \rightarrow \infty} \mathbb{E}^{\Omega_n^N} \left[ \frac{1}{T} \left[ \int_0^T [c_n(\mathbf{S}_n(t), \mathbf{S}_{-n}(t), \Omega_n^N)] dt \right. \right. \\ &\quad \left. \left. + \varrho_1(L_n^x(T) + L_n^y(T) + L_n^z(T)) \right. \right. \\ &\quad \left. \left. + \varrho_2(\mathbf{U}_n^x(T) + \mathbf{U}_n^y(T) + \mathbf{U}_n^z(T)) \right] + \rho_3 L_n^E(T) \right] \\ \text{s.t. } &d\mathbf{S}_n(t) = F_n(\mathbf{S}_n(t), \mathbf{S}_{-n}(t), \Omega_n^N) dt \\ &+ D_n(\mathbf{S}_n(t)) d\mathbf{W}_n(t) + d\mathbf{L}_n(t) - d\mathbf{U}_n(t). \end{aligned} \quad (17)$$

The per-stage cost function  $c_n(\mathbf{S}_n, \mathbf{S}_{-n}, \Omega_n^N)$  is defined as:  $c_n(\mathbf{S}_n, \mathbf{S}_{-n}, \Omega_n^N) = P_n^N + \alpha c_n^{\text{delay}}(\mathbf{S}_n) + \frac{1}{N} \sum_{k \neq n} \frac{\beta}{\gamma + \|\mathbf{X}_n - \mathbf{X}_k\|_2^2}$ ,<sup>5</sup> where the first term in the cost function denotes the transmission power cost, and the second term is the users request delay cost  $c_n^{\text{delay}}(\mathbf{S}_n) = \sum_{m \in \mathcal{M}, j \in \mathcal{J}_n} Q_{nm}^j$  with cost weighting  $\alpha > 0$ . The third term refers to the normalized UAV anti-collision cost and aims to penalize the UAVs collision events<sup>6</sup> with cost weighting  $\beta > 0$  and  $\gamma > 0$ .  $\mathbf{S}_n(t) = [E_n(t), Q_{nm}^j(t), \mathbf{X}_n(t), G_{nm}(t)]_{m \in \mathcal{M}, j \in \mathcal{J}_n}^T$  is the local state of the  $n$ -th UAV and  $\mathbf{S}_{-n}(t)$  is the states of all other UAVs. Note that the last two terms in the ergodic control cost (16) refer to the penalty associated with the UAVs displacement at the lower and upper boundary of the 3D flight region, where  $\varrho_1$  and  $\varrho_2$  are the penalty rates. From (7), (9) and (11), the dynamic drift can be written as  $F_n(\mathbf{S}_n, \mathbf{S}_{-n}, \Omega_n^N) = [(-P_n^N - E_n^p + E_n^{\text{ch}}), A_{nm}^j - s_n^j R_{nm}(\Omega_n^N), \mathbf{v}_n^x, \mathbf{v}_n^y, \mathbf{v}_n^z, -\frac{1}{2} a_H G_{nm}]_{m \in \mathcal{M}, j \in \mathcal{J}_n}^T$ , and the diffusion term  $D_n = \text{diag}(\mathbf{0}_{1+MJ}, \sigma_X^x, \sigma_X^y, \sigma_X^z, 2\sqrt{a_H G_{n1}}, \dots, 2\sqrt{a_H G_{nM}})$ . The reflection process  $\mathbf{L}_n(t)$  and  $\mathbf{U}_n(t)$  of the  $n$ -th UAV are defined in (8)-(14).

**C. Optimality Condition for Multi-UAV Control Problem**

The optimal solution for the stochastic differential game (SDG) in (16)-(17) is defined as the Nash equilibrium (NE):

**Definition 2 (Nash Equilibrium of the Stochastic Differential Game):** The set of controls  $\Omega_n^{N*}(\mathbf{S}^N) = \{P_n^{N*}, \{u_{nm}^{N*}\}_{m \in \mathcal{M}}, \mathbf{v}_n^{N*}\}$ ,  $\forall n \in \mathcal{N}$  constitute the Nash equilibrium for the  $N$ -UAVs SDG Problem (16)-(17) if

$$J^n(\Omega_n^{N*}, \Omega_{-n}^{N*}) \leq J^n(\Omega_n^{N'}, \Omega_{-n}^{N*}), \forall n \in \mathcal{N}, \quad (18)$$

for any other alternative control  $\{\Omega_1^{N'}, \dots, \Omega_N^{N'}\}$ .

Given the problem formulation in (16), the NE policy  $\Omega^{N*}$  defined in (18) can be obtained by solving the following sufficient condition for NE:

<sup>5</sup>In this paper, the per-stage transmission power cost accounts for the penalty of using a large transmission power and can be denoted as the soft power constraint. The weight  $\alpha$  and  $\beta$  can be interpreted as the corresponding Lagrange Multipliers which determine the Pareto trade off between the power, delay and the UAV collision cost. The precise operating point on the Pareto boundary is determined by the application scenario and is out of the scope of the paper.

<sup>6</sup>In practice, there are local anti-collision mechanisms in each UAV and the onboard controller will avoid the crashing of the UAV. However, the cost term is to penalize the triggering of such on-board anti-collision events.

**Theorem 1 Sufficient Condition for the NE :** Assume there exists functions  $V_n^N(\mathbf{S}_n) \in \mathcal{C}^2(\mathbf{S}_n)$  and  $\theta_n^N$ ,  $\forall n$  that solves the following average cost HJB optimality equation:

$$\begin{aligned} \min_{\Omega_n^N \in \mathcal{U}} &\left[ P_n^N + \sum_{m=1}^M \sum_{j=1}^J \alpha Q_{nm}^j + \frac{1}{N} \sum_{k \neq n} \frac{\beta}{\gamma + \|\mathbf{X}_n - \mathbf{X}_k\|_2^2} \right. \\ &+ \sum_{n=1}^N F_n(\mathbf{S}_n, \mathbf{S}_{-n}, \Omega_n^N)^T \cdot \nabla_{\mathbf{S}_n} V_n^N(\mathbf{S}_n) \\ &\left. + \sum_{n=1}^N \frac{1}{2} \text{tr} \left( D_n D_n^T \nabla_{\mathbf{S}_n}^2 V_n^N(\mathbf{S}_n) \right) \right] = \theta_n^N, \forall n \in \mathcal{N} \quad (19) \end{aligned}$$

where  $\nabla_{\mathbf{S}_n} V_n^N(\mathbf{S}_n)$  is the gradient of  $V_n^N(\mathbf{S}_n)$ , and  $\nabla_{\mathbf{S}_n}^2 V_n^N(\mathbf{S}_n)$  is the Hessian of  $V_n^N(\mathbf{S}_n)$ . Then we have the following results:

1.  $\theta_n^N$  is the optimal cost for UAV  $n$  in Problem (16), and  $V_n^N(\mathbf{S}_n)$  is the optimal value function, with the boundary conditions

$$\nabla_{x_n} V_n^N(\mathbf{S}_n)|_{x_n=L_x} = \nabla_{y_n} V_n^N(\mathbf{S}_n)|_{y_n=L_y} \quad (20)$$

$$= \nabla_{z_n} V_n^N(\mathbf{S}_n)|_{z_n=L_z} = -\varrho_1;$$

$$\nabla_{x_n} V_n^N(\mathbf{S}_n)|_{x_n=U_x} = \nabla_{y_n} V_n^N(\mathbf{S}_n)|_{y_n=U_y} \quad (21)$$

$$= \nabla_{z_n} V_n^N(\mathbf{S}_n)|_{z_n=U_z} = \varrho_2;$$

$$\begin{aligned} &\nabla_{Q_{nm}^j} V_n^N(\mathbf{S}_n)|_{Q_{nm}^j=0} \\ &= 0, \forall j \in \mathcal{J}, m \in \mathcal{M}, \end{aligned} \quad (22)$$

$$\nabla_{E_n} V_n^N(\mathbf{S}_n)|_{E_n=N_E} = 0; \quad (23)$$

$$\nabla_{E_n} V_n^N(\mathbf{S}_n)|_{E_n=0} = -\rho_3; \quad (24)$$

2. If admissible control  $\{P_n^{N*}, \{u_{nm}^{N*}\}_{m \in \mathcal{M}}, \mathbf{v}_n^{N*}\}$  attains the minimum of the L.H.S. of (19), then  $\Omega_n^{N*}(\mathbf{S}) = \{P_n^{N*}, \{u_{nm}^{N*}\}_{m \in \mathcal{M}}, \mathbf{v}_n^{N*}\}$  is the optimal control policy for Problem (16).

*Proof:* Please see Appendix A. ■

Based on this, each UAV player  $n$  will focus on optimizing its long-term cost  $J^n(\Omega_n^N, \Omega_{-n}^N)$  in a selfish manner and adapt its trajectory and power control to the system state. However, the problem is rather challenging as each UAV player should have access to the knowledge of its own state  $\mathbf{S}_n(t)$  and other UAV players state feedback  $\mathbf{S}_{-n}(t)$ . Furthermore, in order to build the Nash equilibrium for this  $N$ -UAV SDG, the system of  $N$  coupled HJB equations in (19) should be jointly solved. When  $N \geq 2$ , the computational complexity grows exponentially.

**Challenge 1:** Solving the  $N$ -tuple of HJB equations (19) either analytically or numerically is highly non-trivial due to the huge dimension of the state and mutual coupling via the anti-collision cost.

**Challenge 2 (Centralized Solution):** Brute force optimization of (19) results in a centralized control policy which requires global real-time state observations  $\{\mathbf{S}_1, \dots, \mathbf{S}_N\}$  and they are very difficult to obtain.



#### IV. DECOUPLING BY MEAN-FIELD ANALYSIS

In this section, we shall overcome the aforementioned challenges using mean-field limit analysis [28]–[30]. Specifically, the coupling in the states can be fully characterized by a mean-field equilibrium, which is specified by a pair of backward Hamiltonian-Jacobi-Bellman (HJB) and forward Fokker-Planck-Kolmogorov (FPK) equations. As such, conditioned on the mean-field equilibrium, the  $N$ -dim HJB equations in (19) can be decoupled and this can substantially reduce the complexity. Furthermore, the stationary policies in (19) can be reduced to decentralized policies in which the actions of the  $n$ -th UAV will be adaptive to the local states  $S_n$  only. We will then solve the mean-field equilibrium equation and the HJB equation by deriving a reduced-complexity optimality equation and analyze the structural properties of the mean-field optimal solution.

##### A. Mean-Field Analysis

We first define the following indistinguishable property to characterize the individual behavior in the mean-field equilibrium:

**Definition 2: Indistinguishable Property of multi-UAV SDG Problem** The multi-UAV wireless communication system states  $S^N(t) = \{S_1(t), \dots, S_N(t)\}$  are said to be indistinguishable under the policy  $\Omega^N$ , if the law of  $S^N(t)$ ,  $\mathcal{P}(\cdot)$  is invariant by any permutation of  $N$  elements:

$$\mathcal{P}(S_1(t), \dots, S_N(t)) = \mathcal{P}(S_{\pi(1)}(t), \dots, S_{\pi(N)}(t)), \quad (25)$$

where  $\pi(\cdot)$  denotes the permutation operation over  $\{1, \dots, N\}$ .

**Remark 2:** Note that to achieve the convergence of mean field equilibrium, indistinguishable condition defined in Def.2 should be satisfied, meaning that the mean field game comprises of generic players.

The system of coupled HJB equations in (19) can be rewritten in the form of:

$$\sum_{n=1}^N \mathcal{L}_n^{\Omega_n^{N*}} V_n^N(S_n) + c_n(S_n, S_{-n}, \Omega_n^{N*}) = \theta_n^N, \forall n \in \mathcal{N}, \quad (26)$$

where  $\mathcal{L}_n^{\Omega_n^{N*}}$  is defined as the  $n$ -th UAV dynamic generator with

$$\begin{aligned} \mathcal{L}_n^{\Omega_n^{N*}} V_n^N(S_n) &\stackrel{\text{def}}{=} F_n(S_n, S_{-n}, \Omega_n^{N*})^T \cdot \nabla_{S_n} V_n^N(S_n) \\ &\quad + \frac{1}{2} \text{tr} \left( D_n D_n^T \nabla_{S_n}^2 V_n^N(S_n) \right) \end{aligned} \quad (27)$$

We then have the following results on the optimal admissible policy in (26):

**Lemma 1:** The system states  $S^N(t) = \{S_1(t), \dots, S_N(t)\}$  under the optimal admissible policy in (26) are indistinguishable.

*Proof:* Please see Appendix B. ■

Based on this, each UAV dynamic  $S_n(t)$  in the multi-UAV system has a unique ergodic invariant distribution with  $\mu_n^N(S_n)$ , which satisfies the following Fokker-Planck-Kolmogorov (FPK) equation:

$$\mathcal{L}_n^{*\Omega_n^{N*}} \mu_n^N = 0, \int_{S_n} \mu_n^N(S_n) dS_n = 1, \forall n \in \mathcal{N}, \quad (28)$$

where  $\mathcal{L}_n^{*\Omega_n^{N*}}$  is the formal adjoint operator of  $\mathcal{L}_n^{\Omega_n^{N*}}$  defined in (26). Hence, all UAVs become generic and the problem P1 can be approximated by the following mean-field problem P2:

**Problem 2 Mean-field Problem:**

$$\min_{\Omega \in \mathcal{U}} J(\Omega, \mu) \quad (29)$$

$$\begin{aligned} &= \limsup_{T \rightarrow \infty} \left[ \mathbb{E}^\Omega \left[ \frac{1}{T} \left( \int_0^T c^\Omega(S(t), \mu(t)) dt \right. \right. \right. \\ &\quad \left. \left. \left. + \varrho_1 c_L(S(T)) + \varrho_2 c_U(S(T)) + \varrho_3 c_L^E(S(T)) \right) \right] \right] \\ &\text{s.t. } dS(t) = F^\Omega(S(t))dt + D(S(t))dW(t) + dL(t) \\ &\quad - dU(t), \end{aligned} \quad (30)$$

where  $J(\Omega, \mu)$  is optimized over the admissible control policy set  $\Omega \in \mathcal{U}$ .  $c^\Omega(S, \mu) = P + \alpha \sum_{m \in \mathcal{M}, j \in \mathcal{J}} Q_m^j + \int_{S'} \frac{\beta}{\gamma + \|\mathbf{x} - \mathbf{x}'\|_2^2} \mu(S') dS'$ ,  $c_L(S) = L^x(T) + L^y(T) + L^z(T)$ ,  $c_U(S) = U^x(T) + U^y(T) + U^z(T)$ ,  $c_L^E(S(T)) = L^E(T)$ , and  $F^\Omega(S) = [(-P - E^p + E^{ch}), A_m - s^j R_m, v^x, v^y, v^z, -\frac{1}{2} a_H G_m]^T_{m \in \mathcal{M}, j \in \mathcal{J}}$ , and  $D(S) = \text{diag}(\mathbf{0}_{1+M}, \sigma_x^x, \sigma_x^y, \sigma_x^z, 2\sqrt{a_H G_m}, \dots, 2\sqrt{a_H G_m})$ .  $\mu(t)$  denotes the stationary distribution of the process (30). The reflection processes  $L(t)$  and  $U(t)$  are defined similarly to (8)-(14), which satisfies  $\int_0^t 1_{\{x(t) < L_x\}} dL^x(t) = \int_0^t 1_{\{x(t) > U_x\}} dU^x(t) = \int_0^t 1_{\{y(t) < L_y\}} dL^y(t) = \int_0^t 1_{\{y(t) > U_y\}} dU^y(t) = \int_0^t 1_{\{z(t) < L_z\}} dL^z(t) = \int_0^t 1_{\{z(t) > U_z\}} dU^z(t) = 0$  and  $\int_0^t 1_{\{E(t) > N_E\}} dU^E(t) = 0$ ,  $\int_0^t 1_{\{E(t) < 0\}} dL^E(t)$ .

Note that the original state couplings  $\frac{1}{N} \sum_{k \neq n} \frac{\beta}{\gamma + \|\mathbf{x}_n - \mathbf{x}_k\|_2^2} = \frac{1}{N} \sum_{k \neq n} \int_{S'} \frac{\beta}{\gamma + \|\mathbf{x}'_n - \mathbf{x}'_k\|_2^2} \delta_{S_k} dS'$  are completely characterized by mean-field limiting term  $\int_{S'} \frac{\beta}{\gamma + \|\mathbf{x} - \mathbf{x}'\|_2^2} \mu(S') dS'$  as the empirical distribution produced by all members  $\frac{1}{N} \sum_n \delta_{S_n}(dS)$  can be approximated by single UAV player's statistical property  $\mu(S)$ . In addition, the optimal solution for the mean-field problem in (29)-(30) is defined as the mean-field equilibrium (MFE):

**Definition 3 Mean-Field Equilibrium for P2 :** The controls  $\Omega^*(S) = \{P^*, \{u_m^*\}_{m \in \mathcal{M}}, v^*\}$  achieves the mean-field equilibrium of the mean-field problem P2 if for all alternative policy  $\Omega'$

$$J(\Omega^*, \mu^*) \leq J(\Omega', \mu^*), \quad (31)$$

where  $\mu^*$  denotes the stationary distribution induced by  $\Omega^*$ . The MFE policy  $\Omega^* = \{P^*, \{u_m^*\}_{m \in \mathcal{M}}, v^*\}$  and stationary distribution defined in (31) can be obtained by solving the following sufficient condition for mean field equilibrium (HJB-FPK equations):

$$\min_{\Omega(S)} (\mathcal{L}^\Omega V(S) + c(S, \Omega(S), \mu)) = \theta \quad (32)$$

$$\mathcal{L}^{*\Omega} \mu(S) = 0, \int_S \mu(S) dS = 1, \quad (33)$$

where  $\mathcal{L}^\Omega V(S) = F^\Omega(S) \cdot \nabla_S V(S) + \frac{1}{2} \text{tr} (D D^T \nabla_S^2 V(S))$ , and  $\mathcal{L}^{*\Omega}$  is the formal adjoint operator of  $\mathcal{L}^\Omega$ .

Next, we shall establish that as  $N \rightarrow \infty$ , the original SDG problem P1 is asymptotically equivalent to the mean-field problem P2 in the following sense.

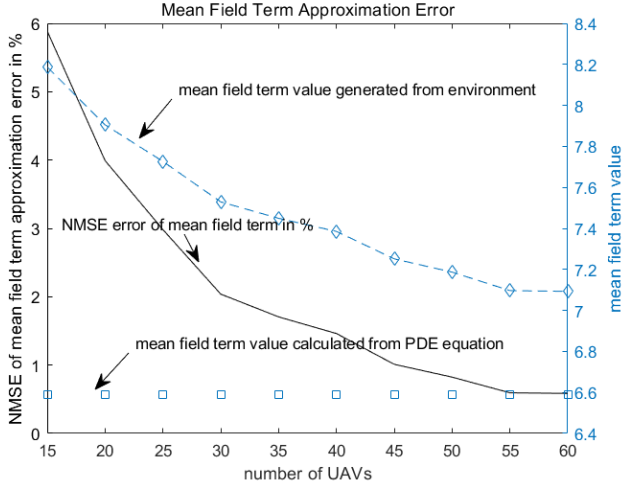


Fig. 4. Normalized mean field term approximation error with  $M = 60$ ,  $J = 10$ .

**Theorem 2 Asymptotic Equivalence of P1 and P2 :** Let  $(\theta_1^N, V_1^N), \dots, (\theta_N^N, V_N^N)$  be the solution of (19) in P1 and  $\mu_1^N, \dots, \mu_N^N$  be the stationary distributions of the states induced by the optimal policy  $\{\Omega_1^{N*}, \dots, \Omega_N^{N*}\}$  of P1. We have the following results:

- 1)  $\sup_{n,k} (|\theta_n^N - \theta_k^N| + \|V_n^N - V_k^N\|_{C^2(S)} + \|\mu_n^N - \mu_k^N\|_\infty) \rightarrow 0$  as  $N \rightarrow \infty$ ;
- 2) any limit point  $(\theta, V, \mu)$  of  $\{(\theta_1^N, V_1^N, \mu_1^N), \dots, (\theta_N^N, V_N^N, \mu_N^N)\}$  satisfies the mean-field equilibrium in (32)-(33).

**Proof:** Please see the Appendix C.

As such, we can focus on solving the mean-field problem P2, which involves one HJB equation in (32) to determine the optimal value  $\theta$  and value function  $V(\cdot)$ , and one FPK equation in (33) to determine stationary distribution  $\mu(\cdot)$ . Figure 4 illustrates the comparison of the actual mean field term  $\frac{1}{N} \sum_{k \neq n} \int_{S'} \frac{\beta}{\gamma + \|\mathbf{X}'_n - \mathbf{X}'_k\|_2} \delta_{S_k} dS'$  and the mean field term computed from (32)-(33). It can be seen that the proposed mean field approximation solution is quite accurate for number of UAV  $N \rightarrow 50$ .

### B. Structural Properties of the Mean-Field Optimal Solution

According to equation (33), we shall derive the optimal solution based on the value function and analyze some key policy structural properties for the multi-UAV mean-field control. One challenge to solve (32)-(33) is the safety constraints in the admissible action space. Specifically, inequality (15) is a compulsory state-action coupled constraint limited on the power and trajectory control, which has no closed-form solution.

**Challenge 3:** There is no closed-form solution for the action given the value function due to the complex state-action coupled constraints (15).

In the following, we shall transform the state-action constraint into a solely state-dependent *inward pointing boundary condition* for the HJB equation so as to ensure that the system

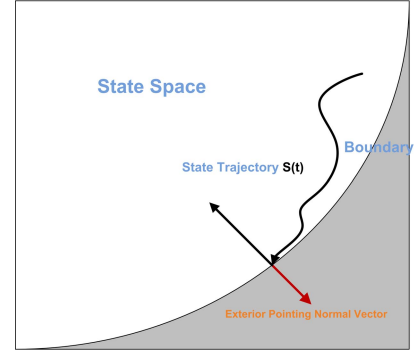


Fig. 5. Illustration of the inward pointing condition when state trajectory reaches the boundary of the state space.

state will not leave the safety region. As Fig. 5 illustrates, the admissible policy should guarantee the state trajectory in the safety domain instead of reaching outwards to the boundary. These are formally stated in the following lemma.

**Lemma 2: Transformation of the Instantaneous Safety Constraint:** The mean-field equilibrium optimization in (32) is equivalent to the following equations:

$$\begin{aligned}
 \min_{\{P, u, v\}} & \left[ P + \sum_{m=1}^M \sum_{j=1}^J \alpha Q_m^j + \int_{S'} \frac{\beta}{\gamma + \|\mathbf{X} - \mathbf{X}'\|} \mu(S') dS' \right. \\
 & + \sum_{m=1}^M \sum_{j=1}^J \nabla_{Q_m^j} V(\mathbf{S}) (A_m^j - \mathbf{s}^j u_m W \log_2(1 + |H_m|^2 P)) \\
 & + \nabla_E V(\mathbf{S}) (E^{ch} - P - E^p) + \mathbf{v}^T \cdot \nabla_{\mathbf{X}} V(\mathbf{S}) \\
 & - \frac{1}{2} a_H \sum_{m=1}^M G_m \nabla_{G_m} V(\mathbf{S}) \\
 & \left. + \frac{1}{2} \text{tr}(\sigma_X \sigma_X^T \nabla_{\mathbf{X}}^2 V(\mathbf{S})) + a_H \sum_{m=1}^M G_m \nabla_{G_m}^2 V(\mathbf{S}) \right] \\
 & = \theta, \mathbf{S} \in \mathcal{S};
 \end{aligned} \tag{34}$$

with boundary conditions

$$\begin{aligned}
 F^{\{P^*, u^*, v^*\}}(\mathbf{S}) \cdot \eta(\mathbf{S}) \Big|_{E = \frac{(c_1 + c_2 v_c^2)}{v_c} (D^{ch}(\mathbf{X}) - D_0)} \\
 \leq 0,
 \end{aligned} \tag{35}$$

$$\nabla_{\mathbf{X}^x} V(\mathbf{S})|_{\mathbf{X}^x = L_x} = \nabla_{\mathbf{X}^y} V(\mathbf{S})|_{\mathbf{X}^y = L_y} \tag{36}$$

$$= \nabla_{\mathbf{X}^z} V(\mathbf{S})|_{\mathbf{X}^z = L_z} = -\varrho_1; \tag{37}$$

$$\nabla_{\mathbf{X}^x} V(\mathbf{S})|_{\mathbf{X}^x = U_x} = \nabla_{\mathbf{X}^y} V(\mathbf{S})|_{\mathbf{X}^y = U_y} \tag{38}$$

$$= \nabla_{\mathbf{X}^z} V(\mathbf{S})|_{\mathbf{X}^z = U_z} = \varrho_2; \tag{39}$$

$$\nabla_{Q_m^j} V(\mathbf{S}) \Big|_{Q_m^j = 0} = 0, \forall j \in \mathcal{J}, m \in \mathcal{M}, \tag{40}$$

$$\nabla_E V(\mathbf{S})|_{E = N_E} = 0; \tag{41}$$

$$\nabla_E V(\mathbf{S})|_{E=0} = -\rho_3; \tag{42}$$

where  $\eta(\mathbf{S})$  denotes the exterior pointing normal vector defined on the boundary  $E = \frac{(c_1 + c_2 v_c^2)}{v_c} (D^{ch}(\mathbf{X}) - D_0)$ .

**Proof:** Please refer to Appendix D.

Based on this, the admissible mean-field optimal control actions can be obtained given the value function as follows:

**Corollary 1 Mean-Field Optimal Control Actions:** The optimal power control policy has the following water-filling



structure (41), as shown at the bottom of this page, where  $u_m^*$  is the optimal user scheduling control as given by:

$$u_m^* = \begin{cases} 1, & m = m^*, \\ 0, & \text{otherwise.} \end{cases} \quad (42)$$

$$m^* = \arg \max_{n \in [1, M]} \sum_{j=1}^J s^j \frac{\partial V(\mathbf{S})}{\partial Q_m^j} \left( W \log_2 \left( \frac{\frac{W}{\ln 2} \frac{\partial V(\mathbf{S})}{\partial Q_m^j} \frac{\sigma^2}{H_m}}{1 - \frac{\partial V(\mathbf{S})}{\partial E}} \right) \right)^+. \quad (43)$$

The optimal UAV velocity policy has the following structure (44), as shown at the bottom of the next page, where  $V(\mathbf{S})$  and  $\mu(\mathbf{S})$  are the joint solutions of HJB-FPK equations (33).

*Proof:* Please refer to Appendix E. ■

## V. DATA-DRIVEN DNN SOLUTION FOR SOLVING HJB-FPK EQUATIONS

The remaining task is to solve the HJB-FPK PDEs (32)-(33). In this section, we will focus on developing a data-driven deep learning algorithm to learn the value function  $V(\mathbf{S})$ , distribution  $\mu(\mathbf{S})$  and control policy  $\Omega(\mathbf{S})$ . Specifically, we construct a model-specific DNN by exploiting the HPM method [33], [34]. The learning agent has system state realizations CSI, QSI, UAV energy queue and UAV position as the real-time input, and produces transmission and trajectory control actions as the real-time output. Compared with generic DNN structures, the proposed structure is more *expressive* in the sense that it can represent the  $V(\mathbf{S})$  and  $\mu(\mathbf{S})$  more efficiently with less number of parameters.

### A. Problem-Specific DNN Architecture Based on Homotopy Perturbation Method

We will first construct the DNN architecture for the PDE learning. Consider the  $I$  layers DNN parameterized approximation structure for the value function and stationary distribution as follows:

$$\hat{V}^I(\mathbf{S}; \mathbf{w}_H) = \sum_{i=0}^I V_i(\mathbf{S}), \quad (45)$$

$$\hat{\mu}^I(\mathbf{S}; \mathbf{w}_F) = \sum_{i=0}^I \mu_i(\mathbf{S}), \quad (46)$$

where  $\mathbf{w}_H$  and  $\mathbf{w}_F$  denote the neural network tunable weighting parameters for HJB and FPK equations respectively.  $\{V_i(\mathbf{S}), \mu_i(\mathbf{S})\}$  denotes the pair of the  $i$ th homotopy deformation layer approximation of value function and distribution  $\{V(\mathbf{S}), \mu(\mathbf{S})\}$ .

Based on the HPM method [33], [34], with the proper selection of the convergence control parameters  $\{\mathbf{w}_H, \mathbf{w}_F\}$ , the homotopy-series solution  $\{\hat{V}^I(x; \mathbf{w}_H), \hat{\mu}^I(x; \mathbf{w}_F)\}$  can be

constructed to approximate the solution of the nonlinear PDE system (32)-(33) accurately

$$\hat{V}^I(\mathbf{S}; \mathbf{w}_H) \approx V(\mathbf{S}), \quad \hat{\mu}^I(\mathbf{S}; \mathbf{w}_F) \approx \mu(\mathbf{S}).$$

The homotopy solutions  $\{\hat{V}^I(x; \mathbf{w}_H), \hat{\mu}^I(x; \mathbf{w}_F)\}$  are governed by the following HPM deformation equations:

$$(1 - \alpha_H(q; \mathbf{w}_H)) \mathcal{I}^H[\hat{V}^I(\mathbf{S}; \mathbf{w}_H) - V_0(\mathbf{S})] = \beta_H(q; \mathbf{w}_H) (\mathcal{B}^H[\hat{V}^I(\mathbf{S}; \mathbf{w}_H), \hat{\mu}^I(\mathbf{S}; \mathbf{w}_F)]) \quad (47)$$

$$(1 - \alpha_F(q; \mathbf{w}_F)) \mathcal{I}^F[\hat{\mu}^I(\mathbf{S}; \mathbf{w}_F) - \mu_0(\mathbf{S})] = \beta_F(q; \mathbf{w}_F) q (\mathcal{B}^F[\hat{V}^I(\mathbf{S}; \mathbf{w}_H), \hat{\mu}^I(\mathbf{S}; \mathbf{w}_F)]) \quad (48)$$

where  $\alpha_H(q; \mathbf{w}_H)$ ,  $\beta_H(q; \mathbf{w}_H)$ ,  $\alpha_F(q; \mathbf{w}_F)$  and  $\beta_F(q; \mathbf{w}_F)$  are the deformation functions with  $\alpha_H(0; \mathbf{w}_H) = \alpha_F(0; \mathbf{w}_F) = \beta_H(0; \mathbf{w}_H) = \beta_F(0; \mathbf{w}_F) = 0$  and  $\alpha_H(1; \mathbf{w}_H) = \alpha_F(1; \mathbf{w}_F) = \beta_H(1; \mathbf{w}_H) = \beta_F(1; \mathbf{w}_F) = 1$ .  $V_0(\mathbf{S})$  and  $\mu_0(\mathbf{S})$  are the initial approximation of the solution to the PDE system.  $q \in [0, 1]$  denotes the embedding perturbation parameter.  $\mathcal{I}^H$  and  $\mathcal{I}^F$  are the identity operators with  $\mathcal{I}^H[V] = V$ ,  $\mathcal{I}^F[\mu] = \mu$ .  $\mathcal{B}^H$  and  $\mathcal{B}^F$  are the nonlinear differential operators of the HJB and FPK equations respectively given by:

$$\begin{aligned} \mathcal{B}^H[V, \mu] &= P^* + \sum_{m=1}^M \sum_{j=1}^J \alpha Q_m^j + \int_{S'} \frac{\beta}{\gamma + \|\mathbf{X} - \mathbf{X}'\|} \mu(S') dS' \\ &+ \sum_{m=1}^M \sum_{j=1}^J \nabla_{Q_m^j} V (A_m^j - s^j u_m^* W \log_2(1 + |H_m|^2 P^*)) \\ &+ \nabla_E V (E^{ch} - P^* - E^p) + \mathbf{v}^{*T} \cdot \nabla_{\mathbf{X}} V \\ &+ \frac{1}{2} \text{tr}(\sigma_X \sigma_X^T \nabla_{\mathbf{X}}^2 V) + \frac{1}{2} a_H \sum_{m=1}^M G_m \nabla_{G_m} V \\ &+ a_H \sum_{m=1}^M G_m \nabla_{G_m}^2 V - \theta \end{aligned} \quad (49)$$

$$\begin{aligned} \mathcal{B}^F[V, \mu] &= \sum_{m=1}^M \sum_{j=1}^J \nabla_{Q_m^j} \left( (A_m^j - s^j u_m^* W \log_2(1 + |H_m|^2 P^*)) \mu \right) \\ &+ \nabla_E ((-P^* - E^p + E^{ch}) \mu) + \nabla_{\mathbf{X}} (\mathbf{v}^{*T} \cdot \mu) \\ &- \frac{1}{2} \text{tr}(\sigma_X \sigma_X^T \nabla_{\mathbf{X}}^2 \mu) + \frac{1}{2} a_H \sum_{m=1}^M \nabla_{G_m} (G_m \mu) \\ &- a_H \sum_{m=1}^M \nabla_{G_m}^2 (G_m \mu) \end{aligned} \quad (50)$$

As the embedding parameter  $q$  varies from 0 to 1, the approximated solutions  $\{\hat{V}^I(\mathbf{S}; \mathbf{w}_H), \hat{\mu}^I(\mathbf{S}; \mathbf{w}_F)\}$

$$P^* = \begin{cases} \sum_{j,m=1}^{J,M} s^j u_m^* \left( \frac{\frac{W}{\ln 2} \frac{\partial V(\mathbf{S})}{\partial Q_m^j}}{1 - \frac{\partial V(\mathbf{S})}{\partial E}} - \frac{\sigma^2}{H_m} \right)^+, & \text{if} \\ 0 & \text{otherwise,} \end{cases} \quad E > D^{ch}(\mathbf{X}) \left( \frac{c_1}{v_c} + c_2 v_c \right), \quad (41)$$

vary from the initial guess  $\{V_0(\mathbf{S}), \mu_0(\mathbf{S})\}$  to the exact solutions  $\{V(\mathbf{S}), \mu(\mathbf{S})\}$ . The approximated solutions  $\{\hat{V}^I(\mathbf{S}; \mathbf{w}_H), \hat{\mu}^I(\mathbf{S}; \mathbf{w}_F)\}$  are subject to the boundary conditions and the normalization condition <sup>7</sup>

$$\begin{aligned} \nabla_{\mathbf{X}^x} \hat{V}^I(\mathbf{S}; \mathbf{w}_H) \Big|_{\mathbf{X}^x=L_x} \\ = \nabla_{\mathbf{X}^y} \hat{V}^I(\mathbf{S}; \mathbf{w}_H) \Big|_{\mathbf{X}^y=L_y}, \end{aligned} \quad (51)$$

$$= \nabla_{\mathbf{X}^z} \hat{V}^I(\mathbf{S}; \mathbf{w}_H) \Big|_{\mathbf{X}^z=L_z} = -\varrho_1;$$

$$\begin{aligned} \nabla_{\mathbf{X}^x} \hat{V}^I(\mathbf{S}; \mathbf{w}_H) \Big|_{\mathbf{X}^x=U_x} \\ = \nabla_{\mathbf{X}^y} \hat{V}^I(\mathbf{S}; \mathbf{w}_H) \Big|_{\mathbf{X}^y=U_y}, \\ = \nabla_{\mathbf{X}^z} \hat{V}^I(\mathbf{S}; \mathbf{w}_H) \Big|_{\mathbf{X}^z=U_z} = \varrho_2; \end{aligned} \quad (52)$$

$$\nabla_{Q_m^j} \hat{V}^I(\mathbf{S}; \mathbf{w}_H) \Big|_{Q_m^j=0} = 0, \forall j \in \mathcal{J}, m \in \mathcal{M}; \quad (53)$$

$$\nabla_E \hat{V}^I(\mathbf{S}; \mathbf{w}_H) \Big|_{E=N_E} = 0; \quad (54)$$

$$\nabla_E \hat{V}^I(\mathbf{S}; \mathbf{w}_H) \Big|_{E=0} = -\rho_3; \quad (55)$$

$$\begin{aligned} \nabla_{\mathbf{X}^x} \hat{\mu}^I(\mathbf{S}; \mathbf{w}_F) \Big|_{\mathbf{X}^x=L_x, U_x} = \nabla_{\mathbf{X}^y} \hat{\mu}^I(\mathbf{S}; \mathbf{w}_F) \Big|_{\mathbf{X}^y=L_y, U_y} \\ = \nabla_{\mathbf{X}^z} \hat{\mu}^I(\mathbf{S}; \mathbf{w}_F) \Big|_{\mathbf{X}^z=L_z, U_z} = 0; \end{aligned} \quad (56)$$

$$\hat{\mu}^I(\mathbf{S}; \mathbf{w}_F) \Big|_{Q_m^j=0} = 0, \forall j \in \mathcal{J}, m \in \mathcal{M}; \quad (57)$$

$$\hat{\mu}^I(\mathbf{S}; \mathbf{w}_F) \Big|_{E=N_E} = 0; \quad (58)$$

$$\begin{aligned} \hat{\mu}^I(\mathbf{S}; \mathbf{w}_F) \Big|_{E=0} = 0; \\ \int_{\mathcal{S}} \hat{\mu}^I(\mathbf{S}; \mathbf{w}_F) d\mathcal{S} = 1 \end{aligned} \quad (59)$$

By substituting the  $I$ th order homotopy-series expansion w.r.t. embedding perturbation parameter  $\{\sum_{i=0}^I V_i(\mathbf{S})q^i, \sum_{i=0}^I \mu_i(\mathbf{S})q^i\}$  into the (47)-(48), and equating the like-power of  $q$ , we have the deformation equation to construct the recursion relationship of the series expansion. By using the neural network to represent the unrolled recursion, the output of the  $i$ th layer can be represented as:

$$V_i(\mathbf{S}) = \chi_i \sum_{j=1}^{i-1} \mathbf{w}_H^{i,j} V_j(\mathbf{S}) \quad (60)$$

$$+ \mathbf{w}_H^{i,i} \mathcal{B}^H[V_{i-1}(\mathbf{S}), \mu_{i-1}(\mathbf{S})],$$

$$\mu_i(\mathbf{S}) = \chi_i \sum_{j=1}^{i-1} \mathbf{w}_F^{i,j} \mu_j(\mathbf{S}) \quad (61)$$

$$+ \mathbf{w}_F^{i,i} \mathcal{B}^F[V_{i-1}(\mathbf{S}), \mu_{i-1}(\mathbf{S})],$$

<sup>7</sup>The boundary conditions of the value function and stationary distribution are derived by the reflecting boundary on the bounded state space.

where

$$\chi_i = \begin{cases} 1, & i > 1, \\ 0, & i \leq 1, \end{cases} \quad (62)$$

The weightings  $\mathbf{w}_H = \{\mathbf{w}_H^{i,j}, 1 \leq j \leq i \leq I\}$  and  $\mathbf{w}_F = \{\mathbf{w}_F^{i,j}, 1 \leq j \leq i \leq I\}$  are the DNN tunable parameters to control the convergence during the learning. Fig.6 illustrates an example of the  $I = 3$ -layer DNN architecture, where each hidden layer is constructed to approximate the recursion relationship in (60)-(61). The final layer of the DNN generates the control actions output based on the real time state input and the approximated value function and stationary distribution.

### B. Online Training Algorithm

Based on the problem-specific DNN architecture, we now describe our DNN based data-driven learning algorithm. Denote  $\{\mathbf{S}^i, i \in \{1, \dots, d\}\}$  and  $\{D_v^i, D_\mu^i, i \in \{1, \dots, d\}\}$  as the finite set of points and boundary operators on the value function and stationary distribution at boundaries according to (51)-(58). The proposed DNN aims to find appropriate DNN weightings  $\{\mathbf{w}_H, \mathbf{w}_F\}$  to minimize the loss function as follows:

$$\begin{aligned} \min_{\{\mathbf{w}_H, \mathbf{w}_F\}} L(\mathbf{w}_H, \mathbf{w}_F) \\ = E^{\hat{\Omega}^*} \left[ \left| \mathcal{B}^H[\hat{V}^I(\mathbf{S}; \mathbf{w}_H), \hat{\mu}^I(\mathbf{S}; \mathbf{w}_F)] \right|^2 \right. \\ + \left| \mathcal{B}^F[\hat{V}^I(\mathbf{S}; \mathbf{w}_H), \hat{\mu}^I(\mathbf{S}; \mathbf{w}_F)] \right|^2 \\ + \sum_{i=1}^d \left( \left| D_v^i(\hat{V}^I(\mathbf{S}^i; \mathbf{w}_H)) \right|^2 \right. \\ + \left. \left| D_\mu^i(\hat{\mu}^I(\mathbf{S}^i; \mathbf{w}_F)) \right|^2 \right) + \left| \int_{\mathcal{S}} \hat{\mu}^I(\mathbf{S}; \mathbf{w}_F) d\mathcal{S} - 1 \right|^2 \Big] \end{aligned} \quad (63)$$

where  $\hat{V}^I(\mathbf{S}; \mathbf{w}_H)$  and  $\hat{\mu}^I(\mathbf{S}; \mathbf{w}_F)$  are the  $I$  layers DNN approximated solution of the HJB-FPK PDE equations. The first two terms in the loss function measure how closely the DNN approximated solution satisfies the PDE operator, and the second two terms  $\sum_{i=1}^d \left( \left| D_v^i(\hat{V}^I) \right|^2 + \left| D_\mu^i(\hat{\mu}^I) \right|^2 \right)$  denote the approximation error of the boundary conditions (51)-(58). The last term  $\left| \int_{\mathcal{S}} \hat{\mu}^I(\mathbf{S}; \mathbf{w}_F) d\mathcal{S} - 1 \right|^2$  denotes the approximation error of the distribution normalization condition (59).  $\hat{\Omega}^*$  denotes the derived optimal control solution based on  $\hat{V}^I(\mathbf{S}; \mathbf{w}_H)$  and  $\hat{\mu}^I(\mathbf{S}; \mathbf{w}_F)$  according to (41)-(44). Note that the loss function depends on the online sampled realizations of system state only and is independent of the target optimal action. As such, the learning process is completely online and autonomous without requiring knowledge of the

$$\mathbf{v}^* = \begin{cases} \frac{-\nabla_{\mathbf{X}} V(\mathbf{S})}{\|\nabla_{\mathbf{X}} V(\mathbf{S})\|_2} \min(v_c, \left\| \frac{\nabla_{\mathbf{X}} V(\mathbf{S})}{2c_2 \frac{\partial V(\mathbf{S})}{\partial E}} \right\|_2), & \text{if} \\ \frac{-v_c(\mathbf{X} - \mathbf{X}^{ch})}{\|\mathbf{X} - \mathbf{X}^{ch}\|_2}, & E > D^{ch}(\mathbf{X})(\frac{c_1}{v_c} + c_2 v_c), \\ \text{otherwise,} & \end{cases} \quad (44)$$

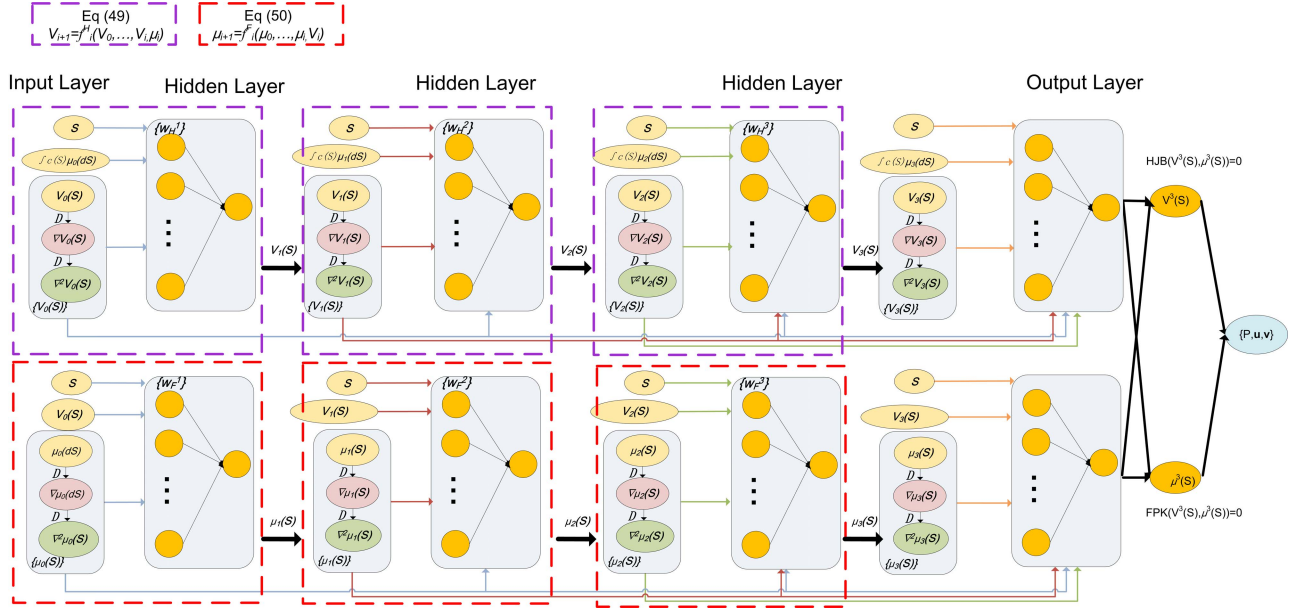


Fig. 6. Problem-specific deep neural network architecture based on (60)-(61) with the approximation order  $I = 3$ .

target optimal action to train the DNN. Based on this, we will apply a stochastic gradient descent approach to update the DNN parameters  $\{w_H, w_F\}$  using the online sampled data. The detailed algorithm is summarized in Fig. 7.

### C. Convergence and Complexity Analysis

We now establish the convergence condition and the approximation power of the proposed DNN learning algorithm for the nonlinear PDEs by proving the approximation error bounds.

**Theorem 3: Convergence of the Proposed Learning Algorithm:** Assume that the initial approximation  $V_0(S)$  and  $\mu_0(S)$  satisfy the boundary conditions (51)-(58). Then we have  $\lim_{I \rightarrow \infty} \hat{V}^I(S) = V(S)$ ,  $\lim_{I \rightarrow \infty} \hat{\mu}^I(S) = \mu(S)$ , and

$$\|\hat{V}^I(S) - V(S)\| \leq \frac{\gamma_H^{I+1}}{1 - \gamma_H} \|V_0(S)\|, \quad (64)$$

$$\|\hat{\mu}^I(S) - \mu(S)\| \leq \frac{\gamma_F^{I+1}}{1 - \gamma_F} \|\mu_0(S)\|, \quad (65)$$

where  $0 < \gamma_H < 1$  and  $0 < \gamma_F < 1$ .

*Proof:* Please refer to Appendix F. ■

**Remark 3:** Note that the proposed DNN learning structure is capable of approximating the high dimensional HJB-FPK PDE solutions in the sense that the approximation error bound (convergence rate) is independent of the number of UAVs  $N$ , the number of users  $M$  and number of content files  $J$ . As a result, the complexity and training time of the proposed solution is free from the curse of dimensionality.

## VI. RESULTS AND DISCUSSIONS

In this section, we shall evaluate the performance of the proposed online data-driven learning scheme for the multi-UAV control through simulation results. As will be described later, the corresponding system setup for the simulation is summarized in Table II, and several benchmark schemes are considered for the performance comparison.

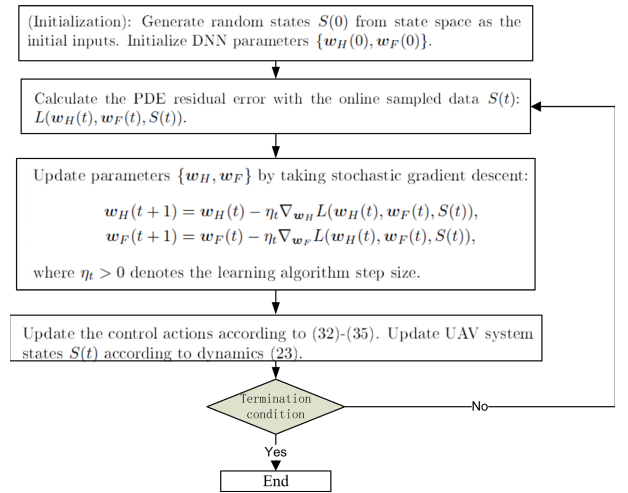


Fig. 7. Algorithm flowchart of the proposed data-driven learning algorithm.

### A. Simulation Setup

In our simulation, the UAVs are constrained in the cube target region with a size of  $200 \times 200 \times 50$  m, and the initial point is set at the center of the region. We consider a multi-users scenario that the  $M$  users are grouped into several user clusters/hotspots, where each cluster is of the same size and uniformly distributed in the target region. The heights of the charging station and each user are assumed to be  $z^{ch} = 50$  m and  $\bar{z}_m = 1$  m,  $\forall m \in \mathcal{M}$ , respectively. We apply 12V Lipo battery for the UAV energy charging with charging current 5A. The overall system setup for the simulation is given in Table II.

In the simulation, the continuous dynamic is discretized by the Euler-Maruyama method with the time slot  $\tau = 5$  ms. We compare the proposed learning scheme with the following baselines using numerical simulations:



TABLE II  
SYSTEM SETUP FOR SIMULATIONS

Parameter	Value	Parameter	Value
$W$	20 MHz	$D_0$	20 m
$[L_x, L_y, L_z]$	$[-100, -100, 50]$	$\rho_0$	1
$[U_x, U_y, U_z]$	$[100, 100, 100]$	$A_{nm}^j, \forall n, m, j$	5Mbps
$\sigma_X^x, \sigma_X^y, \sigma_X^z$	0.5	$P_{\max}$	10 dBW
$a_H$	1	$a$	60W
$v_c$	50m/s	$N_E$	37000 Joule
$M_{\text{UAV}}$	4 kg	$\alpha, \beta, \gamma$	1
$c_1$	272.6	$c_2$	0.0337
$\varrho_1$	50	$\varrho_2$	50

- **(1) Throughput maximization offline design [19], [46]–[48]:** In this scheme, the UAV trajectory, power and user scheduling control are jointly optimized by maximizing the overall throughput:  $\max \mathbf{X}(t), 0 \leq P_n(t) \leq P_{\max}, \sum_m u_{nm}(t) = 1 \sum_{t,m,j} s_n^j(t) R_{nm}(t)$ , s.t.  $\|\mathbf{X}_n(t+1) - \mathbf{X}_n(t)\|_2 \leq v_c, \|\mathbf{X}_n(t) - \mathbf{X}_m(t)\|_2 \geq d_{\min}, \forall n, m$ , where  $P_{\max}$  denotes the maximum power consumption at each time slot, and  $d_{\min}$  denotes the collision avoidance distance between different UAVs.
- **(2) Energy-efficient UAV trajectory offline design [10], [49]:** In this scheme, the UAV trajectory and power control are optimized by solving the following energy-efficient maximization problem:  $\max \mathbf{X}(t), 0 \leq P_n(t) \leq P_{\max}, \sum_m u_{nm}(t) = 1 \frac{\sum_{t,m,j} s_n^j(t) R_{nm}(t)}{\sum_{t,m,j} E_n^p(t)}$ , s.t.  $\|\mathbf{X}_n(t+1) - \mathbf{X}_n(t)\|_2 \leq v_c, \|\mathbf{X}_n(t) - \mathbf{X}_m(t)\|_2 \geq d_{\min}, \forall n, m$  where  $E_n^p(t)$  denotes the per-stage UAV propulsion energy consumption.
- **(3) Queue weighted design:** In this scheme, the UAV moves sequentially and selects the user with the highest request queue length for transmission scheduling. The transmission power is adaptive to both CSI and QSI by solving the following optimization problem:  $\min_{P_n} P_n(t) - \alpha \sum_{m,j} s_n^j(t) Q_{nm}(t) R_{nm}(t)$ .

For baseline (1) and (2), we apply block SCA based iterative algorithm to solve the non-convex optimization problem. For proposed HPM-based learning scheme, we first choose

$$\begin{aligned}
 V_0(\mathbf{S}) &= 1 - \left[ 1 - \frac{\rho_3}{2N_E} E^2 + \rho_3 E \right] \prod_{j,m} [1 - (Q_m^j)^2] \\
 &\times \left[ 1 + \frac{1}{2} \frac{\varrho_1 + \varrho_2}{L_x - U_x} (\mathbf{X}^x)^2 - \frac{L_x \varrho_2 + U_x \varrho_1}{L_x - U_x} \mathbf{X}^x \right] \\
 &\times \left[ 1 + \frac{1}{2} \frac{\varrho_1 + \varrho_2}{L_y - U_y} (\mathbf{X}^y)^2 - \frac{L_y \varrho_2 + U_y \varrho_1}{L_y - U_y} \mathbf{X}^y \right] \\
 &\times \left[ 1 + \frac{1}{2} \frac{\varrho_1 + \varrho_2}{L_z - U_z} (\mathbf{X}^z)^2 - \frac{L_z \varrho_2 + U_z \varrho_1}{L_z - U_z} \mathbf{X}^z \right] \\
 \mu_0(\mathbf{S}) &= \prod_{j,m} \left( 1 - e^{-Q_m^j} \right) (1 - e^{-E - N_E}) (1 - e^E)
 \end{aligned}$$

as the initial guess of  $V(\mathbf{S})$  and  $\mu(\mathbf{S})$  to satisfy the conditions in (51)–(58), where  $\varrho_1, \varrho_2 > 0$  are the parameters specifying boundary conditions in (51)–(52)

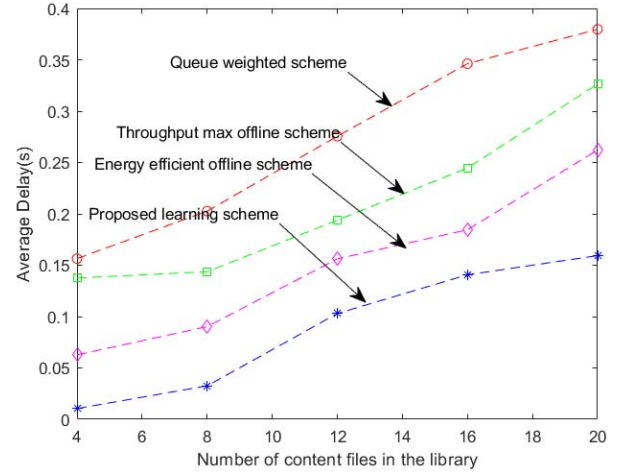


Fig. 8. Long-term average delay performance comparison vs number of content files in the content library with  $M = 60$ ,  $N = 10$  with single caching capacity in UAV.

### B. Multi-UAV Control Performance

Fig. 8 shows the average delay cost (seconds) versus number of content files in the library. As number of content files in the library increases, the average delay increases because each UAV can only serve single content at each time. In Fig. 9, we extended the simulation to consider multiple content files per UAV cache. In both cases, the proposed solution has significant gain over the baselines..

In Fig. 10, we compare the performance of proposed solution applied on heterogeneous and homogeneous UAVs respectively, where in heterogeneous scenario, UAVs are grouped into 3 types with different flight height region and power cost consumption. It can be seen that the proposed solution is quite robust for the case of heterogeneous UAVs.

Fig. 11 illustrates the average delay cost per user versus the UAV numbers. It can be noted that as the density of the multi-UAV system increases, the proposed scheme will be more advantageous when compared to the offline baseline schemes. This is because the adopted mean-field approach benefits multi-UAV system with high density, while the baseline schemes cannot handle the UAVs trajectory and transmission scheduling well in a denser system.

Fig. 12 shows the optimized UAV trajectory obtained by using the proposed online learning scheme with UAV number  $N = 5$ . To be specific, there are 5 user clusters (in the red circles), and each user cluster has 12 users. We assume that users in the different clusters request for the different content files cached in the specific UAVs. For instance, UAV1 only serve the users in Cluster 1, 2 and 5, while UAV2 serve the users in cluster 2, 3 and 4. Each UAV under the proposed learning scheme can control radio resources and trajectory that are adaptive to its energy consumption state and its serving group QoS state, and won't violate its safety constraint due to the boundary conditions in (35)–(39). Note that the queue states carry a higher weight in the value function compared with other state components. Hence, UAV1 flies around from cluster 1 to cluster 5 rather than directly flying to cluster

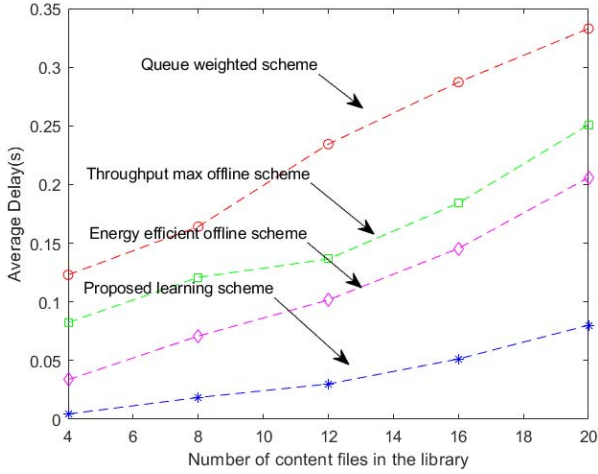


Fig. 9. Long-term average delay cost vs number of content files in the content library with  $M = 60$ ,  $N = 10$ , each UAV can cache 10 files.

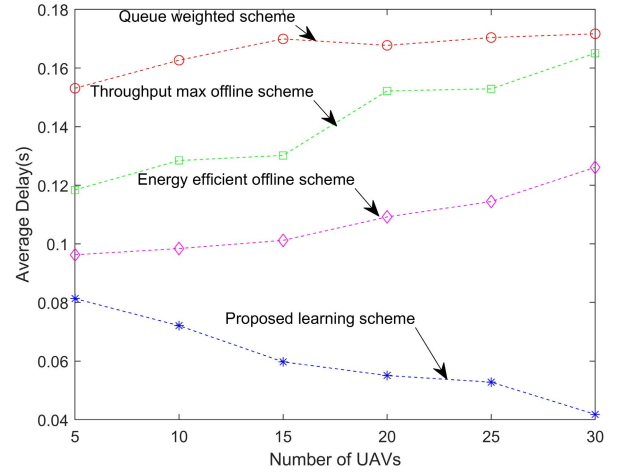


Fig. 11. Average delay vs UAV numbers with  $M = 60$ ,  $J = 10$ .

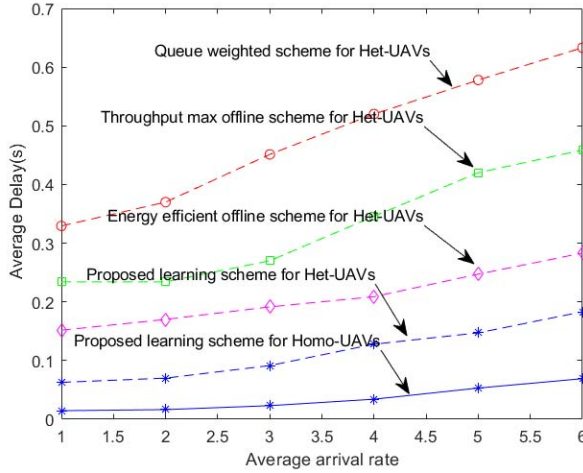


Fig. 10. Performance comparison of proposed solution and baselines applied on heterogeneous and homogeneous UAVs with  $M = 60$ ,  $N = 10$ ,  $J = 10$ .

5 because the request queue state of the users in cluster 1 is much larger than the other clusters.

### C. Convergence and Complexity Analysis

Figure 13 illustrates the approximation errors versus the number of neurons in the DNN. As illustrated, we can achieve similar approximation quality with much smaller size of weight vectors, illustrating the superiority in the expressive power of the proposed neural network. Fig.14 illustrates convergence curve of the proposed learning algorithm. It can be found that the performance converges to the steady point quite fast. Fig.15 shows the performance gap of average delay vs average transmission power between brute-force optimal, proposed learning solution and baselines. It can be seen that the performance of the proposed scheme is very close to the brute-force solution performance, thus the quality of value function  $V(\mathbf{S})$  and stationary distribution  $\mu(\mathbf{S})$  approximation is good w.r.t. the optimal one under the proposed deep learning algorithm. Furthermore, Table III shows the computation time comparison of the baselines, brute-force optimal and proposed scheme. We can see that the computation time of the queue

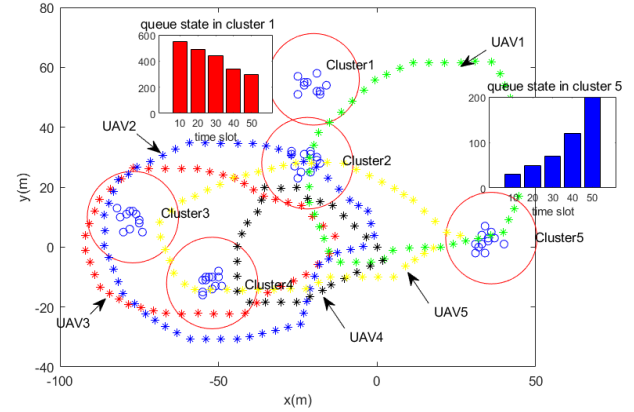


Fig. 12. UAV trajectory under the proposed scheme with number of UAVs  $N = 5$ , number of users  $M = 60$ .

TABLE III

COMPARISON OF AVERAGE COMPUTATION TIME TO COMPUTE THE CONTROL ACTION IN ONE SLOT

	$M = 30$ $N = 5$	$M = 60$ $N = 5$	$M = 60$ $N = 10$
Queue Weighted	15.2ms	16.2ms	18.4ms
Throughput Max	5.621s	5.778s	10.271s
Energy Efficient	5.74s	5.733s	10.134s
Proposed Scheme	11.232ms	11.26ms	11.4ms
B-F Algorithm	$> 10^2$ s	$> 10^2$ s	$> 10^2$ s

weighted baseline is smaller, but it has the worst performance than the other schemes. The brute-force numerical solution has a close performance to that of the proposed scheme, but with the largest computational complexity. While the offline schemes are of high time-consuming and worse than the proposed scheme.

### VII. CONCLUSION

In this paper, we propose an online data-driven multi-UAV trajectory and transmission control scheme. We formulate the optimization problem as an ergodic stochastic differential game to optimize the users' quality-of-experience. By using mean-field game analysis, we derive a reduced-complexity optimality condition, and analyze the corresponding feasible

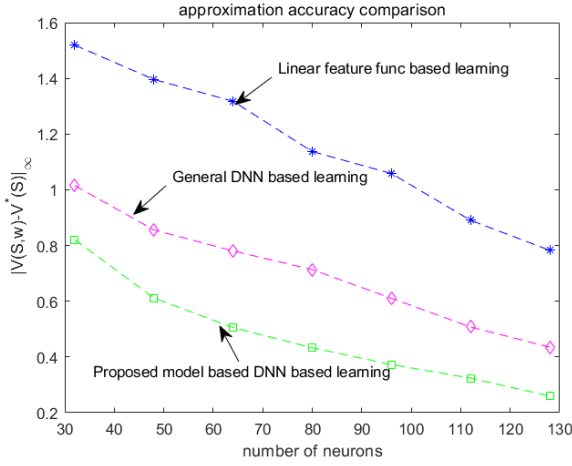


Fig. 13. Approximation accuracy comparison of generic DNN and proposed problem specific DNN illustrating the proposed DNN architecture is more expressive than generic DNN structure.

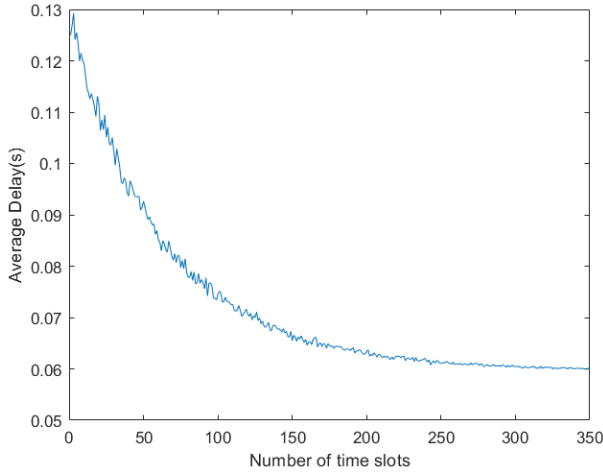


Fig. 14. Convergence of the proposed learning algorithm.

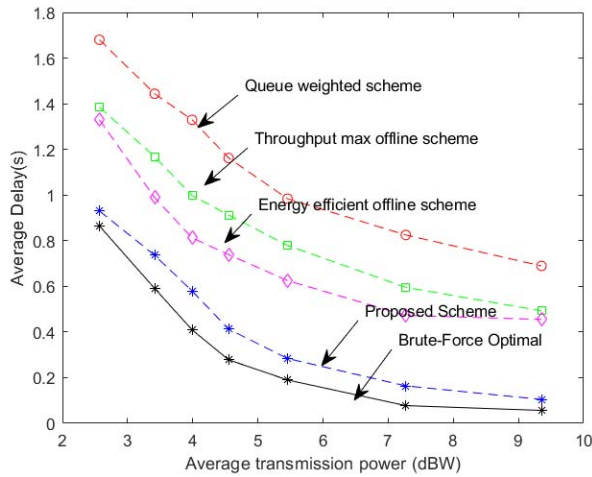


Fig. 15. Average delay performance comparison of proposed learning algorithm and brute-force optimal solution.

control solution under a complex coupled state-action constraint. As solving such a high dimensional PDE is still highly challenging numerically and analytically, we apply a novel

data-driven online learning approach to resolve the problem. Numerical results show that the proposed scheme has much better performance than the existing baselines.

## VIII. APPENDIX A

We have the differential equations as  $dV_n^N(S_n) = \sum_l \frac{dV_n^N(S_n)}{dS_n^l} dS_n^l + \frac{1}{2} \sum_l \sum_k \frac{d^2 V_n^N(S_n)}{dS_n^l dS_n^k} dS_n^l dS_n^k$ , where  $S_n^l$  and  $S_n^k$  denote the state elements in  $S_n$ . By substituting the differential dynamics into this equation with SDE properties, we have

$$\begin{aligned}
 & dV_n^N(S_n(t)) \\
 &= \sum_{j,m} \left[ [A_{nm}^j(t) - s_n^j(t)R_{nm}(t)] \frac{dV_n^N(S_n)}{dQ_{nm}^j} dt \right. \\
 &\quad \left. + \frac{dV_n^N(S_n)}{dQ_{nm}^j} dL_{nm}^j(t) \right] + (-P_n(t) - E_n^p(t) + \\
 &\quad + E_n^{ch}(t)) \frac{dV_n^N(S_n)}{dE_n} dt - \frac{dV_n^N(S_n)}{dE_n} dU_n^E(t) \\
 &\quad + \frac{dV_n^N(S_n)}{dE_n} dL_n^E(t) \\
 &\quad + v_n^T(t) \nabla_{\mathbf{X}_n} V_n^N(S_n) dt + \frac{1}{2} \text{tr}(\sigma_{\mathbf{X}} \sigma_{\mathbf{X}}^T \nabla_{\mathbf{X}_n}^2 V_n^N(S_n)) dt \\
 &\quad + \nabla_{\mathbf{X}_n}^T V_n^N(S_n) dL_n^X(t) - \nabla_{\mathbf{X}_n}^T V_n^N(S_n) dU_n^X(t) \\
 &\quad + \sum_m ((a_H - a_H G_{nm}(t)) \frac{dV_n^N(S_n)}{dG_{nm}} \\
 &\quad + a_H G_{nm}(t) \frac{d^2 V_n^N(S_n)}{dG_{nm}^2}) dt
 \end{aligned} \tag{66}$$

Integrating this equation over the interval  $(0, T)$  under the given policy  $\Omega_n^N$ , dividing  $T$  and subtracting  $\theta_n^N$  on both sides, we can obtain

$$\begin{aligned}
 & \frac{\mathbb{E}^{\Omega_n^N} [V_n^N(S_n(T))] - V_n^N(S_n(0))}{T} - \theta_n^N \\
 &+ \mathbb{E}^{\Omega_n^N} \left[ \frac{1}{T} \int_0^T [c_n(S_n(t), S_{-n}(t), \Omega_n^N)] dt \right] \\
 &= -\theta_n^N + \mathbb{E}^{\Omega_n^N} \left[ \frac{1}{T} \int_0^T \left[ \sum_i \mathcal{L}_i^{\Omega_n^N} (V(S_n(t))) \right] dt \right] \\
 &+ \mathbb{E}^{\Omega_n^N} \left[ \frac{1}{T} \int_0^T [c_n(S_n(t), S_{-n}(t), \Omega_n^N)] dt \right] \\
 &+ \mathbb{E}^{\Omega_n^N} \left[ \frac{1}{T} \int_0^T \sum_{j,m} \frac{dV_n^N(S_n)}{dQ_{nm}^j} \bigg|_{Q_{nm}^j=0} dL_{nm}^j(t) \right. \\
 &\quad \left. - \int_0^T \frac{dV_n^N(S_n)}{dE_n} \bigg|_{E_n=N_E} dU_n^E(t) \right. \\
 &\quad \left. + \int_0^T \frac{dV_n^N(S_n)}{dE_n} \bigg|_{E_n=0} dL_n^E(t) \right. \\
 &\quad \left. + \int_0^T \nabla_{\mathbf{X}_n}^T V_n^N(S_n) \bigg|_{\mathbf{X}_n=[L_x, L_y, L_z]^T} dL_n^X(t) \right. \\
 &\quad \left. - \int_0^T \nabla_{\mathbf{X}_n}^T V_n^N(S_n) \bigg|_{\mathbf{X}_n=[U_x, U_y, U_z]^T} dU_n^X(t) \right],
 \end{aligned}$$



where  $\mathcal{L}_i^{\Omega_i^N}$  is defined as the  $i$ -th UAV dynamic generator with (27). If (19) and the boundary conditions in (20)-(23) are satisfied, for any admissible control policy, we have (after rearranging)  $\frac{\mathbb{E}^{\Omega_n^N} [V_n^N(\mathbf{S}_n(T)) - V_n^N(\mathbf{S}_n(0))]}{T} + \mathbb{E}^{\Omega_n^N} [\frac{1}{T} \int_0^T [c_n(\mathbf{S}_n(t), \mathbf{S}_{-n}(t), \Omega_n^N) dt] + \mathbb{E}^{\Omega_n^N} [\frac{1}{T} \rho_1 (L_n^x(T) + L_n^y(T) + L_n^z(T)) + \frac{1}{T} \rho_2 (U_n^x(T) + U_n^y(T) + U_n^z(T)) + \frac{1}{T} \rho_3 (L_n^E(T))] \geq \theta_n^N$ . After taking  $\limsup$  as  $\frac{1}{T} \rightarrow 0$  on both sides, we have  $\limsup_{T \rightarrow \infty} \frac{\mathbb{E}^{\Omega_n^N} [V_n^N(\mathbf{S}_n(T)) - V_n^N(\mathbf{S}_n(0))]}{T} = 0$ . Therefore,  $\limsup_{T \rightarrow \infty} \mathbb{E}^{\Omega_n^N} [\frac{1}{T} \int_0^T [c_n(\mathbf{S}_n(t), \mathbf{S}_{-n}(t), \Omega_n^N) dt + \rho_1 (L_n^x(T) + L_n^y(T) + L_n^z(T)) + \rho_2 (U_n^x(T) + U_n^y(T) + U_n^z(T)) + \rho_3 (L_n^E(T))] \geq \theta_n^N$ , where the equality can be realized if there exists an admissible policy  $\Omega_n^{N*}$  such that it attains the minimum of the L.H.S. of (19). Hence, policy  $\Omega_n^{N*}(\mathbf{S}) = \{P_n^{N*}, \{u_{nm}^{N*}\}_{m \in \mathcal{M}}, v_n^{N*}\}$  is the optimal control policy for Problem 2.

## IX. APPENDIX B

Based on HJB equations (19) and the admissible control policy, all UAV individuals have the symmetric cost structure and dynamic construction, e.g., state boundary, diffusion coefficients, and it can be noted that the obtained optimal feasible UAV controls are homogeneous. Hence we have the result that the states of UAVs are indistinguishable (exchangeable) [29], [30]: when the number of UAV players tends to infinity, all UAVs exhibit identical behavior and the game can be approximated by mean-field game.

## X. APPENDIX C

(1) We first rewrite the per stage cost as  $c(\mathbf{S}_n, \mathbf{S}_{-n}, \Omega_n^N) = c^1(\mathbf{S}_n, \Omega_n^N) + c^2(\mathbf{S}_n, \frac{1}{N} \sum_{k \neq n}^N \delta \mathbf{S}_k)$ , where  $c^1(\mathbf{S}_n, \Omega_n^N) = \alpha \sum_{m=1}^M \sum_{j=1}^J Q_{nm}^j + P_n^N$ , and  $c^2(\mathbf{S}_n, \frac{1}{N} \sum_{k \neq n}^N \delta \mathbf{S}_k) = \frac{1}{N} \sum_{k \neq n}^N \int_{\mathbf{S}'} \frac{\beta}{\gamma + \|\mathbf{X}'_n - \mathbf{X}'_k\|_2} \delta \mathbf{S}_k d\mathbf{S}'$  with  $\delta \mathbf{S}_k = \delta(\mathbf{S}' = \mathbf{S}_k)$ . Then based on the property of  $N$  independent UAV dynamics, we have the stationary distributions of  $N$  UAVs should be decoupled with  $\mu_n^N(d\mathbf{S}_n)$ ,  $\forall n \in \mathcal{N}$ .

Define that  $\check{c}_n^\mu(\mathbf{S}_n) = \int_{\mathbf{S}^N} c^2(\mathbf{S}_n, \frac{1}{N} \sum_{k \neq n}^N \delta \mathbf{S}_k) \prod_{k \neq n}^N \mu_k^N(d\mathbf{S}_k)$ , we have

$$\min_{\Omega_n^N = \{P_n^N, \{u_{nm}^N\}_{m \in \mathcal{M}}, v_n^N\}} \left( \mathcal{L}_n^{\Omega_n^N} V_n^N + c^1(\mathbf{S}, \Omega_n^N) \right) + \check{c}_n^\mu(\mathbf{S}) = \theta_n^N, \forall n \in \mathcal{N}$$

as  $\int \mathcal{L}_n^{\Omega_n^N} f(\mathbf{S}) \mu_n^N(d\mathbf{S}) = 0$ , for all  $f \in C^2(\mathcal{S})$ ,  $\forall n \in \mathcal{N}$ . Then consider the following equation with the solution  $(\check{V}^N, \check{\theta}^N, \mu^N)$ :

$$\min_{\Omega = \{P, \{u_m\}_{m \in \mathcal{M}}, v\}} (\mathcal{L}^\Omega \check{V}^N + c^1(\mathbf{S}, \Omega)) + \check{c}^\mu(\mathbf{S}) = \check{\theta}^N, \forall n \in \mathcal{N},$$

$$\mathcal{L}^{*\Omega} \mu^N(\mathbf{S}) = 0, \int_{\mathcal{S}} \mu^N(\mathbf{S}) d\mathbf{S} = 1,$$

where  $\check{c}^\mu(\mathbf{S}) = \int_{\mathcal{S}^N} c^2(\mathbf{S}, \frac{1}{N} \sum_{k=1}^N \delta \mathbf{S}_k) \prod_{k=1}^N \mu_k^N(d\mathbf{S}_k)$ . It can be deduced that

$$\|\check{c}_n^\mu - \check{c}^\mu\|_\infty = \sup \left| \int_{\mathcal{S}^N} \left( c^2(\mathbf{S}_n, \frac{1}{N} \sum_{k \neq n}^N \delta \mathbf{S}_k) - c^2(\mathbf{S}_n, \frac{1}{N} \sum_{k=1}^N \delta \mathbf{S}_k) \right) \prod_{k=1}^N \mu_k^N(d\mathbf{S}_k) \right| \leq \frac{\epsilon_1}{N}$$

with constant  $\epsilon_1$  independent of  $N$  as  $\left| c^2(\mathbf{S}_n, \frac{1}{N} \sum_{k \neq n}^N \delta \mathbf{S}_k) - c^2(\mathbf{S}_n, \frac{1}{N} \sum_{k=1}^N \delta \mathbf{S}_k) \right| = \frac{\beta/\gamma}{N}$ . Then we have  $\|\check{c}_n^\mu - \check{c}^\mu\|_\infty \rightarrow 0$  as  $N \rightarrow \infty$ . Based on this, by considering the continuous dependence of HJB equation and FPK equation, we can conclude that as  $N \rightarrow \infty$

$$\sup_n \left( |\theta_n^N - \check{\theta}^N| + \|V_n^N - \check{V}^N\|_{C^2(\mathcal{S})} + \|\mu_n^N - \mu^N\|_\infty \right) \rightarrow 0.$$

Then the convergence in Theorem (2) (1) can be established.

(2) We first have the following result from [50]

$$\lim_{N \rightarrow \infty} \int_{\mathcal{S}^N} c^2(\mathbf{S}_n, \frac{1}{N} \sum_{k=1}^N \delta \mathbf{S}_k) \prod_{k=1}^N \mu(d\mathbf{S}_k) = c^2(\mathbf{S}, \mu). \quad (67)$$

Since  $\lim_{N \rightarrow \infty} \|\mu_n^N - \mu\|_\infty = 0$  according to the Lemma 3.2.5 in [51], by using the result in (1) and triangle inequality, we can conclude that

$$\lim_{N \rightarrow \infty} \left\| \int_{\mathcal{S}^N} c^2(\mathbf{S}_n, \frac{1}{N} \sum_{k=1}^N \delta \mathbf{S}_k) \prod_{k=1}^N \mu(d\mathbf{S}_k) \right. \quad (68)$$

$$\left. - \int_{\mathcal{S}^N} c^2(\mathbf{S}_n, \frac{1}{N} \sum_{k=1}^N \delta \mathbf{S}_k) \prod_{k=1}^N \mu^N(d\mathbf{S}_k) \right\|_\infty = 0. \quad (69)$$

Then we have

$$\lim_{N \rightarrow \infty} \left\| \int_{\mathcal{S}^N} c^2(\mathbf{S}_n, \frac{1}{N} \sum_{k=1}^N \delta \mathbf{S}_k) \prod_{k=1}^N \mu^N(d\mathbf{S}_k) - c^2(\mathbf{S}, \mu) \right\|_\infty = 0.$$

Finally, it is sufficed to conclude that the limit point of (26)-(28) is the solution of

$$\min_{\Omega = \{P, \{u_m\}_{m \in \mathcal{M}}, v\}} (\mathcal{L}^\Omega V(\mathbf{S}) + c^1(\mathbf{S}, \Omega) + c^2(\mathbf{S}, \mu)) = \theta$$

$$\mathcal{L}^{*\Omega} \mu(\mathbf{S}) = 0, \int_{\mathcal{S}} \mu(\mathbf{S}) d\mathbf{S} = 1.$$

## XI. APPENDIX D

The coupled state-action constraint in (15) can be transformed to a state only constraint w.r.t. the UAV energy state and position state:  $E \geq \frac{(c_1 + c_2 v_c^2)}{v_c} (D^{ch}(\mathbf{X}) - D_0)$ . Based on this, ever since the system state reaches the safety boundary  $E = \frac{(c_1 + c_2 v_c^2)}{v_c} (D^{ch}(\mathbf{X}) - D_0)$ , the control action should guarantee the energy and position state pointing inward to the safety region [52], [53], thus we have the equation (35). The value function boundary conditions (36)-(39) come from reflecting boundary property in (20)-(23).

## XII. APPENDIX E

According to the HJB-FPK equation in (33), by solving  $\min_{\Omega(\mathbf{S})=\{P, \mathbf{u}, \mathbf{v}\}} (P + \mathcal{L}^\Omega V(\mathbf{S}))$ , we can obtain the optimal solution in the following form:

$$P^* = \sum_{j,m=1}^{J,M} s^j u_m^* \left( \frac{\frac{W}{\ln 2} \frac{\partial V(\mathbf{S})}{\partial Q_m^j}}{1 - \frac{\partial V(\mathbf{S})}{\partial E}} - \frac{\sigma^2}{H_m} \right)^+,$$

$$u_m^* = \begin{cases} 1, & m = m^*, \\ 0, & \text{otherwise}, \end{cases} \quad (70)$$

where  $m^*$  is defined as (43).

$$\mathbf{v}^* = \frac{v_c \nabla_{\mathbf{X}} V(\mathbf{S}) / (2c_2 \frac{\partial V(\mathbf{S})}{\partial E})}{\left\| \nabla_{\mathbf{X}} V(\mathbf{S}) / (2c_2 \frac{\partial V(\mathbf{S})}{\partial E}) \right\|_2}, \quad (71)$$

Based on the boundary condition in (35), we impose that  $\mathbf{v}^* = \frac{v_c (\mathbf{X} - \mathbf{X}^{ch})}{\left\| \mathbf{X} - \mathbf{X}^{ch} \right\|_2}$ ,  $P^* = 0$  at the state space boundary  $E = \frac{(c_1 + c_2 v_c^2)}{v_c} (D^{ch} - D_0)$ . Specifically, the transmission is suspended, and energy are all utilized for the UAV flight such that it can be guaranteed to get to the recharge station in case of running out of on-board energy.

## XIII. APPENDIX F

We first show that the series constructed by DNN in (60)-(61) converges to the PDE solutions  $V(\mathbf{S})$ ,  $\mu(\mathbf{S})$ . The proof is based on the Banach's fixed point theorem. According to the recursion equations (60)-(61), we have  $\left\| \hat{V}^{I+1}(\mathbf{S}) - \hat{V}^I(\mathbf{S}) \right\| = V_{I+1}(\mathbf{S})$  and  $\left\| \hat{\mu}^{I+1}(\mathbf{S}) - \hat{\mu}^I(\mathbf{S}) \right\| = \mu_{I+1}(\mathbf{S})$ , and we need to show that the series  $\hat{V}^{I+1}(\mathbf{S})$  and  $\hat{\mu}^I(\mathbf{S})$  are the Cauchy sequences in the Banach space. When the auxiliary convergence parameters  $\{\mathbf{w}_H, \mathbf{w}_F\}$  are properly chosen at each recursion update, it can be shown that the series generated by (60)-(61) satisfy

$$\left\| \hat{V}^{I+1}(\mathbf{S}) - \hat{V}^I(\mathbf{S}) \right\| = \left\| V_{I+1}(\mathbf{S}) \right\| \leq \gamma_H \left\| V_I(\mathbf{S}) \right\| \quad (72)$$

$$\leq \gamma_H^2 \left\| V_{I+1}(\mathbf{S}) \right\| \leq \dots$$

$$\leq \gamma_H^{I+1} \left\| V_0(\mathbf{S}) \right\|,$$

$$\left\| \hat{\mu}^{I+1}(\mathbf{S}) - \hat{\mu}^I(\mathbf{S}) \right\| = \left\| \mu_{I+1}(\mathbf{S}) \right\| \leq \gamma_F \left\| \mu_I(\mathbf{S}) \right\| \quad (73)$$

$$\leq \gamma_F^2 \left\| \mu_{I+1}(\mathbf{S}) \right\| \leq \dots$$

$$\leq \gamma_F^{I+1} \left\| \mu_0(\mathbf{S}) \right\|,$$

where  $0 \leq \gamma_H < 1$ ,  $0 \leq \gamma_F < 1$  are the constant independent of the state dimensions. Then we have that  $\lim_{I,J \rightarrow \infty} \left\| \hat{V}^I(\mathbf{S}) - \hat{V}^J(\mathbf{S}) \right\| = 0$ ,  $\lim_{I,J \rightarrow \infty} \left\| \hat{\mu}^I(\mathbf{S}) - \hat{\mu}^J(\mathbf{S}) \right\| = 0$  such that  $\hat{V}^I(\mathbf{S})$  and  $\hat{\mu}^I(\mathbf{S})$  are the Cauchy sequences. As discussed of homotopy analysis method in [34], if the homotopy series solutions  $\hat{V}^I(\mathbf{S})$  and  $\hat{\mu}^I(\mathbf{S})$  are convergent, and the initial guess  $V_0(\mathbf{S})$ ,  $\mu_0(\mathbf{S})$  satisfy the boundary conditions (51)-(58), then the solutions  $\lim_{I \rightarrow \infty} \hat{V}^I(\mathbf{S})$  and  $\lim_{I \rightarrow \infty} \hat{\mu}^I(\mathbf{S})$  must be the exact solutions of nonlinear PDE equations (32)-(33), and satisfy the boundary conditions.

Next we will show the approximation error bound between the proposed solutions and the true solutions in (32)-(33). According to the triangle inequality and the equations (67)-(68), we can derive that  $\forall I, J \in \mathbb{N}$ ,  $J \geq I > 0$ ,

$$\left\| \hat{V}^J(\mathbf{S}) - \hat{V}^I(\mathbf{S}) \right\| \leq \frac{1 - \gamma_H^{J-I}}{1 - \gamma_H} \gamma_H^{I+1} \left\| V_0(\mathbf{S}) \right\|$$

$$\left\| \hat{\mu}^J(\mathbf{S}) - \hat{\mu}^I(\mathbf{S}) \right\| \leq \frac{1 - \gamma_F^{J-I}}{1 - \gamma_F} \gamma_F^{I+1} \left\| \mu_0(\mathbf{S}) \right\|.$$

Then based on that the series generated by the homotopy perturbation converge to the exact PDE solutions as discussed above:  $\lim_{I \rightarrow \infty} \hat{V}^I(\mathbf{S}) = V(\mathbf{S})$  and  $\lim_{I \rightarrow \infty} \hat{\mu}^I(\mathbf{S}) = \mu(\mathbf{S})$ , we can have  $\left\| V(\mathbf{S}) - \hat{V}^I(\mathbf{S}) \right\| \leq \frac{1}{1 - \gamma_H} \gamma_H^{I+1} \left\| V_0(\mathbf{S}) \right\|$  and  $\left\| \mu(\mathbf{S}) - \hat{\mu}^I(\mathbf{S}) \right\| \leq \frac{1}{1 - \gamma_F} \gamma_F^{I+1} \left\| \mu_0(\mathbf{S}) \right\|$ . This completes the theorem (3) proof.

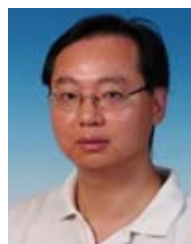
## REFERENCES

- [1] H. Wang, G. Ding, F. Gao, J. Chen, J. Wang, and L. Wang, "Power control in UAV-supported ultra dense networks: Communications, caching, and energy transfer," *IEEE Commun. Mag.*, vol. 56, no. 6, pp. 28–34, Jun. 2018.
- [2] M. Mozaffari, W. Saad, M. Bennis, Y.-H. Nam, and M. Debbah, "A tutorial on UAVs for wireless networks: Applications, challenges, and open problems," 2018, *arXiv:1803.00680*. [Online]. Available: <http://arxiv.org/abs/1803.00680>
- [3] Y. Zeng, Q. Wu, and R. Zhang, "Accessing from the sky: A tutorial on UAV communications for 5G and beyond," 2019, *arXiv:1903.05289*. [Online]. Available: <http://arxiv.org/abs/1903.05289>
- [4] Y. Zeng, J. Lyu, and R. Zhang, "Cellular-connected UAV: Potential, challenges, and promising technologies," *IEEE Wireless Commun.*, vol. 26, no. 1, pp. 120–127, Feb. 2019.
- [5] S. Zhang, H. Zhang, Q. He, K. Bian, and L. Song, "Joint trajectory and power optimization for UAV relay networks," *IEEE Commun. Lett.*, vol. 22, no. 1, pp. 161–164, Jan. 2018.
- [6] 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, Nov. 2016.
- [7] U. Challita and W. Saad, "Network formation in the sky: Unmanned aerial vehicles for multi-hop wireless backhauling," 2017, *arXiv:1707.09132*. [Online]. Available: <http://arxiv.org/abs/1707.09132>
- [8] M. Mozaffari, W. Saad, M. Bennis, and M. Debbah, "Mobile unmanned aerial vehicles (UAVs) for energy-efficient Internet of Things communications," *IEEE Trans. Wireless Commun.*, vol. 16, no. 11, pp. 7574–7589, Nov. 2017.
- [9] M. Jiang, Y. Li, Q. Zhang, and J. Qin, "Joint position and time allocation optimization of UAV enabled time allocation optimization networks," *IEEE Trans. Commun.*, vol. 67, no. 5, pp. 3806–3816, May 2019.
- [10] Y. Zeng and R. Zhang, "Energy-efficient UAV communication with trajectory optimization," *IEEE Trans. Wireless Commun.*, vol. 16, no. 6, pp. 3747–3760, Jun. 2017.
- [11] M. M. Azari, F. Rosas, K.-C. Chen, and S. Pollin, "Ultra reliable UAV communication using altitude and cooperation diversity," *IEEE Trans. Commun.*, vol. 66, no. 1, pp. 330–344, Jan. 2018.
- [12] M. Mozaffari, W. Saad, M. Bennis, and M. Debbah, "Efficient deployment of multiple unmanned aerial vehicles for optimal wireless coverage," *IEEE Commun. Lett.*, vol. 20, no. 8, pp. 1647–1650, Aug. 2016.
- [13] J. Lyu, Y. Zeng, R. Zhang, and T. J. Lim, "Placement optimization of UAV-mounted mobile base stations," *IEEE Commun. Lett.*, vol. 21, no. 3, pp. 604–607, Mar. 2017.
- [14] Y. Chen, W. Feng, and G. Zheng, "Optimum placement of UAV as relays," *IEEE Commun. Lett.*, vol. 22, no. 2, pp. 248–251, Feb. 2018.
- [15] S. Jeong, O. Simeone, and J. Kang, "Mobile edge computing via a UAV-mounted cloudlet: Optimization of bit allocation and path planning," *IEEE Trans. Veh. Technol.*, vol. 67, no. 3, pp. 2049–2063, Mar. 2018.

- [16] 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, Feb. 2019.
- [17] E. Koyuncu, M. Shabanighazikelayeh, and H. Seferoglu, "Deployment and trajectory optimization of UAVs: A quantization theory approach," *IEEE Trans. Wireless Commun.*, vol. 17, no. 12, pp. 8531–8546, Dec. 2018.
- [18] S. Roth, A. Karimenezhad, and A. Sezgin, "Base-stations up in the air: Multi-UAV trajectory control for min-rate maximization in uplink C-RAN," in *Proc. IEEE Int. Conf. Commun. (ICC)*, May 2019, pp. 1–6.
- [19] 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, Mar. 2018.
- [20] W. Han, A. Liu, and V. K. N. Lau, "PHY-caching in 5G wireless networks: Design and analysis," *IEEE Commun. Mag.*, vol. 54, no. 8, pp. 30–36, Aug. 2016.
- [21] M. Chen, M. Mozaffari, W. Saad, C. Yin, M. Debbah, and C. S. Hong, "Caching in the sky: Proactive deployment of cache-enabled unmanned aerial vehicles for optimized quality-of-experience," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 5, pp. 1046–1061, May 2017.
- [22] A. Liu, V. Lau, W. Ding, and E. Yeh, "Mixed-timescale online PHY caching for dual-mode MIMO cooperative networks," 2018, *arXiv:1810.07503*. [Online]. Available: <http://arxiv.org/abs/1810.07503>
- [23] M. Chen, W. Saad, and C. Yin, "Echo-liquid state deep learning for 360° content transmission and caching in wireless VR networks with cellular-connected UAVs," 2018, *arXiv:1804.03284*. [Online]. Available: <https://arxiv.org/abs/1804.03284>
- [24] N. Zhao *et al.*, "Caching unmanned aerial vehicle-enabled small-cell networks: Employing energy-efficient methods that store and retrieve popular content," *IEEE Veh. Technol. Mag.*, vol. 14, no. 1, pp. 71–79, Mar. 2019.
- [25] X. Xu, Y. Zeng, Y. L. Guan, and R. Zhang, "Overcoming endurance issue: UAV-enabled communications with proactive caching," *IEEE J. Sel. Areas Commun.*, vol. 36, no. 6, pp. 1231–1244, Jun. 2018.
- [26] S. Chai and V. K. N. Lau, "Online trajectory and radio resource optimization of cache-enabled UAV wireless networks with content and energy recharging," *IEEE Trans. Signal Process.*, vol. 68, pp. 1286–1299, 2020.
- [27] D. P. Bertsekas, *Dynamic Programming and Optimal Control*, vol. 2, no. 3. Belmont, MA, USA: Athena scientific, 2011.
- [28] O. Guéant, J.-M. Lasry, and P.-L. Lions, "Mean field games and applications," in *Paris-Princeton Lectures on Mathematical Finance 2010*. Berlin, Germany: Springer, 2011, pp. 205–266.
- [29] H. Tembine and M. Huang, "Mean field difference games: McKean-Vlasov dynamics," in *Proc. IEEE Conf. Decis. Control Eur. Control Conf.*, Dec. 2011, pp. 1006–1011.
- [30] K. Hamidouche, W. Saad, M. Debbah, and H. V. Poor, "Mean-field games for distributed caching in ultra-dense small cell networks," in *Proc. Amer. Control Conf. (ACC)*, Jul. 2016, pp. 4699–4704.
- [31] H. Wu, F. Lyu, C. Zhou, J. Chen, L. Wang, and X. Shen, "Optimal UAV caching and trajectory in aerial-assisted vehicular networks: A learning-based approach," *IEEE J. Sel. Areas Commun.*, vol. 38, no. 12, pp. 2783–2797, Dec. 2020.
- [32] W. Shi, J. Li, H. Wu, C. Zhou, N. Cheng, and X. Shen, "Drone-cell trajectory planning and resource allocation for highly mobile networks: A hierarchical DRL approach," *IEEE Internet Things J.*, vol. 8, no. 12, pp. 9800–9813, Jun. 2021.
- [33] J.-H. He, "Homotopy perturbation method for solving boundary value problems," *Phys. Lett. A*, vol. 350, nos. 1–2, pp. 87–88, Jan. 2006.
- [34] S. Liao, *Beyond Perturbation: Introduction to the Homotopy Analysis Method*. Boca Raton, FL, USA: CRC Press, 2003.
- [35] P. Grohs, F. Hornung, A. Jentzen, and P. von Wurstemberger, "A proof that artificial neural networks overcome the curse of dimensionality in the numerical approximation of black-scholes partial differential equations," 2018, *arXiv:1809.02362*. [Online]. Available: <http://arxiv.org/abs/1809.02362>
- [36] Y. Khoo, J. Lu, and L. Ying, "Solving parametric PDE problems with artificial neural networks," 2017, *arXiv:1707.03351*. [Online]. Available: <http://arxiv.org/abs/1707.03351>
- [37] I. E. Lagaris, A. Likas, and D. I. Fotiadis, "Artificial neural networks for solving ordinary and partial differential equations," *IEEE Trans. Neural Netw.*, vol. 9, no. 5, pp. 987–1000, Sep. 1998.
- [38] T. Feng, T. R. Field, and S. Haykin, "Stochastic differential equation theory applied to wireless channels," *IEEE Trans. Commun.*, vol. 55, no. 8, pp. 1478–1483, Aug. 2007.
- [39] Y. Zeng, J. Xu, and R. Zhang, "Energy minimization for wireless communication with rotary-wing UAV," 2018, *arXiv:1804.02238*. [Online]. Available: <http://arxiv.org/abs/1804.02238>
- [40] C. Di Franco and G. Buttazzo, "Energy-aware coverage path planning of UAVs," in *Proc. IEEE Int. Conf. Auto. Robot. Syst. Competitions*, Apr. 2015, pp. 111–117.
- [41] A. Filippone, *Flight Performance of Fixed and Rotary Wing Aircraft*. Amsterdam, The Netherlands: Elsevier, 2006.
- [42] A. R. S. Bramwell, D. Balmford, and G. Done, *Bramwell's Helicopter Dynamics*. Amsterdam, The Netherlands: Elsevier, 2001.
- [43] M. Simic, C. Bil, and V. Vojisavljevic, "Investigation in wireless power transmission for UAV charging," *Procedia Comput. Sci.*, vol. 60, pp. 1846–1855, 2015.
- [44] B. Zhang, C. H. Liu, J. Tang, Z. Xu, J. Ma, and W. Wang, "Learning-based energy-efficient data collection by unmanned vehicles in smart cities," *IEEE Trans. Ind. Informat.*, vol. 14, no. 4, pp. 1666–1676, Apr. 2018.
- [45] T. J. Nugent, Jr. and J. T. Kare, "Laser power beaming for defense and security applications," *Proc. SPIE*, vol. 8045, May 2011, Art. no. 804514.
- [46] F. Cheng *et al.*, "UAV trajectory optimization for data offloading at the edge of multiple cells," *IEEE Trans. Veh. Technol.*, vol. 67, no. 7, pp. 6732–6736, Jul. 2018.
- [47] Y. Zeng, R. Zhang, and T. J. Lim, "Throughput maximization for UAV-enabled mobile relaying systems," *IEEE Trans. Commun.*, vol. 64, no. 12, pp. 4983–4996, Dec. 2016.
- [48] A. Liu and V. K. N. Lau, "Optimization of multi-UAV-aided wireless networking over a ray-tracing channel model," *IEEE Trans. Wireless Commun.*, vol. 18, no. 9, pp. 4518–4530, Sep. 2019.
- [49] J. Zhang, Y. Zeng, and R. Zhang, "Spectrum and energy efficiency maximization in UAV-enabled mobile relaying," in *Proc. IEEE Int. Conf. Commun. (ICC)*, May 2017, pp. 1–6.
- [50] E. Hewitt and L. J. Savage, "Symmetric measures on Cartesian products," *Trans. Amer. Math. Soc.*, vol. 80, no. 2, pp. 470–501, 1955.
- [51] A. Arapostathis, V. S. Borkar, and M. K. Ghosh, *Ergodic Control of Diffusion Processes*, no. 143. Cambridge, U.K.: Cambridge Univ. Press, 2012.
- [52] H. M. Soner, "Optimal control with state-space constraint I," *SIAM J. Control Optim.*, vol. 24, no. 3, pp. 552–561, May 1986.
- [53] I. Capuzzo-Dolcetta and P.-L. Lions, "Hamilton-Jacobi equations with state constraints," *Trans. Amer. Math. Soc.*, vol. 318, no. 2, pp. 643–683, Feb. 1990.



**Shuqi Chai** (Associate Member, IEEE) received the B.Eng. degree in electronic and information engineering from Xi'an Jiaotong University, Xi'an, China, in 2014. She is currently pursuing the Ph.D. degree with the Department of Electronic and Computer Engineering, The Hong Kong University of Science and Technology. Her research interests include dynamic programming in wireless communication systems, reinforcement learning, and UAV communication.



**Vincent K. N. Lau** (Fellow, IEEE) received the B.Eng. degree (Hons.) from The University of Hong Kong, Hong Kong, in 1992, and the Ph.D. degree from Cambridge University, Cambridge, U.K., in 1997. From 1997 to 2004, he was with Bell Labs. In 2004, he was with the Department of Electronic and Computer Engineering, The Hong Kong University of Science and Technology (HKUST), Hong Kong. He is currently a Chair Professor and the Founding Director of the Huawei-HKUST Joint Innovation Laboratory, HKUST. His research interests include robust and delay-optimal cross layer optimization for MIMO/OFDM wireless systems, interference mitigation techniques for wireless networks, massive MIMO, M2M, and network control systems.