Communication-Efficient Massive UAV Online Path Control: Federated Learning Meets Mean-Field Game Theory

Hamid Shiri, Jihong Park, and Mehdi Bennis (Invited Paper)

Abstract—This paper investigates the control of a massive population of UAVs such as drones. The straightforward method of control of UAVs by considering the interactions among them to make a flock requires a huge inter-UAV communication which is impossible to implement in real-time applications. One method of control is to apply the mean field game (MFG) framework which substantially reduces communications among the UAVs. However, to realize this framework, powerful processors are required to obtain the control laws at different UAVs. This requirement limits the usage of the MFG framework for real-time applications such as massive UAV control. Thus, a function approximator based on neural networks (NN) is utilized to approximate the solutions of Hamilton-Jacobi-Bellman (HJB) and Fokker-Planck-Kolmogorov (FPK) equations. Nevertheless, using an approximate solution can violate the conditions for convergence of the MFG framework. Therefore, the federated learning (FL) approach which can share the model parameters of NNs at drones, is proposed with NN based MFG to satisfy the required conditions. The stability analysis of the NN based MFG approach is presented and the performance of the proposed FL-MFG is elaborated by the simulations.

Index Terms—Autonomous UAV, communication-efficient online path control, mean-field game, federated learning.

I. INTRODUCTION

Real-time control of a large number of unmanned aerial vehicles (UAVs) is instrumental in enabling mission-critical applications, such as covering wide disaster sites in emergency cell networks [1], search-and-rescue missions to deliver first-aid packets, and firefighting scenarios [2], [3]. One key challenge is inter-UAV collision, notably under random wind dynamics [1], [4]. A straightforward solution is to exchange instantaneous UAV locations, incurring huge communication overhead, which is thus unfit for real-time operations. Alternatively, in this article we propose a novel real-time massive UAV control framework leveraging mean-field game (MFG) theory [5]–[7] and federated learning (FL) [8], [9].

In our proposed *FL-MFG control* method, each UAV determines its optimal control decision (e.g., acceleration) not by exchanging UAV states (e.g., position and velocity), but by locally estimating the entire UAV population's state distribution, hereafter referred to as *MF distribution*. According to MFG [6], such a distributed control decision asymptotically achieves the epsilon-Nash equilibrium as the number of UAVs goes to infinity. To implement this, the UAV needs to solve a pair of coupled stochastic differential equations (SDEs), namely, the Fokker-Plank-Kolmogorov (FPK) and Hamilton–Jacobi–Bellman (HJB) equations, for the population

Hamid Shiri and Mehdi Bennis are with the Faculty of Information Technology and Electrical Engineering, University of Oulu, 90570 Oulu, Finland (email: {hamid.shiri, mehdi.bennis}@oulu.fi).

Jihong Park is with the School of Information Technology, Deakin University, Geelong, VIC 3220, Australia (email: {jihong.park}@deakin.edu.au).



1

Fig. 1. An illustration of dispatching massive UAVs from a source point to a destination site. Each UAV communicates with neighboring UAVs for achieving: 1) the fastest travel, while jointly minimizing 2) motion energy and 3) inter-UAV collision, under wind perturbations.

distribution estimation and optimal control decision, respectively. The complexity of solving FPK and HJB increases with the state dimension, creating another bottleneck in real-time applications.

To resolve this complexity issue, FL-MFG control utilizes neural-network (NN) based approximations [4], [10], [11] and FL [8]. Specifically, instead of solving HJB and FPK equations, every UAV runs a pair of two NNs, HJB NN and FPK NN whose outputs approximate the solutions of HJB and FPK equations, respectively. The approximation accuracy increases with the number of UAV state observations, i.e., NN training samples. To accelerate the NN training speed, by leveraging FL, each UAV periodically exchanges the HJB NN and FPK NN model parameters with other UAVs, thereby reflecting the locally non-observable training samples. In a source-destination UAV dispatching scenario shown in Fig. 1, simulation results corroborate that FL-MFG control achieves up to 50% shorter travel time, 25% less motion energy, 75%less total transmitted bits, and better collision avoidance measured by 50% lower collision probability number, compared to baseline schemes: FL-MFG exchanging only either HJB NN or FPK NN model parameters, and a control scheme running only HJB NN while exchanging state observations.

A. Background and Related Works

1. UAV Path Planning: Path planning control is about controlling the movement of UAVs to accomplish a target mission. The mission can be broadly categorized into two scenarios. One is to control the UAVs to provision a service, such as ground surveillance [12], emergency networks [1], and hotspot aerial cellular networks [13]. In this case, the mission can be achieved by maximizing the provisioning network coverage [14], surveillance range [15], and traffic offloading rate [16], while minimizing a cost such as the motion and communication energy consumption. The energy consumption is highly dependent on the altitude of UAVs and the types of UAVs (e.g., fixed-wing or multi-copter) [17], as well as communication channels as we shall briefly review in the next part of this subsection.

The other mission is to reach a target destination for the purpose of disaster aid delivery, firefighting, or search and

This work was supported in part by Academy of Finland under Grant 294128, in part by the 6Genesis Flagship under Grant 318927, in part by the Kvantum Institute Strategic Project NOOR, in part by the EU-CHISTERA projects LeadingEdge and CONNECT, and in part by the Academy of Finland through the MISSION Project under Grant 319759.

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication. Citation information: DOI 10.1109/TCOMM.2020.3017281, IEEE Transactions on Communications

2

rescue [2], [3]. To this end, the objective is to minimize the travel time/path, while also minimizing a cost function of their motion and communication energy consumption [17], the risk of collision [4], base-station-UAV disconnection [18] and interference of UAVs on the ground base-stations (BS) [19]. Particularly, the collision avoidance is one key issue in such a mission wherein inter-UAV distances can be very short when collectively lifting a heavy load and/or sharing the same shortest path to the destination. Furthermore, the mission should be executed in real time under a harsh environment [20], [21], e.g., disturbed by wind due to the low altitude in a rescue mission, aggravating the mission completion difficulty, compared to the aforementioned UAVas-an-infrastructure missions wherein the optimal UAV paths can be pre-programmed for their long-term operations under stationary environments. The latter type mission is of prime interest throughout this article.

For a given mission, the UAV path planning can be implemented in an offline or online way. In an offline method, the optimal path is pre-programmed [1], [22], so is vulnerable to environmental dynamics such as random wind perturbations and moving obstacles. Online methods [11], [23] resolve such a problem by continuously updating their models during exploitation, at the cost of the difficulty in training, in terms of accuracy and convergence speed due to the lack of training samples. To accelerate the training speed to achieve higher accuracy in real time, in this work we apply federated learning across multiple UAVs through UAV-to-UAV communications. The communication models will be briefly reviewed in the next part of this subsection.

2. UAV Communication: In the existing literature, UAV communications are grouped into UAV-to-ground communication and inter-UAV communication. For the UAV-to-ground communication, the channel characteristics are sensitive to the blockages and fading near the ground, as well as to the interference from UAVs, ground base stations (BSs), and ground users. This communication scenario includes a remotely controlled UAV path planning wherein UAVs' states are downloaded by a ground station, and the control commands are uploaded to UAVs [11]. Since direct wireless communications between the UAV and its ground controller are limited by the signal attenuation in long-range UAV operations, the 3GPP Release 15 has enabled long-range UAV-to-ground communications through cellular networks [23]-[25]. In this context, interference management is one major issue in UAVto-ground communication. Recent studies have addressed this issue by adjusting the UAV altitude [19], ground BS height and antenna tilting angle [26], while exploiting multi-antenna UAVs [27], transmit power control [28], and optimal subcarrier assignment [29]. For given channel conditions, the UAV-to-ground connectivity should be ensured by limiting the range of UAV operations [18], [30]. Alternatively, one can partly offload the controlling operations to the UAVs so that they can locally carry out decisions even when the connectivity is temporarily lost [11].

On the other hand, inter-UAV communications have different channel characteristics compared to UAV-to-ground communications. For instance, in a high altitude, signal scattering becomes sparse due to the thin atmosphere and lack of obstacles. In this case, only path loss may dominate the channel quality without fading [31]. An example of this case is autonomously controlled UAVs that communicate with each other for collective path planning while avoiding inter-UAV collision [1], [4]. Some of these UAVs may communicate with the ground BS or satellite system [32]. In this case, by balancing the UAV-to-UAV communications, UAV-to-ground communications, and UAV action optimization, in terms of the energy budget and spectrum bandwidth, the UAVs can achieve efficient data transmission rate [33]. In this study, we focus only on the case where UAVs communicate only with neighbors for avoiding collision while accelerating their online training speed.

3. (Windy) Environment: Most works in UAV control rely on the knowledge of the environment and the quality of measurements [34]–[36]. Wind profile information is one of the most significant challenging requirements for the UAV control, which can be obtained from various sources such as land-based systems (e.g., weather station, costly meteorological mast, or portable sonar/lidar sensors), airborne platforms (e.g., tethered balloon, or kite), and UAVs (e.g., by installing anemometers on them or by calculations on the state of the UAV) [37], [38]. Still, by utilizing the most accurate tools, obtaining the exact information about the environment in realistic real-time control applications is not possible.

In addition to wind, presence of static and moving obstacles, threats, other UAVs, and the uncertainties of state observation and the action generators make the real-time path planning of UAV a complex and challenging problem [39]–[41]. For uncertain and dynamic environments, the learning tools such as reinforcement learning models [42], [43], Bayesian methods [44], and models based on MFG [1], [45] can be promising control alternatives since they can adapt themselves within the environment. To the best of our knowledge, there are not many works considering the massive number of UAVs scenario in a windy environment in an online manner all together for real-time applications, and the existing works usually suffer from complex computations or communications cost.

4. UAV Control: In general, two models are used to handle the massive UAV flocking problem. First, the direct control model where a stochastic differential equation (SDE) at each UAV should be solved to obtain the control inputs. Second, the mean-field approximation approach, where the global behavior of the agents is obtained to control the UAVs.

Following [45], in the direct control model, agents can obtain the optimal solution by exchanging the exact information of their states with each other and solving their SDEs. Nevertheless, due to high complexity and cost of communications for high number of agents, the direct SDE method is only limited to a small number of agents.

Two mean-field based approaches, i.e., the mean-field game (MFG) theory and mean-field control (MFC) theory, are developed to study the behavior of a large population by approximating the behavior with mean-field under exchangeability assumptions, and describing it by a pair of partial differential equations [5]–[7].

In MFC, an optimal action rule for the whole population is obtained by solving the optimization problem in a collaborative way. However, in MFG theory which is developed mainly in [5]–[7], the agents are competing in the *N*-player non-cooperative game, where the population behavior is approximated by a mean-field and as a result, the game problem becomes tractable.

In MFG, the optimal action at each agent is obtained by solving a pair of coupled partial differential equations (PDEs): one, called HJB equation which depends on the agent's own state and the interaction term which depends on the distribution of the states of all the agents; and the other one, called FPK equation, describes the evolution of the distribution of the agents depending on the general control rule obtained from HJB equation.

5. ML Aided Control: Many of the numerical methods to solve the HJB-FPK equations are computationally expensive especially because of curse of dimensionality, and are illsuited for real-time applications [46]–[52]. In [46] a numerical approximation method is proposed to approximate the Kolmogorov PDE considering the curse of dimensionality issue. In [47] a type of iterative method for discrete HJB is obtained and its convergence is investigated. In [48] a probabilistic numerical method based on least-square regressions to solve nonlinear HJB equation together with its analysis is obtained. Besides, there are many more numerical methods to solve PDEs as in [49]–[52] with high processing requirements. Therefore, new methods based on machine learning (ML), e.g., (deep) reinforcement learning and neural-network-based online methods, are important to obtain the solution for PDEs with more accuracy or speed [10], [53]–[55]. The (deep) reinforcement learning methods are developed to learn the solution of the HJB equation and control rule in [53]–[55]. In [10] a neural-network-based online solution of the HJB equation is used to explore the infinite horizon optimal robust guaranteed cost control of uncertain nonlinear systems.

Classical centralized multi-agent learning requires communication resources, such as bandwidth and energy, in order to gather all data samples from the agents in the central unit or server. Federated learning is proposed in [8], [9] to enable model training without sharing data. Instead, every agent trains its local model with its local data samples, and the global training model is obtained by sharing and averaging the local models (rather than the local samples). FL has several benefits such as communication cost reduction and privacy preservation. FL was also shown to be useful in enabling URLLC [56]. In addition, there are several research works considering communication cost in control, such as the work in [57] that investigates channel delays in swarm stability for autonomous vehicular platoon systems, [58] that improves communication efficiency of a control system, and [59] which improves communication cost and adaptive estimation error by leveraging sparsity. In summary, there is a need for communication-efficient methods to enable ML and control in autonomous applications such as UAVs.

6. Major Challenges: Based on the research works mentioned above, still there are some important challenges of implementing MFG framework for real-time control of UAVs in a dynamic environment including: 1) except for few cases obtaining an analytical solution for HJB-FPK differential equations is impossible due to its untractability, and the available numerical methods incur high processing power which is not suitable for real-time applications; 2) the approximation methods depending on each agents solution might result in different control rules at the agents which violates the interchangeability condition of MFG; 3) for distributed approximation methods, the effect of communication channel and the payload size is an important issue in multi agents systems; 4) ML based methods require enough samples for training, and have convergence concerns, which are not considered in MFG frameworks.

B. Contributions and Organization

In this paper, we study the real-time control of a large population of UAVs in a windy environment. We start this by explaining the scenario of multiple UAVs to be moved from a starting region to a destination region as shown in Fig.1. Then, two control methods are described in detail, i.e. the HJB control method and the MFG control method to dispatch multiple UAVs quickly, safely, and with low energy consumption. The main contribution of this paper are summarized as follows:

- To reduce the computational cost of HJB control and MFG control, an NN-based function approximator is utilized to approximate the solution of HJB and FPK equations adaptively. The method used here is a variant of our previous work in [4], where one single-layer network with two outputs is utilized to estimate the solution of each HJB and FPK equations. This method gives an approximate solution to the HJB and MFG framework. It is shown in [4] that MFG-based learning method requires less communications cost than HJB-based control.
- To validate the feasibility of the proposed method, the Lyapunov stability analysis is used for HJB and FPK approximation error. These analyses show that the error of approximate solutions for HJB and FPK is bounded, which means that the obtained approximate control actions from NNs are an approximation of the optimal control actions. However, one main requirement for the stability of the approximate solution of MFG is that enough data samples should be provided to update NN's weights, which is challenging in real-time applications.
- To make sure the UAVs do not lack data samples and to mitigate the communication costs and stability concerns of MFG, an FL based MFG strategy is proposed, which will be named as MfgFL-HF control in this paper. This method benefits from the communication reduction property of the FL method for better training of the NN models. The performance and stability of MfgFL-HF are verified by the simulations. It is shown that adopting FL can yield faster, safer, energy-efficient, and communication-efficient control over the baseline methods.

The remainder of this paper is organized as follows. Section II describes the system model for controlling the population of UAVs. Section III explains the HJB and MFG control methods and the necessity to propose an alternative method. Section IV proposes the online NN-base method to obtain an approximate solution for HJB and FPK equations, with their stability analysis brought in the appendices. Section V proposes FL-based MFG methods in detail. Section VI validates the performance of the proposed method by simulations, followed by our conclusions in Section VII.

II. SYSTEM MODEL

Consider the scenario of Fig.1, where a set \mathcal{N} of N UAVs are set to go from a starting position to a specified destination in a windy environment. There are three major issues in this problem as: A) Dynamics of the control system, which reflects the relationship between the parameters of the system, and also the effect of the environment in the system. The more information we have about the environment, the better model we can utilize for control. Here we will assume the wind perturbations as the main source of randomness in the system. B) Control problem, which will consider the costs and interests to formulate a problem where its solution can control the UAVs to the destinations. Here, one major assumption for control is that the number of UAVs, i.e. N, is large. When the number of UAVs is getting larger, the complexity and the risk of the problem increases consequently, especially in the real-time application with expensive UAVs such as UAVs.

C) Channel, in a multi-UAV control, the communications among the UAVs are of critical importance to achieving the control objectives. Here, following the explanations in the Introduction, we will only consider inter-UAV channels, which are modeled as Rician in [60]. In the following, we will consider these three challenges to address the objective of the paper. However, the more focus will be on control with more details on the next sections.

A. Dynamics of Control System

In order to solve the UAV control problem, we should obtain the relationships among its location, velocity, acceleration, and effect of wind on them at the coordinate system. Then let us use a Cartesian coordinate system with the origin at the target position as the global reference coordinate. We define $r_i(t) \in \mathbb{R}^2$ as the vector from the target destination to the current position of *i*-th UAV u_i at time $t \ge 0$. Therefore, the objective of each u_i , $1 \le i \le N$, is to gradually reduce the distance between destination point and the u_i 's current position, by tuning its velocity $v_i(t) \in \mathbb{R}^2$ by controlling the acceleration $a_i(t) \in \mathbb{R}^2$ under random wind dynamics. Following [61], the wind dynamics are assumed to follow an Ornstein-Uhlenbeck process with an average wind velocity v_o . The temporal state dynamics are thereby given as:

$$v_i(t) = a_i(t)dt - c_0 (v_i(t) - v_o) dt + V_o dW_i(t)$$
(1a)
$$dr_i(t) = v_i(t)dt,$$
(1b)

where c_0 is a positive constant, $V_o \in \mathbb{R}^{2 \times 2}$ is the covariance matrix of the wind velocity, and $W_i(t) \in \mathbb{R}^2$ is the standard Wiener process independently and identically distributed (i.i.d.) across UAVs.

Now in order to write the dynamics of the controlled system (1a-1b) in a compact form, let us define the state of each u_i as $s_i(t) \triangleq [r_i(t)^{\mathsf{T}}, v_i(t)^{\mathsf{T}}]^{\mathsf{T}} \in \mathbb{R}^4$; so the SDEs (1a-1b) can be rewritten as

$$ds_i(t) = (As_i(t) + B(a_i(t) + c_0 v_o)) dt + G dW_i(t),$$
(2)

where $A \stackrel{\Delta}{=} \begin{pmatrix} 0 & I \\ 0 & -c_0 I \end{pmatrix}$, $B \stackrel{\Delta}{=} \begin{pmatrix} 0 \\ I \end{pmatrix}$, $G \stackrel{\Delta}{=} \begin{pmatrix} 0 \\ V_o \end{pmatrix}$, and I denotes the two-dimensional identity matrix. Furthermore, by defining $f(s_i(t)) = As_i(t) + c_0 Bv_o$, we rewrite the equation (2), in a compact form as ds

$$a_i(t) = (f(s_i(t)) + Ba_i(t)) dt + G dW_i(t).$$
(3)

B. Control Problem

ď

In general, a model to solve the mentioned control problem should consider three high-level interests for the UAVs. First, travel time minimization: each UAV i should increase speed in the direction to the destination point, to reduce the remaining distance to the destination point while considering to limit the total speed of the UAV. Second, motion energy: each UAV i should reduce the (motion) energy consumption since the UAVs flight time depends on its battery capacity. Lastly, collision avoidance: the collective interest of the whole population is to make a flock of the UAVs traveling together to avoid UAVs colliding each other and also to complete the mission quickly. Nevertheless, there is a trade-off among the interests which should be considered in the control problem.

To achieve the aforementioned points, UAV u_i at time $t < T_f$ aims to minimize its average cost $\psi_i^{a_i}(s_i, t; s_{-i})$, where T_f is the terminal control time, and $s_i(t) = s_i, s_{-i}(t) = s_{-i}$ are the state of UAV u_i and the set of states of all UAVs excluding UAV u_i at time t, respectively. The average is taken with respect to a measure (of the integral inside the expectation)

with a probability distribution depending on (s_i, t) , and it is calculated for the trajectory $\{s_i(\tau)\}_{[t,T]}$ obtained by the control law a_i . The cost $\psi_i^{a_i}(s_i, t; s_{-i})$ consists of the term $g(a_i(\tau), s_i(\tau); s_{-i}(\tau))$ depending on the local state $s_i(\tau)$ and the control action $a_i(\tau)$ with given states of other UAV's as $s_{-i}(\tau)$, i.e.,

$$\psi^{a_i}(s_i, t; s_{-i}) = \mathbb{E}\left[\int_t^{T_f} g(a_i(\tau), s_i(\tau); s_{-i}(\tau)) \mathrm{d}\tau\right], \qquad (4)$$

where \mathbb{E} is the expectation operator, and $g(a_i(\tau), s_i(\tau); s_{-i}(\tau))$

$$g(a_i(\tau), s_i(\tau); s_{-i}(\tau)) = \phi_L(s_i(\tau)) + c_3 ||a_i(\tau)||^2 + c_2 \phi_G(s_i(\tau); s_{-i}(\tau))$$
(5)

in which, the term $\phi_L(s_i(\tau))$ depends only on the local state $s_i(t)$ and the term $\phi_G(s_i(\tau); s_{-i}(\tau))$ relies on the global state $\{s_i(t), s_{-i}(t)\}$, given as:

$$\phi_L(s_i(\tau)) = \frac{v_i(\tau) \cdot r_i(\tau)}{\|r_i(\tau)\|} + c_1 \|v_i(\tau)\|^2, \tag{6}$$

$$\phi_G(s_i(\tau); s_{-i}(\tau)) = \frac{1}{N} \sum_{u_j \in \mathcal{N}} \frac{\|v_j(\tau) - v_i(\tau)\|^2}{\left(\varepsilon + \|r_j(\tau) - r_i(\tau)\|^2\right)^{\beta}}, \quad (7)$$

and the terms c_1 , c_2 , c_3 , β , and ε are positive constants.

The local term $\phi_L(s_i(\tau))$ and the second term in (5) focus on the two objectives, i.e. travel time and motion energy minimization. It is intended to minimize the remaining travel distance $||r_i(\tau)||$ by maximizing the speed towards the destination, i.e., minimizing the projected speed $v_i(\tau) \cdot r_i(\tau) / ||r_i(\tau)||$ towards the opposite direction to the destination. Also, we minimize the kinetic energy and the acceleration control energy by minimizing proxy terms $||v_i(\tau)||^2$ (speed) and $||a_i(\tau)||^2$ (acceleration), respectively [45], [62]. The actual instantaneous motion power consumption $P(\tau)$ of a UAV in the environment, knowing the UAV's speed $||v(\tau)||$, and characteristics of the UAV and air, is calculated by

$$P(\tau) = \lambda_0 \left(1 + \frac{3\|v(\tau)\|^2}{\omega_{\rm tip}^2}\right) + \lambda_1 \left(\sqrt{1 + \frac{\|v(\tau)\|^4}{4\chi_o^4}} - \frac{\|v(\tau)\|^2}{2\chi_o^2}\right)^{\frac{1}{2}} + \lambda_2 \frac{\|v(\tau)\|^3}{2},\tag{8}$$

where λ_0 , λ_1 , and λ_2 , rotor blade tip speed as ω_{tip} , and mean rotor induced velocity in hovering χ_o are the physical characteristics of UAV in the environment [63]. Then, the motion energy E(t) for each UAV *i* at time *t*, is defined as

$$E(t) = \int_{\tau=0}^{t} P(\tau) \mathrm{d}\tau, \qquad (9)$$

which is used as the energy metric to compare different algorithms in our work.

In addition, for comparison purposes, the mission completion metric for each agent i is defined as the time $t = T_i$, when the state $s_i(t)$ of the agent enters the area defined as $\{s_{dest} :$ $||s_{dest}|| = const.$ for the first time, i.e., $||s_i(T_i)|| \le ||s_{dest}||$. Then, the average travel time is defined as $T_{avg} = \sum_{u_i \in \mathcal{N}} T_i$, and the mission completion time is defined as $T_{\max} = \arg \max_{T \in \{T_i\}} T$.

The global term $\phi_G(s_i(\tau); s_{-i}(\tau))$ in (5) refers to collision avoidance and is intended to form a flock of UAVs moving together [64]. The flocking leads to small relative inter-UAV velocities for avoiding collision even when their controlled velocities are slightly perturbed by wind dynamics. A collision happens when the inter-UAV distance is less than a defined distance r_{coll} . Furthermore, the flocking yields closer inter-UAV distances without collision. This is beneficial for allowing more UAVs to exchange their local states through better channel quality, thereby contributing also to collision avoidance. The formation of a flock as mentioned in [65]

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication. Citation information: DOI 10.1109/TCOMM.2020.3017281, IEEE Transactions on Communications

5

is a result of three components: a) separation, i.e. steer to avoid crowding; b) alignment, i.e. steer toward the average heading of neighbors; c) cohesion, i.e. steer toward the average position of neighbors. In view of this, we adopt the Cucker-Smale flocking [1], [64] that reduces the relative speeds for the UAVs. The relative speed $||v_j(\tau) - v_i(\tau)||$ and the inter-UAV distance $||r_j(\tau) - r_i(\tau)||$ are thus incorporated in the numerator and denominator of $\phi_G(s_i(\tau); s_{-i}(\tau))$, respectively. In addition, inspired by [66], the velocity alignment $\phi_A(t)$ and number of collision risks $\phi_C(t)$ as metrics to compare different algorithms, are defined as

$$\phi_{\rm A}(t) = \frac{1}{tN^2} \int_0^t \sum_{u_i \in \mathcal{N}} \sum_{u_j \in \mathcal{N}} \|v_j(\tau) - v_i(\tau)\| \mathrm{d}\tau, \tag{10}$$

$$\phi_{\mathsf{C}}(t) = \frac{1}{tN^2} \int_0^t \sum_{u_i \in \mathcal{N}} \sum_{u_j \in \mathcal{N}} \mathbb{1}_{\|r_j(\tau) - r_i(\tau)\| \le r_{\mathsf{C}}} \mathrm{d}\tau, \qquad (11)$$

where the hazard radius $r_{\rm C}$ defines a dangerous potential collision zone around the UAV. Lower values of $\phi_{\rm A}(t)$ means that the amplitudes of velocity differences between UAVs are small and hence they have made a better flock to travel together. Lower values of $\phi_{\rm C}(t)$ mean that the UAVs do not tend to be too close to each other and the risk of them colliding each other is smaller.

Incorporating the cost (4) under the temporal dynamics (3), the control problem of UAV u_i at time t is formulated as:

$$\psi(s_i, t; s_{-i}) = \min_{a_i} \psi^{a_i}(s_i, t; s_{-i}) \tag{12}$$

s.t.
$$ds_i(t) = (f(s_i(t)) + Ba_i(t)) dt + G dW_i(t),$$
 (13)

The minimum cost $\psi(s_i, t; s_{-i})$ is referred to as the *value* function of the optimal control, and should be derived to obtain the optimal action $a_i(t)$ for UAV u_i . The methods to encounter this problem will be introduced in the following sections.

C. UAV to UAV Wireless Channel Model

In many multi-UAV control problems, communication among the UAVs is a critical condition. However, in the problem of this paper the UAVs will be required to communicate their data with each other while they are moving at a height h. Following [60] for the UAV to UAV communication channels, the Rice model can be used to model both dominant LOS and NLOS paths. The Rice distribution is given by:

$$p_z(\zeta) = \frac{\zeta}{\chi^2} \exp\left(\frac{-\zeta^2 - \xi^2}{2\chi^2}\right) I_0\left(\frac{z\xi}{\chi^2}\right),\tag{14}$$

where $\zeta \geq 0$, and ξ and χ are the strength of LOS and NLOS paths respectively. However, when there is no LOS path between UAVs, this model is reduced to Rayleigh model by setting $\zeta = 0$. Therefore, depending on the scenario, the channel model parameters can be set. Assuming the frequency division multiple access (FDMA) is used for each UAV to UAV communication, with the transmission power P_o , and the distance r_d from a UAV to another UAV, the received signalto-noise (SNR) at each time is

$$SNR = \frac{P_o z r_d^{-\alpha}}{W_o \sigma_n},$$
(15)

where σ_n is the noise power, W_o is the bandwidth, $\alpha \ge 2$ is the path loss exponent, and z is the Rice random variable defined in (14) and is assumed to be independent and identically distributed (i.i.d) across the different UAVs and times.

Another parameter which we will use in this paper is the communication latency. A signal is successfully decoded if the SNR at time t is greater than a target SNR η , i.e., SNR $(t) \ge \eta$.

The number of bits $b_i(D_0)$, transmitted during D_0 time slots, is given as:

$$b_i(D_0) = \theta \sum_{t=1}^{D_0} \mathbb{1}_{\text{SNR}(t) \ge \eta} W_o \log_2(1+\eta),$$
(16)

where θ is the channel coherence time. The latency of transmitting *b* bits is the minimum D_0 , i.e. D_m that satisfies $b_i(D_0) \ge b$. A latency outage occurs when D_m is greater than a predefined threshold D_M , i.e., $D_m > D_M$ in an algorithm.

III. HJB AND MFG CONTROL

At each time instant, each UAV seeks a solution for the defined objective function. The first intuitive method is to analyze the problem directly as solving an HJB equation and then see if there is a need for other alternatives and also have the basis for the proposed methods. However, in the following, it will be clarified that when the number of UAVs is high, the communications and processing complexity will increase to the extent that the real-time implementation will not be possible. Therefore, an alternative MFG method which can reduce the number of communications significantly will be explained.

A. HJB Control

In this method, we assume that all N UAVs perform their optimal action based on the observed local states of all other UAVs. Then according to Bellman's principle for optimal control, for each time $t \leq T_f$ and state s_i of agent *i*, we can rewrite (12) as

$$\psi(s_{i}, t; s_{-i}) = \min_{a_{i}} \mathbb{E} \left[\int_{t}^{T_{f}} g(a_{i}(\tau), s_{i}(\tau); s_{-i}(\tau)) d\tau \right]$$
(17)
$$= \min_{a_{i}} \left\{ g(a_{i}(t), s_{i}; s_{-i}) + \mathbb{E} \left[\psi(s_{i} + ds_{i}, t + dt; s_{-i}) \right] \right\}$$
(18)

By utilizing (13), the second order Taylor expansion of the second term in (18) is calculated as

$$\mathbb{E}\left[\psi(s_{i} + ds_{i}, t + dt; s_{-i})\right] = \left\{\left[\nabla_{s_{i}}\psi(s_{i}, t; s_{-i})\right]^{\mathsf{T}}(f(s_{i}) + Ba_{i}(t)) + \frac{1}{2}\mathrm{tr}(GG^{\mathsf{T}}[\Delta_{s_{i}}\psi(s_{i}, t; s_{-i})]) + \dot{\psi}(s_{i}, t; s_{-i})\right\} \mathrm{d}t + \psi(s_{i}, t; s_{-i}).$$
(19)

Therefore, the HJB equation with this method is obtained as

$$\dot{\psi}(s_i, t; s_{-i}) + \min_{a_i} \{ g(a_i(t), s_i; s_{-i}) + [\nabla_{s_i} \psi(s_i, t; s_{-i})]^{\mathsf{T}} \\ \times (f(s_i) + Ba_i(t)) + \frac{1}{2} \mathrm{tr}(GG^{\mathsf{T}}[\Delta_{s_i} \psi(s_i, t; s_{-i})]) \} = 0,$$
 (20)

where ∇ and Δ denote the gradient and Laplacian operators, respectively. The optimal action at UAV *i* is obtained as

$$a_i(t) = -\frac{1}{2c_3} B^{\mathsf{T}} \nabla_{s_i} \psi(s_i, t; s_{-i}).$$
(21)

Therefore, substituting (21) in (20) yields the HJB equation

$$\dot{\psi}(s_i, t; s_{-i}) + \left(f(s_i) - \frac{1}{4c_3} B B^{\mathsf{T}} \nabla_{s_i} \psi(s_i, t; s_{-i})\right)^{\mathsf{T}} \left[\nabla_{s_i} \psi(s_i, t; s_{-i})\right] \\ + \frac{1}{2} \operatorname{tr}(G G^{\mathsf{T}}[\Delta_{s_i} \psi(s_i, t; s_{-i})]) + \phi_L(s_i) + c_2 \phi_G(s_i; s_{-i}) = 0.$$
(22)

Solving this HJB equation requires an enormous exchange of states among the UAVs, which becomes impossible when the number of UAVs, i.e. N, is high. In order to address this challenge, we leverage the capabilities of the MFG framework explained in the next subsection. Therefore, the UAVs will need to exchange the states only at the beginning of the mission, and after that, they will calculate the optimal actions based on their own state.

^{0090-6778 (}c) 2020 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications_standards/publications/rights/index.html for more information. Authorized licensed use limited to: Oulu University. Downloaded on August 21,2020 at 06:53:25 UTC from IEEE Xplore. Restrictions apply.

B. MFG Control

Another method to encounter this problem is to use MFG framework when the number of UAVs is very high. Let $m_N(s,t)$ be the empirical state distribution function of the UAVs at time instant t defined as

$$m_N(s,t) \stackrel{\Delta}{=} \frac{1}{N} \sum_{j=1}^N \delta(s - s_j(t)), \tag{23}$$

where $\delta(\cdot)$ is the Dirac delta function. Then, the interaction term (7) can be rewritten as

$$\phi_G(s_i(t); s_{-i}(t)) = \int_s m_N(s, t) \ \frac{\|v_i(t) - v\|^2}{(\varepsilon + \|r_i(t) - r\|)^2)^\beta} \mathrm{d}s.$$
(24)

Since the states $s_i(t)$ for i = 1, ..., N are independent and identically distributed which evolve according to SDE (13), utilizing the ergodic theory gives

$$\lim_{N \to \infty} m_N(s,t) = m(s,t), \tag{25}$$

where m(s,t) is the distribution of generic UAV's state, i.e., s, corresponding to the SDE (13) with optimal policy a(t) obtained by (21). The distribution m(s,t), which is called mean field (MF), is the solution of the Fokker-Plank-Kolmogorov (FPK) equation as

$$\dot{m}(s,t) + \nabla_s \cdot [(f(s) + Ba(t))m(s,t)] - \frac{1}{2} \text{tr}(GG^{\mathsf{T}}\Delta_s m(s,t)) = 0, \quad (26)$$

where ∇_{s} denotes the divergence operator, and the initial distribution of the UAVs is given as $m(s,0) = \frac{1}{N} \sum_{j=1}^{N} \delta(s - s_j(0))$. For an optimal action rule a(t), the FPK equation (26) can give the distribution m(s,t) of state of a typical agent s at each time t. In MFG framework, FPK and HJB equations are coupled through the optimal action a(t) obtained by (21) and the interaction term $\phi_G(s_i; s_{-i})$ which can be approximated by $\phi_G(s_i(t); m(s, t))$ defined in (27) (since according to (25), for large N, the actual distribution $m_N(s, t)$ can be approximated by m(s, t)).

$$\phi_G(s_i(t); m(s, t)) \triangleq \int_s m(s, t) \ \frac{\|v_i(t) - v\|^2}{(\varepsilon + \|r_i(t) - r\|)^2)^\beta} \mathrm{d}s$$
(27)

By this definition, the HJB equation (22) can be rewritten as $\dot{\psi}(s_i, t; m(s, t)) + \frac{1}{-} \operatorname{tr}(GG^{\mathsf{T}}[\Delta_{s_i}\psi(s_i, t; m(s, t))])$

$$+ \left(f(s_i) - \frac{1}{4c_3} B B^{\mathsf{T}} \nabla_{s_i} \psi(s_i, t; m(s, t)) \right)^{\mathsf{T}} \nabla_{s_i} \psi(s_i, t; m(s, t)) + \phi_L(s_i) + c_2 \phi_G(s_i; m(s, t)) = 0,$$
(28)

and the corresponding action is

$$a_{i}(t) = -\frac{1}{2c_{3}}B^{\mathsf{T}}\nabla_{s_{i}}\psi(s_{i}, t; m(s, t)),$$
(29)

Therefore, by substituting (29) in (26) the FPK equation can also be rewritten in the following form:

$$\dot{m}(s,t) + \nabla_s \cdot \left[(f(s) - \frac{1}{2c_3} B B^{\mathsf{T}} \nabla_s \psi(s,t;m(s,t))) m(s,t) \right] - \frac{1}{2} \text{tr}(G G^{\mathsf{T}} \Delta_s m(s,t)) = 0$$
(30)

By solving HJB and FPK equation pairs, i.e., (28) and (30), the optimal action for each UAV can be calculated.

IV. NN-BASED HJB AND MFG LEARNING CONTROL -STATE SHARING METHODS

In this section, inspired by [67], we discuss NN-based methods to obtain approximate solutions for HJB and FPK equations for the multiple UAV control application. However, to have better readability and due to lack of space, we will use the simplifications of Table I unless there is a need for complete forms to avoid confusion.

 TABLE I

 LIST OF SIMPLIFIED NOTATIONS.

6

Notation	Simplified Definition	Notation	Simplified Definition
$\psi(s_i, t; m(s, t))$	2/1	m(s,t)	m
$\psi(s_i,t;s_{-i})$	Ψ	$w_{H_d}(t)$	w_{H_d}
$\phi_L(s_i)$	ϕ_L	$\sigma_{H}(s_i; m(s, t))$	σu
$\phi_G(s_i; m(s, t))$	ϕ_G	$\sigma_{H}(s_i; s_{-i})$	Ч
$\phi_G(s_i;s_{-i})$		$\sigma_{F}(s)$	σ_{F}
$f(s_i)$	f	$e_{H}(s_i,t)$	e_{H}
∇_{s_i}	∇	$e_{F}(s,t)$	e_{F}
∇_s		$J_{H}(\hat{w}_{H_{0}}, \hat{w}_{H_{1}})$	J_{H}
Δ_s	Δ	$J_{F}(\hat{w}_{F_{0}}, \hat{w}_{F_{1}})$	J_{F}
Δ_s		$a_i(t)$	a_i
∇_s	$\nabla \cdot$	$\varepsilon_{H}(s_i, t)$	ε _H
$H(s_i, t; m(s, t))$	н	$\varepsilon_{F}(s,t)$	ε _F
$H(s_i, t; s_{-i})$		$L_s(s_i(t))$	L_s
F(s,t;a(t))	F	$R_i(t)$	R_i

A. HJB Learning Control

Here, following our previous work [4], we find an approximate solution to the HJB equation to obtain the corresponding action. Any approximate solution will result in some error, and the approximated HJB equation may not be exactly equal to zero. Therefore, first, based on the defined simplifications above, we represent the HJB equation (22) and (28) by H as

$$\mathsf{H} \triangleq \dot{\psi} + \left(f - \frac{BB^{\mathsf{T}} \nabla \psi}{4c_3} \right)^{\mathsf{T}} \nabla \psi + \phi_L + c_2 \phi_G + \frac{1}{2} \mathrm{tr} (GG^{\mathsf{T}} \Delta \psi) = 0,$$
(31)

where we obtain the ϕ_G empirically by (7) in this subsection.

Similar to [67], given the state distribution of UAVs at each time t, let the function $\psi(s_i, t; s_{-i})$ and its derivative correspondingly be approximated by functions as

$$\hat{\psi}(s_i, t; s_{-i}) \stackrel{\Delta}{=} \hat{w}_{\mathsf{H}_0}(t)^{\mathsf{T}} \sigma_{\mathsf{H}}(s_i; s_{-i}), \qquad (32)$$

$$\hat{\psi}(s_i, t; s_{-i}) \triangleq \hat{w}_{\mathsf{H}_1}(t)^{\mathsf{T}} \sigma_{\mathsf{H}}(s_i; s_{-i}), \qquad (33)$$

where vector functions $\hat{w}_{H_0}(t)$ and $\hat{w}_{H_1}(t)$ are approximations to the optimal vector weight functions $w_{H_0}(t)$ and $w_{H_1}(t)$, respectively, and the value error of these approximations are

$$\varepsilon_{\mathsf{H}_0}(s_i, t) \stackrel{\Delta}{=} \psi(s_i, t; s_{-i}) - w_{\mathsf{H}_0}(t)^{\mathsf{T}} \sigma_{\mathsf{H}}(s_i; s_{-i}), \qquad (34)$$

$$\varepsilon_{\mathsf{H}_{1}}(s_{i},t) \stackrel{\Delta}{=} \dot{\psi}(s_{i},t;s_{-i}) - w_{\mathsf{H}_{1}}(t)^{\mathsf{T}} \sigma_{\mathsf{H}}(s_{i};s_{-i}) \,. \tag{35}$$

Then, using these definitions, and notation simplifications as in Table I, the HJB equation (31) and its approximatation are written as

$$\mathbf{H} = w_{\mathsf{H}_{1}}^{\mathsf{T}} \sigma_{\mathsf{H}} + \left(f - \frac{1}{4c_{3}} BB^{\mathsf{T}} ([\nabla \sigma_{\mathsf{H}}]^{\mathsf{T}} w_{\mathsf{H}_{0}}) \right)^{\mathsf{T}} [\nabla \sigma_{\mathsf{H}}]^{\mathsf{T}} w_{\mathsf{H}_{0}}$$

$$+ \frac{1}{2} \left[\sum_{k=1}^{N} \operatorname{tr} (GG^{T} [\Delta \sigma_{\mathsf{H}}^{[k]}]) \mathbf{e}_{k} \right]^{\mathsf{T}} w_{\mathsf{H}_{0}} + \phi_{L} + c_{2} \hat{\phi}_{G} + \varepsilon_{\mathsf{H}} = 0, \quad (36)$$

$$\hat{\mathsf{H}} = \hat{w}_{\mathsf{H}_{1}}^{\mathsf{T}} \sigma_{\mathsf{H}} + \left(f - \frac{1}{4c_{3}} BB^{\mathsf{T}} ([\nabla \sigma_{\mathsf{H}}]^{\mathsf{T}} \hat{w}_{\mathsf{H}_{0}}) \right)^{\mathsf{T}} [\nabla \sigma_{\mathsf{H}}]^{\mathsf{T}} \hat{w}_{\mathsf{H}_{0}}$$

$$+ \frac{1}{2} \left[\sum_{k=1}^{N} \operatorname{tr} (GG^{T} [\Delta \sigma_{\mathsf{H}}^{[k]}]) \mathbf{e}_{k} \right]^{\mathsf{T}} \hat{w}_{\mathsf{H}_{0}} + \phi_{L} + c_{2} \hat{\phi}_{G}, \quad (37)$$

where the superscript [k] shows the k's element of the corresponding vector, \mathbf{e}_k is a vector with k's element equal to 1 and other elements equal to zero, and $\varepsilon_{\rm H}$ is the error of HJB equation with the function approximator defined as

$$\varepsilon_{\mathsf{H}} \stackrel{\Delta}{=} c_{2} \varepsilon_{\phi_{G}} + \varepsilon_{\mathsf{H}_{1}} - \frac{1}{4c_{3}} [\nabla \varepsilon_{\mathsf{H}_{0}}]^{\mathsf{T}} B B^{\mathsf{T}} [\nabla \sigma_{\mathsf{H}}]^{\mathsf{T}} w_{\mathsf{H}_{0}} - \frac{1}{4c_{3}} [\nabla \varepsilon_{\mathsf{H}_{0}}]^{\mathsf{T}} B B^{\mathsf{T}} \nabla \varepsilon_{\mathsf{H}_{0}} + \left(f - \frac{1}{4c_{3}} B B^{\mathsf{T}} [\nabla \sigma_{\mathsf{H}}]^{\mathsf{T}} w_{\mathsf{H}_{0}} \right)^{\mathsf{T}} \nabla \varepsilon_{\mathsf{H}_{0}} + \frac{1}{2} \operatorname{tr} (G G^{\mathsf{T}} \Delta \varepsilon_{\mathsf{H}_{0}})$$
(38)

^{0090-6778 (}c) 2020 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications_standards/publications/rights/index.html for more information. Authorized licensed use limited to: Oulu University. Downloaded on August 21,2020 at 06:53:25 UTC from IEEE Xplore. Restrictions apply.

Algorithm 1 Hjb control

- 1: **Initialization:** $\hat{w}_{H_0}(0) = 0$ and $\hat{w}_{H_1}(0) = 0$.
- 2: for Each UAV $i = 1, \ldots, N$, in parallel, do
- *Collect* the states $s_{-i}(t)$ from neighboring UAVs. 3:
- 4: Update the weights $\hat{w}_{H_0}(n)$ and $\hat{w}_{H_1}(n)$ by (43) and (44).
- 5:
- Calculate the value $\hat{\psi} = \hat{w}_{H_0}^{\mathsf{T}} \sigma_{\mathsf{H}}.$ Take the optimal action $\hat{a} = -\frac{1}{2c_3} B^{\mathsf{T}} [\nabla \sigma_{\mathsf{H}}]^{\mathsf{T}} \hat{w}_{\mathsf{H}_0}.$ 6:
- 7: end for

where ε_{ϕ_G} is the uncertainty of the interaction term. Then, the corresponding approximate action can be obtained by

$$a = -\frac{1}{2c_3} B^{\mathsf{T}} [\nabla \sigma_{\mathsf{H}}]^{\mathsf{T}} w_{\mathsf{H}_0} - \frac{1}{2c_3} B^{\mathsf{T}} [\nabla \varepsilon_{\mathsf{H}_0}], \tag{39}$$

$$\hat{a} = -\frac{1}{2c_3} B^{\mathsf{T}} [\nabla \sigma_{\mathsf{H}}]^{\mathsf{T}} \hat{w}_{\mathsf{H}_0}.$$

$$\tag{40}$$

Therefore, the error of approximating HJB by NNs is

$$e_{\mathsf{H}} \triangleq \hat{\mathsf{H}} - \mathsf{H}$$

$$= \frac{1}{2c_3} \tilde{w}_{\mathsf{H}_0}^{\mathsf{T}} [\nabla \sigma_{\mathsf{H}}] B B^{\mathsf{T}} [\nabla \sigma_{\mathsf{H}}]^{\mathsf{T}} w_{\mathsf{H}_0} - \frac{1}{2} \left[\sum_{k=1}^{N} \operatorname{tr} (G G^T [\Delta \sigma_{\mathsf{H}}^{[k]}]) \mathbf{e}_k \right]^{\mathsf{T}} \tilde{w}_{\mathsf{H}_0}$$
$$= \tilde{w}_{\mathsf{H}_0}^{\mathsf{T}} [\nabla \sigma_{\mathsf{H}_0}] \mathbf{e}_k [\nabla \sigma_{\mathsf{H}_0}] \mathbf$$

 $\tilde{w}_{\mathsf{H}_{1}}^{\mathsf{I}}\sigma_{\mathsf{H}} - \tilde{w}_{\mathsf{H}_{0}}^{\mathsf{I}} [\nabla\sigma_{\mathsf{H}}]f - \frac{1}{4c_{3}} \tilde{w}_{\mathsf{H}_{0}}^{\mathsf{I}} [\nabla\sigma_{\mathsf{H}}]BB^{\mathsf{T}} [\nabla\sigma_{\mathsf{H}}]^{\mathsf{T}} \tilde{w}_{\mathsf{H}_{0}} - \varepsilon_{\mathsf{H}}, \quad (41)$ where $\tilde{w}_{H_0} \triangleq w_{H_0} - \hat{w}_{H_0}$, and $\tilde{w}_{H_1} \triangleq w_{H_1} - \hat{w}_{H_1}$. The optimal weights w_{H_1} and w_{H_0} should minimize the loss function defined as

$$J_{\mathsf{H}}(\hat{w}_{\mathsf{H}_{0}}, \hat{w}_{\mathsf{H}_{1}}) \stackrel{\Delta}{=} \frac{1}{2} e_{\mathsf{H}}^{\mathsf{T}} e_{\mathsf{H}} + c_{\mathsf{H}} \underbrace{\max\left\{0, \dot{L}_{s}\right\}}_{R_{i}} \mathbb{1}_{\|s_{i}(t)\| \ge s_{\mathsf{dest}}}_{R_{i}} \tag{42}$$

where $c_{\rm H}$ is a positive constant, L_s as the simplified notation of $L_s(s_i(t))$ is a Lyapunov candidate function, and L_s is its derivative with respect to time. The regularizer term shown as as R_i or $R_i(t)$ is meant to stop the movement when reaching the destination, i.e., $||s_i(t)|| = ||[r_i(t)^{\mathsf{T}}, v_i(t)^{\mathsf{T}}]^{\mathsf{T}}|| \le ||s_{\text{dest}}||$. Then, by discretizating the time with small dt steps, the gradient descent updates are written as

$$\hat{w}_{\mathsf{H}_{0}}(n+1) = \hat{w}_{\mathsf{H}_{0}}(n) - \mu_{\mathsf{H}}(\nabla_{\hat{w}_{\mathsf{H}_{0}}}e_{\mathsf{H}})e_{\mathsf{H}} - \mu_{\mathsf{H}}c_{\mathsf{H}}\nabla_{\hat{w}_{\mathsf{H}_{0}}}R_{i}, \qquad (43)$$

$$\hat{w}_{\mathsf{H}_{1}}(n+1) = \hat{w}_{\mathsf{H}_{1}}(n) - \mu_{\mathsf{H}}(\nabla_{\hat{w}_{\mathsf{H}_{1}}}e_{\mathsf{H}})e_{\mathsf{H}}, \tag{44}$$

where the gradients $\nabla_{\hat{w}_{H_0}} e_{H}$ and $\nabla_{\hat{w}_{H_1}} e_{H}$ are obtained as

$$\nabla_{\hat{w}_{\mathsf{H}_{0}}} e_{\mathsf{H}} = \frac{1}{2c_{3}} [\nabla\sigma_{\mathsf{H}}] BB^{\mathsf{T}} [\nabla\sigma_{\mathsf{H}}]^{\mathsf{T}} \tilde{w}_{\mathsf{H}_{0}} - \frac{1}{2c_{3}} [\nabla\sigma_{\mathsf{H}}] BB^{\mathsf{T}} [\nabla\sigma_{\mathsf{H}}]^{\mathsf{T}} w_{\mathsf{H}_{0}} + [\nabla\sigma_{\mathsf{H}}] f + \frac{1}{2} \left[\sum_{k=1}^{N} \operatorname{tr} (GG^{T} [\Delta\sigma_{\mathsf{H}}^{[k]}]) \mathbf{e}_{k} \right], \quad (45)$$

$$\nabla_{\hat{w}_{\mathsf{H}_{1}}} e_{\mathsf{H}} = \sigma_{\mathsf{H}}.\tag{46}$$

The corresponding Hjb learning control based on these update equations is described in Algorithm 1. After initialization of the weights, each UAV has to collect instantaneous states of other UAVs and use it to update the equations (43) and (44). Then, it uses the updated model to take the proper action. However, the stability of this algorithm is explored by the following Proposition 1.

Proposition 1 (HJB Lyapunov stability): For small uncertainty of interaction term, i.e., $\|\varepsilon_{\phi_G}\| \ll 1$, and a bounded interaction term, i.e., $\|\phi_G\| \leq M_1$, the system state and the model weights of constructed adaptive HJB neural network obtained by Algorithm 1 are uniformly ultimately bounded (UUB), i.e., there exist s_{dest} , w_0 , and w_1 at time T such that $||s(t)|| \leq s_{\text{dest}}$, $||w_{\mathsf{H}_0}(t) - \hat{w}_{\mathsf{H}_0}(t)|| \le w_0$, and $||w_{\mathsf{H}_1}(t) - \hat{w}_{\mathsf{H}_1}(t)|| \le w_1$ for all $t \ge T + T'$.

Proof. See Appendix A.

B. MFG Learning Control

Here, we find an approximate solution to the pair of HJB-FPK equations in MFG framework. Regarding the HJB equation, we follow the method explained in previous subsection and by considering that the interaction term is obtained using (27). Then, we follow the similar approximation procedure to approximate the solution for FPK equation. Let us first rewrite the FPK equation (30) by using the simplified notations in Table I, and define F as

$$\mathbf{F} \stackrel{\Delta}{=} \dot{m} + \nabla \cdot \left[(f - \frac{1}{2c_3} B B^{\mathsf{T}} \nabla \psi) m \right] - \frac{1}{2} \operatorname{tr}(G G^{\mathsf{T}} \Delta m) = 0 \tag{47}$$

Using the equality $\nabla \cdot [a\vec{b}] = a\nabla \cdot \vec{b} + \vec{b}^{\mathsf{T}} \nabla a$, where \vec{b} is a vector and a is scalar, we rewrite this FPK equation as

$$\begin{aligned} \mathsf{F} &= \dot{m} + m\nabla \cdot \left[(f - \frac{1}{2c_3} B B^{\mathsf{T}} \nabla \psi) \right] \\ &+ \left[(f - \frac{1}{2c_3} B B^{\mathsf{T}} \nabla \psi) \right]^{\mathsf{T}} \nabla m - \frac{1}{2} \mathrm{tr} (G G^{\mathsf{T}} \Delta m) = 0 \end{aligned} \tag{48}$$

Now, we seek to find an approximate solution to the equation (48). Let us define the linear function approximator $\hat{m}(s,t)$, which approximates the density function $m^*(s,t)$, as

$$\hat{m}(s,t) \stackrel{\Delta}{=} \hat{w}_{\mathsf{F}_0}(t)^{\mathsf{T}} \sigma_{\mathsf{F}}(s), \qquad (49)$$

$$\hat{\dot{n}}(s,t) \stackrel{\Delta}{=} \hat{w}_{\mathsf{F}_{1}}(t)^{\mathsf{T}} \sigma_{\mathsf{F}}(s), \qquad (50)$$

where $\sigma_{\mathsf{F}}(s)$ is a vector of linear or nonlinear functions, and $\hat{w}_{F_0}(t)$ and $\hat{w}_{F_1}(t)$ are the approximation to the optimal weight functions $w_{F_0}(t)$ and $w_{F_1}(t)$ respectively. Then, the errors of approximating the distribution function m(s,t) and its derivative $\dot{m}(s,t)$ are

$$\varepsilon_{\mathsf{F}_0}(s,t) \stackrel{\Delta}{=} m(s,t) - w_{\mathsf{F}_0}(t)^{\mathsf{T}} \sigma_{\mathsf{F}}(s), \qquad (51)$$

$$\varepsilon_{\mathsf{F}_1}(s,t) \stackrel{\Delta}{=} \dot{m}(s,t) - w_{\mathsf{F}_1}(t)^{\mathsf{T}} \sigma_{\mathsf{F}}(s) \,. \tag{52}$$

Considering this definition, and notation simplifications of Table I, the FPK equation (47) and its corresponding approximatation are written as

$$\mathbf{F} = w_{\mathsf{F}_0}^{\mathsf{T}} \sigma_{\mathsf{F}} \nabla \cdot \left[(f - \frac{1}{2c_3} B B^{\mathsf{T}} \nabla \hat{\psi}) \right] + \left[(f - \frac{1}{2c_3} B B^{\mathsf{T}} \nabla \hat{\psi}) \right]^{\mathsf{T}} [\nabla \sigma_{\mathsf{F}}]^{\mathsf{T}} w_{\mathsf{F}_0} - \frac{1}{2} \left[\sum_{k=1}^{N} \operatorname{tr} (G G^T [\Delta \sigma_{\mathsf{F}}^{[k]}]) \mathbf{e}_k \right]^{\mathsf{T}} w_{\mathsf{F}_0} + w_{\mathsf{F}_1}^{\mathsf{T}} \sigma_{\mathsf{F}} + \varepsilon_{\mathsf{F}} = 0,$$
(53)

$$\hat{\boldsymbol{\Xi}} = \hat{w}_{\mathsf{F}_{0}}^{\mathsf{T}} \sigma_{\mathsf{F}} \nabla \cdot \left[(f - \frac{1}{2c_{3}} B B^{\mathsf{T}} \nabla \hat{\psi}) \right] + \left[(f - \frac{1}{2c_{3}} B B^{\mathsf{T}} \nabla \hat{\psi}) \right]^{\mathsf{T}} [\nabla \sigma_{\mathsf{F}}]^{\mathsf{T}} \hat{w}_{\mathsf{F}_{0}} - \frac{1}{2} \left[\sum_{i=1}^{N} \operatorname{tr} (G G^{T} [\Delta \sigma_{\mathsf{F}}^{[k]}]) \mathbf{e}_{k} \right]^{\mathsf{T}} \hat{w}_{\mathsf{F}_{0}} + \hat{w}_{\mathsf{F}_{1}}^{\mathsf{T}} \sigma_{\mathsf{F}},$$
(54)

where $\varepsilon_{\rm F}$ denotes the error of FPK equation caused by the NN, and it is defined as

$$\varepsilon_{\mathsf{F}} \stackrel{\Delta}{=} \nabla \cdot \left[(f - \frac{1}{2c_3} B B^{\mathsf{T}} \nabla \varepsilon_{\psi}) (\hat{w}_{\mathsf{F}_0}^{\mathsf{T}} \sigma_{\mathsf{F}} + \varepsilon_{\mathsf{F}_0}) \right] - \frac{1}{2} \operatorname{tr} (G G^{\mathsf{T}} \Delta \varepsilon_{\mathsf{F}_0}) + \varepsilon_{\mathsf{F}_1} \\ + \varepsilon_{\mathsf{F}_0} \nabla \cdot \left[(f - \frac{1}{2c_3} B B^{\mathsf{T}} \nabla \hat{\psi}) \right] + \left[(f - \frac{1}{2c_3} B B^{\mathsf{T}} \nabla \hat{\psi}) \right]^{\mathsf{T}} [\nabla \varepsilon_{\mathsf{F}_0}],$$
(55)

where ε_{ψ} is the uncertainty in finding ψ . Therefore, the error of approximating FPK equation by neural networks is

$$e_{\mathsf{F}} \stackrel{\Delta}{=} \hat{\mathsf{F}} - \mathsf{F}$$

$$= -\tilde{w}_{\mathsf{F}_{0}}^{\mathsf{T}} \sigma_{\mathsf{F}} \nabla \cdot \left[(f - \frac{BB^{\mathsf{T}} \nabla \hat{\psi}}{2c_{3}}) \right] - \tilde{w}_{\mathsf{F}_{0}}^{\mathsf{T}} [\nabla \sigma_{\mathsf{F}}] \left[(f - \frac{BB^{\mathsf{T}} \nabla \hat{\psi}}{2c_{3}}) \right]$$

$$+ \frac{1}{2} \tilde{w}_{\mathsf{F}_{0}}^{\mathsf{T}} \left[\sum_{k=1}^{N} \operatorname{tr}(GG^{T} [\Delta \sigma_{\mathsf{F}}^{[k]}]) \mathbf{e}_{k} \right] - \tilde{w}_{\mathsf{F}_{1}}^{\mathsf{T}} \sigma_{\mathsf{F}} - \varepsilon_{\mathsf{F}}, \quad (56)$$

where $\tilde{w}_{F_0} \stackrel{\Delta}{=} w_{F_0} - \hat{w}_{F_0}$, and $\tilde{w}_{F_1} \stackrel{\Delta}{=} w_{F_1} - \hat{w}_{F_1}$. Based on these definitions, the optimal weights \hat{w}_{F_1} and \hat{w}_{F_0} should minimize the loss function defined as

Algorithm 2 Mfg control

- 1: **Initialization:** $\hat{m}(\overline{s,0}) = \frac{1}{N} \sum_{j=1}^{N} \delta(s-s_j(0)), \ \hat{w}_{\mathsf{H}_0}(0) = 0, \ \hat{w}_{\mathsf{H}_1}(0) = 0, \ \hat{w}_{\mathsf{F}_0}(0) = 0, \ \hat{w}_{\mathsf{F}_1}(0) = 0.$
- 2: for Each UAV i = 1, ..., N, in parallel, do
- for $n = 1, ..., T_0$ do 3:
- Update weights $\hat{w}_{H_0}(n)$ and $\hat{w}_{H_1}(n)$ by (43) and (44). 4: 5:
- Calculate value $\hat{\psi} = \hat{w}_{\mathsf{H}_0}^{\mathsf{T}} \sigma_{\mathsf{H}}$.
- Update weight $\hat{w}_{\mathsf{F}_0}(n)$ and $\hat{w}_{\mathsf{F}_1}(n)$ by (58) and (59). 6. Obtain MF distribution $\hat{m} = \hat{w}_{\mathsf{F}_0}^{\mathsf{T}} \sigma_{\mathsf{F}}$. 7:
- end for 8:
- Take the optimal action $\hat{a} = -\frac{1}{2c_3}B^{\mathsf{T}}[\nabla\sigma_{\mathsf{H}}]^{\mathsf{T}}\hat{w}_{\mathsf{H}_0}$. 9: 10: end for

$$J_{\mathsf{F}}(\hat{w}_{\mathsf{F}_{0}}, \hat{w}_{\mathsf{F}_{1}}) = \frac{1}{2} e_{\mathsf{F}}^{\mathsf{T}} e_{\mathsf{F}}$$
(57)

Therefore, by discretizating the time with dt time steps, the gradient descent updates for FPK weights are obtained as

$$\hat{w}_{\mathsf{F}_{0}}(n+1) = \hat{w}_{\mathsf{F}_{0}}(n) - \mu_{\mathsf{F}}(\nabla_{\hat{w}_{\mathsf{F}_{0}}}e_{\mathsf{F}})e_{\mathsf{F}}$$
(58)

$$\hat{w}_{\mathsf{F}_{1}}(n+1) = \hat{w}_{\mathsf{F}_{1}}(n) - \mu_{\mathsf{F}}(\nabla_{\hat{w}_{\mathsf{F}_{1}}}e_{\mathsf{H}})e_{\mathsf{F}}$$
(59)

where the gradients $\nabla_{\hat{w}_{\mathsf{F}_0}} e_{\mathsf{F}}$ and $\nabla_{\hat{w}_{\mathsf{F}_1}} e_{\mathsf{F}}$ are calculated as

$$\nabla_{\hat{w}_{\mathsf{F}_{0}}} e_{\mathsf{F}} = \sigma_{\mathsf{F}} \nabla \cdot \left[\left(f - \frac{1}{2c_{3}} B B^{\mathsf{T}} \nabla \hat{\psi} \right) \right] + \left[\nabla \sigma_{\mathsf{F}} \right] \left[\left(f - \frac{1}{2c_{3}} B B^{\mathsf{T}} \nabla \hat{\psi} \right) \right] - \frac{1}{2} \left[\sum_{k=1}^{N} \operatorname{tr} (G G^{T} [\Delta \sigma_{\mathsf{F}}^{[k]}]) \mathbf{e}_{k} \right]$$
(60)

$$\nabla_{\hat{w}_{\mathsf{F}_{1}}} e_{\mathsf{F}} = \sigma_{\mathsf{F}}.\tag{61}$$

However, based on these update pairs and update pairs for HJB equation, the Mfg learning algorithm is described as in Algorithm 2. In this algorithm, first, the UAVs share their states $s_i(0)$ at time t = 0 to obtain the distribution of the population and initial samples. Then, the UAVs start collecting samples and and updating theirs weights until they reach the destination at time $T = T_0 dt$.

We have shown in our previous works [4] that Mfg control algorithm provides better results in terms of energy consumption, communications cost, and flocking of UAVs when sufficient samples are used in model training. However, there are stability concerns which is analyzed in this section. The MFG learning solution consists of two coupled NNs coined HJB NN and FPK NN, which should be stable to ensure stability of MFG learning. Proving stability of the coupled HJB-FPK equations is a challenging issue. To simplify the stability analysis, we decouple the HJB and FPK equations and separately obtain the stability conditions for each of them. Then stability can be proved, when the initial conditions of HJB NN and FPK NN meet the stability conditions required for each of them separately, and if the output of each HJB NN (and FPK NN) falls into the stability space of the other FPK NN (and HJB NN). Starting with the HJB NN, we obtained the stability condition for HJB NN in Proposition 1. In Propositions 2 and 3, we show the stability and convergence conditions required for FPK NN. Then we conclude the stability of MFG NN in Corollary 1.

Proposition 2 (FPK Lyapunov stability): For almost certain ψ , i.e., $\|\varepsilon_{\psi}\| \ll 1$, and differentiable and bounded value function ψ , i.e., $\|\psi\| \leq M_2$, the weights of constructed adaptive FPK neural network obtained in Algorithm 2, which is controlled by its corresponding HJB equation, are UUB, i.e., there exist w_2 , and w_3 at time T such that $||w_{\mathsf{F}_0}(t) - \hat{w}_{\mathsf{F}_0}(t)|| \le w_2$, and $||w_{\mathsf{F}_1}(t) - \hat{w}_{\mathsf{F}_1}(t)|| \le w_3 \text{ for all } t \ge T + T'.$

Proof. See Appendix B.

Proposition 3 (FPK Convergence): Under the assumptions of Proposition 2 and with small step-sizes $\mu_{\rm F}$, the weights of FPK neural network function approximator converges to its optimal weights in mean with no bias and it is stable in mean square deviation sense.

Proof. See Appendix C.
$$\Box$$

Corollary 1: Considering Propositions 1, 2, and 3, we can conclude that the system state and weights of constructed HJB-FPK neural networks obtained by Algorithm 2 are UUB.

Proof. The Algorithm 2 has two parts as HJB part and FPK parts. It is initialized by the states of the UAVs at the source region. At the initial iterations the states are used directly in the algorithm to update the weights of HJB and FPK neural networks, and hence the uncertainty of interaction term is small, i.e., $\|\varepsilon_{\phi_G}\| \ll 1$. Also, by starting with well trained or zero initialized neural network weights, both the interaction term and value functions are upper-bounded, i.e., $\|\phi_G\| \leq M_1$ and $\|\psi\| \leq M_2$. Hence, there is a design, i.e., a choice of parameters in proofs in Appendices A, B, and C, such that all assumptions necessary for Propositions 1, 2, and 3 hold together and completely

V. FEDERATED MFG LEARNING - MODEL SHARING **METHODS**

In this section, we propose federated mean field game learning strategy (MfgFL) and its different implementations, i.e., MfgFL-H, MfgFL-F, and MfgFL-HF, to make the UAVs' control models close to each other and to use sample diversity among the UAVs efficiently. In MfgFL-H the model parameters of HJB neural network are shared with central unit to obtain the global HJB NN model, so the action rules of UAVs are close to each other. In MfgFL-F the FPK NN models are transmitted to the central unit to obtain the global FPK NN model, so the estimation of the population density function at UAVs be more accurate. In MfgFL-HF, both HJB and FPK neural network models are averaged to obtain better global online MFG learning model. In the following, we explain general form of MfgFL strategy, which covers three different implementations.

Although Algorithm 2 can reduce the communications cost of the control algorithm by leveraging the MFG framework, it still requires big sample sets to train and provide conditions of stability. In other words, there is still a need to share a subset of samples among the UAVs or with a central unit, which requires extra communication costs in addition to privacy concerns. Therefore, instead of state sharing, we adopt the federated learning method to address these issues.

In the MfgFL algorithm, one UAV out of all is set to act as a control center, which we call is as leader (or header) UAV. This leader UAV depending on the application may be chosen randomly or considering UAVs power consumption or flight time, which is beyond the scope of this work. Then, we simply set one of the UAVs as the leader, i.e., u_h , and indicate it by index h in the algorithm.

The proposed MfgFL learning control is described in Algorithm 3. Following the FedAvg algorithm, the leader collects models $\hat{w}_{i,d}(n)$ of N_h UAVs at times $n = kn_0$, where $d \in \{H, F, HF\}$ which corresponds to three types of implementations as

Algorithm 3 MfgFL control

- 1: **Initialization:** $\hat{m}(s,0) = \frac{1}{N} \sum_{j=1}^{N} \delta(s-s_j(0)), \ \hat{w}_{\mathsf{H}_0}(0) = 0, \ \hat{w}_{\mathsf{F}_0}(0) = 0, \ \hat{w}_{\mathsf{F}_1}(0) = 0.$
- 2: for $n = 0, 1, 2, \dots, T_0$ do
- 3: if $n = kn_0$ then
- 4: N_h UAVs, in parallel, *send* their model $\hat{w}_{i,d}(kn_0)$ to the leader.
- 5: Leader *updates* the model parameters $\hat{w}_{h,d}(k)$, via (62).
- 6: Leader *broadcasts* the model $\hat{w}_{h,d}(k)$.
- 7: **end if**
- 8: for each UAV i = 1, ..., N, in parallel, do
- 9: **if** UAV *i* receives $\hat{w}_{h,d}(k)$ **then**
- 10: Update $\hat{w}_{i,d}(n)$, as $\hat{w}_{i,d}(n) \leftarrow \hat{w}_{h,d}(k)$.
- 11: end if
- 12: Update $\hat{w}_{H_0}(n)$, $\hat{w}_{H_1}(n)$, $\hat{w}_{F_0}(n)$, and $\hat{w}_{F_1}(n)$ by (43) and (44), (58), and (59). *Take* the optimal action $\hat{a} = -\frac{1}{2c_3}B^{\mathsf{T}}[\nabla \sigma_{\mathsf{H}}]^{\mathsf{T}}\hat{w}_{\mathsf{H}_0}$.
- 13: end for
- 14: **end for**
 - In MfgFL-H (d=H), the model ŵ_{i,d}(n), which is shared with the leader, is equal to the set {ŵ_{H0}(n), ŵ_{H1}(n)} at the UAV *i*. However, the FPK models ŵ_{F0}(n), and ŵ_{F1}(n) are not shared in this implementation.
 - In MfgFL-F (d=F), the model $\hat{w}_{i,d}(n)$, which is shared with the leader, is equal to the set $\{\hat{w}_{F_0}(n), \hat{w}_{F_1}(n)\}$ at the UAV *i*. However, the HJB models $\hat{w}_{H_0}(n)$, and $\hat{w}_{H_1}(n)$ are not shared in this implementation.
 - In MfgFL-HF, (d = HF), the model $\hat{w}_{i,d}(n)$, which is shared with the leader, is equal to the set $\{\hat{w}_{F_0}(n), \hat{w}_{F_1}(n), \hat{w}_{H_0}(n), \hat{w}_{H_1}(n)\}$ at the UAV *i*. In other words, the complete model of MFG is shared with the leader.

It should be noted that the model parameters in different implementations of MfgFL are only shared between the leader and swam of UAVs, and only the models are transmitted rather than raw data samples. This saves significant communication energy as we will see in Section VI. At each iteration n in Algorithm 3, the leader obtains the average model $\hat{w}_{h,d}$ after collecting models $\hat{w}_{i,d}(kn_0)$ from N_h UAVs of the swarm as

$$\hat{w}_{\mathbf{h},d}(k) \leftarrow \frac{1}{N_h} \sum_{i \in \mathcal{N}_h} \hat{w}_{i,d}(kn_0), \tag{62}$$

However, after the average model is calculated at the leader and broadcasted to the UAVs, the updates are done locally by the set local samples S_i at each UAV i, which is the set of local states sampled from the mission starting time to the current time t = ndt, i.e., $S_i = \{s_i(j) | j = 0, ..., n\}$. This procedure is repeated until all the UAVs reach the destination at repetition T_0 .

One main benefit of this procedure is that the UAVs do not need to exchange their local states to update the HJB-FPK models, and communication cost is reduced compared to Mfg and Hjb control methods. An approximate communications payload up to time t = ndt for the MfgFL methods is $N \times \frac{n}{n_0} \times L(\hat{w}_{i,d}) \times b$, where $L(\hat{w}_{i,d})$ is the size of $\hat{w}_{i,d}$ and b is the resolution in bits, and for Mfg and Hjb methods, it is $N \times n \times L(s_i) \times b \times N_s$, where $L(s_i)$ is the size of states and N_s is the number of samples at each time interval dt. Then, the Mfg control requires less N_s than Hjb which corresponds to smaller communications cost of Mfg control, and the MfgFL requires smaller payload than Mfg control because $N_s \gg 1$ and $n_0 \gg 1$. This result will be evaluated in Section VI. In addition, it is not always safe to share the state information, and privacy is also preserved for the UAVs with MfgFL method.

In addition to reducing communication costs and increasing the privacy of the UAVs, the MfgFL method can provide other benefits as well such ensuring stability conditions for MFG framework and increasing training speed. One major condition for the MFG based approach is that the UAVs are indistinguishable. It means that the UAVs should have the same action rule, and hence, it is reasonable that they are trained by big enough samples. Nonetheless, due to energy/bandwidth limitations, it is not possible to provide this huge samples for model training. From this viewpoint, FL-based approaches can increase the model similarity among the UAVs and make them indistinguishable by efficiently using their samples for training.

Another benefit of using the MfgFL approach is the increased training speed of the models at UAVs. This is closely related to the communication cost of the algorithm, since utilizing model averaging means that the algorithm benefit from the various sample of UAVs in a shorter time span. Therefore, it is safe to say that it can provide higher model training speed. However, the performance of MfgFL is explored in next section.

VI. NUMERICAL RESULTS

In this section, we numerically validate the effectiveness of the proposed algorithm MfgFL-HF compared to the baseline methods Hjb control, MfgFL-F, and MfgFL-H, in terms of travel time, motion energy, collision avoidance, and communications cost. Throughout the simulations, we consider N UAVs controlled in a two-dimensional plane at the fixed altitude of h = 40m. Initially, the UAVs are equally separated with the distance $\sqrt{2m}$ each other, and located at a source, which is a square region centered at (150, 100)m in a 2-dimensional plane (see Fig. 2-a). Each UAV aims to reach the destination at the origin, under the wind dynamics described by $V_o = \sigma_{wind}I$ and $v_o = (1, -1)m/s$ (see Sec. II).

Following [4], [10] single hidden layer models (32) and (33) are considered for HJB model, where each hidden node's activation function, i.e., $\sigma_{\text{H},j}(s_i(t))$ for $j = 1, \dots, M_{\text{H}}$, corresponds to each scalar term in a polynomial expansion. The polynomial for $\sigma_{\text{H}}(s_i(t))$ is heuristically chosen as: $(1+x_i(t)+v_{x,i}(t))^6+(1+y_i(t)+v_{y,i}(t))^6$, where $r_i(t) = [x_i(t), y_i(t)]^{\mathsf{T}}$ and $v_i(t) = [v_{x,i}(t), v_{y,i}(t)]^{\mathsf{T}}$, thus the model size for HJB model is $M_{\text{H}} = 54$.

For the MFG based methods, the same neural network structure described above is considered to approximate HJB model. In a similar way, single hidden layer models (51) and (52) are considered for FPK model, where each hidden node's activation function, i.e., $\sigma_{\text{F},j}(s_i(t))$ for $j = 1, \dots, M_{\text{F}}$, corresponds to each scalar term in a polynomial expansion. The polynomial for $\sigma_{\text{F}}(s_i(t))$ is heuristically chosen as: $(1 + x_i(t) + v_{x,i}(t))^6 + (1 + y_i(t) + v_{y,i}(t))^6$. Thus the model size for FPK model is $M_{\text{F}} = 69$.

Unless otherwise stated, the default simulation parameters are: $P_o = 20$ dBm, $W_o = 2$ MHz, $\sigma_n = -110$ dBm/Hz, $\alpha = 0$, $\chi = 1.347$, $\xi = 6.649$; $\sigma_{wind} = 0.1$; $N_h = 0.8N$; $r_{coll} = 0.1$ m, $r_C = \sqrt{2}/2$ m; $c_0 = 0.1$, $c_1 = c_2 = 0.015$, $c_3 = 0.005$, $\mu_H = \mu_F = 0.01$, $c_H = 0.5$, $n_0 = 100$ for MfgFL-H and MfgFL-F, $n_0 = 200$ for MfgFL-HF, and dt = 0.1s for the purpose of discretizing time in simulations. In addition the physical characteristics of UAVs are $\lambda_0 = 0.0049$, $\lambda_1 = 0.0887$ and $\lambda_2 = 0.0092$, $\omega_{tip} = 15$ m/s, and $\chi_o = 1.6120$ m/s.

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication. Citation information: DOI 10.1109/TCOMM.2020.3017281, IEEE Transactions on Communications



Fig. 2. Trajectory snapshots (left, 4 subplots for each control method) of 25 UAVs under (a) Hjb: HJB learning control with the communication range d = 100m, (b) MfgFL-H: MFG learning control with HJB model averaging, (c) MfgFL-F: MFG learning control with FPK model averaging, and (d) MfgFL-HF: MFG learning control with both HJB and FPK model averaging. During the travel time $t = 0 \sim 125$ s, MfgFL-HF shows the best flocking behavior and the most stable HJB model parameters $w_{1,H}$ (rightmost subplot for each control method) of a randomly selected reference UAV u_1 . Consequently, MfgFL-HF yields no collision during its entire travel, in sharp contrast to the others.

Fig. 2 shows the trajectories of N = 25 UAVs under Hjb, MfgFL-H, MfgFL-F, and MfgFL-HF control methods. With Hjb control, all the UAVs should communicate instantaneous states with each other, and use the received states to update their local HJB model. Therefore, Hjb control is extremely costly to be implement in real-time. However, for comparison purposes, it is assumed that the UAVs communicate their states at each time step to calculate the instantaneous interaction term, but the processing at each UAV is limited to one update of (43) and (44) per time step. This results in a fair comparison with FL-based methods, as they are also limited to one update of (43) and (44) per time step.

In all the methods, at first, the untrained UAVs follow the average wind direction while they train the models until the models are trained to the extend that their output commands turn the UAVs towards the destination. Then, the differences among algorithms in terms of collision, model weights, and interaction terms become observable from the trajectory and model weight plots, as explained in the following.

Collision occurrences is shown by star marks in the trajectories. It can be seen that in the proposed MfgFL-HF method, no collision has happened thanks to more sample utilization for HJB and FPK model training by adopting FL averaging for both models. Unlike MfgFL-HF, only one of the HJB or FPK models in MfgFL-H and MfgFL-F methods is trained with enough samples by utilizing FL method. The less-trained model results in more collisions of MfgFL-H and MfgFL-F methods as seen in the trajectory plots Fig. 2-b and Fig. 2c. In Hib method, although enough samples are provided, the UAVs can not use them to train the model in real-time due to limited processing power of the UAVs. Therefore, the models are not trained with enough samples, and a few collisions occur on the path to the destination as seen in the trajectory plots Fig. 2-a. These training behaviors can also be seen in HJB model parameters on the most right side of the Fig. 2,

where in comparison to the other control methods, the model parameters in MfgFL-HF are less divergent after a period of time.

Fig. 2 shows the interaction term ϕ_G for each UAV using the color map on the trajectories. The bluer trajectories of MfgFL-HF method compared to other methods indicates lower interaction term values and better alignment of UAVs on the path to the destination. The reason is better training of the models in the proposed method as explained above. One main benefit of flocking of the UAVs instead of traveling individually or in different clusters is that it results in better communication channels among the UAV due to shorter distances, which can help in better model training and control. Further features of the proposed MfgFL-HF method corresponding to Fig. 2 is explained below using Fig. 3 and Fig. 4.

Fig. 3 represents the motion energy, communications payload, velocity alignment, number of collision risks, speed, and travel distance of the UAVs corresponding to the scenario and methods in Fig. 2. Fig. 3-a represents the average motion energy and its variance among the UAVs. The proposed method MfgFL-HF consumes at least 16% less energy than the other methods, and requires at least 4 times less communication costs than Hjb method (see Fig. 3-b) at the cost of 10%and 6% more average travel time T_{avg} compared to MfgFL-F and MfgFL-FH, respectively. The reason for less energy consumption of MfgFL-HF is that the UAVs can travel in a flock with smaller speed (Fig. 3-e) and smaller interaction term on the trajectory as we observed in Fig. 2. Furthermore, the reason for less communication costs for MfgFL-HF, MfgFL-F, and MfgFL-H methods is due to adopting the FL method. In FL-based methods, at every $n_0 = 100$ time steps, 80% of the UAVs transmit their models to the leader and the leader broadcasts it to all the UAVs. However, in Hjb method, all 25 UAVs broadcasts their states to all the neighbor UAVs at each

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication. Citation information: DOI 10.1109/TCOMM.2020.3017281, IEEE Transactions on Communications



Fig. 3. The comparison of different methods in terms of (a) motion energy, (b) communications payload (c) velocity alignment, (d) number of collision risks, (e) mean/max/min speed, and (f) traveled distance.

time step.

Despite the disadvantage of more travel time for the MfgFL-HF method in this scenario, it demonstrates better velocity alignment and collision avoidance properties than the other defined FL-based methods as shown in Fig. 3-c and Fig. 3d. As explained in definition of metrics $\phi_A(t)$ in (10) and $\phi_C(t)$ (11), their smaller values corresponds to better flocking behavior and lower probability of collision occurrence, respectively. Clearly, the cumulative value of $\phi_A(t)$ at time $T_{avg} = 175$ s for MfgFL-HF method is at least 7% less than the other methods, which means better velocity alignment of MfgFL-HF. Furthermore, the cumulative value of $\phi_C(t)$ at time $T_{avg} = 175$ s for MfgFL-HF is at least 8% less than the other methods, which means lower risk of collision occurrence of the proposed method. This complies with the training discussion above for Fig. 2.

Furthermore, the average, maximum and minimum speed of the UAVs are shown in Fig. 3-e and Fig. 3-f. The obtained maximum speed for the mentioned algorithms are: 12.9m/s for MfgFL-HF, 16.0m/s for MfgFL-F, 14.3m/s for MfgFL-H, 16.2m/s for Hjb. The lower maximum speed of MfgFL-HF is because of less interaction among the UAVs which is a direct result of better flocking. The better flocking is obtained at the cost of increasing the travel distance about 10% compared to the Hjb as shown in Fig. 3-f. This is also consequently because of better model training of MfgFL-HF by utilizing FL method.

Fig. 4 represents the absolute values of approximation errors of HJB in (41) and FPK in (56) equations corresponding to the scenario and methods in Fig. 2. It is noticeable that the values of approximation errors of HJB and FPK, despite not being too small, are acceptably small. This is in compliance with the analysis that the model weights are UUB. The approximation errors of HJB, i.e., $e_{\rm H}$ in (41), depends on the error value of HJB model weights which is proved to be UUB in Proposition 1, and The approximation errors of FPK, i.e., $e_{\rm F}$ in (56), depends on the error value of FPK model weights which is proved to be UUB in Proposition 2. Then, when the error value of model weights is below a threshold, the corresponding absolute values of approximate error of HJB in (41) and/or FPK in (56) will be bounded. This can be seen in Fig. 4-a



Fig. 4. Approximation error values of (a) HJB model, (b) FPK model, are small during the training for all methods.



Fig. 5. Performance of different methods vs number of UAVs in terms of (a) motion energy, (b) travel time (c) velocity alignment, and (d) number of collision risks.

and Fig. 4-b that the corresponding absolute error values are below 1.5 and 0.02, respectively.

Fig. 5 illustrates the performance of different methods versus N number of UAVs. Clearly, MfgFL-HF requires less motion energy for $N = 16, \ldots, 64$, and its performance in terms of travel time $T \leq T_{\text{max}}$, velocity alignment $\phi_A(T_{\text{avg}})$, and number of collision risks $\phi_{\rm C}(T_{\rm avg})$, improves as the number of UAVs increases. This is because, for higher N, more samples can be provided for both HJB and FPK models in MfgFL-HF which results in better training of both models. However, for the other two FL base methods, i.e., MfgFL-H and MfgFL-F, provided samples due to averaging improves only one of the HJB or FPK models, and the other corresponding model still remains less trained. Therefore, the coupled HJB-FPK equation in these two methods still is not well trained and non of MfgFL-H and MfgFL-F can benefit much when number of UAVs increases. Regarding Hjb, increasing the number of UAVs does not improve the performance much since it does not utilize the more provided samples for training the model due to the processing power limitations of the UAVs.

Fig. 6 shows the impact of model update period n_0 on the performance criteria for the FL-based methods. For n_0 larger than 100, the MfgFL-HF method consumes less energy than the other FL-based methods to complete the travel (see Fig. 6-a), while its travel time is only more than MfgFL-H for most of the choices of n_0 (see Fig. 6-b). Regarding the velocity alignment $\phi_A(T_{avg})$ and number of collision risks $\phi_{\rm C}(T_{\rm avg})$, for n_0 in the interval $n_0 \in \{100, \cdots, 400\}$, MfgFL-HF method has lower velocity alignment and number of collision risks than other FL-based methods (see Fig. 6-c and Fig. 6-d). Additionally, there is an acceptable trade-off among the performance criteria for MfgFL-HF method in this interval. This is due to the fact that, for small values of n_0 , e.g., 50, fewer UAVs can successfully transmit their data to the leader because of communication costs such as limited transmission power of UAVs. On the other hand, for very high values of n_0 , e.g., 500, the algorithms cannot benefit from adopting the FL method in real-time application, because when n_0 increases the models at UAVs rely mostly on local samples for larger

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication. Citation information: DOI 10.1109/TCOMM.2020.3017281, IEEE Transactions on Communications



Fig. 6. Performance of different FL-based methods vs model update period in terms of (a) motion energy, (b) travel time (c) velocity alignment, and (d) number of collision risks.



Fig. 7. Performance of different FL-based methods vs wind variance in terms of (a) motion energy, (b) travel time (c) velocity alignment, and (d) number of collision risks.

amount of time and become less trained.

Fig. 7 represents the effect of wind perturbations on the online control of UAVs in the scenario of Fig. 2. Here, the UAVs are set to move under different wind perturbation variances $\sigma_{wind} = 0.1, \dots, 0.5$ and the various comparison criteria are calculated for the algorithms. From Fig. 7-a, the proposed algorithm MfgFL-HF consumes less control energy than the other baselines, and the energy consumption variance is increasing with the wind variance. This is because the agents with MfgFL-HF keep a distance away from each other to avoid collision. Also, Fig. 7-b to Fig. 7-d show smaller travel time of MfgFL-HF than MfgFL-F and Hjb, and better collision avoidance of MfgFL-HF than all other mentioned methods. Overall, these figures emphasize that the algorithm MfgFL-HF is more robust against wind dynamics, thanks to better learning capability of the proposed method.

All the previous figures show the performance of untrained UAVs in a windy environment. However, there is still an important concern remaining: How the proposed method behaves in comparison to the offline methods in windy environments? This comparison is shown in Fig. 8. The term offline method here means that the UAVs are separated enough at the starting point and they are programmed to follow some pre-defined actions to reach the destination when there is no randomness in the environment, i.e., $\sigma_{wind} = 0$ as in [68] and without any collaboration among UAVs. Fig. 8(a) shows the optimal shortest path which the UAVs can follow in an imaginary perfect environment without any collision occurrences. Nevertheless, the offline method fails to reach the destination with no collision in the presence of random wind dynamics, i.e., $\sigma_{wind} = 0.1$, as shown in Fig. 8(b). Moreover, Fig. 8(c) shows



12

Fig. 8. The comparison of online training of UAV by MfgFL-HF algorithm with the UAVs pre-programmed for deterministic environment.

that the trajectory of the trained UAVs with the proposed method in the windy environment with $\sigma_{wind} = 1.5$ is higher than the training environment with $\sigma_{wind} = 0.1$. The proposed method MfgFL-HF is much more robust to random wind perturbations and can reach the destination with no collision on the path towards the destination, while there is no such collision avoidance guarantee in the offline method.

VII. CONCLUSION

In this paper, a novel path planning approach is proposed for a population of UAVs being effected by random wind perturbations in the environment. To this end, the objective is to minimize the transmission time, motion energy, and the interactions among the UAVs. First, the MFG framework is applied in order to reduce the high amount of communications required to control a massive number of UAVs. Next, a function approximator based on neural networks is proposed to approximate the solution of the HJB and FPK equations. The Lyapunov stability analysis for MFG learning is provided to show that the approximate solution for HJB and FPK equations are bounded. Then, on the bases of these assumptions and analyses, an FL-based MFG learning method named MfgFL-HF is proposed to use the samples of UAVs more efficiently for the purpose of training the model weights of neural networks at UAVs. The numerical results confirm the stability of the proposed method and show that it can be used to control a massive UAV population in a windy environment efficiently.

APPENDIX A PROOF OF PROPOSITION 1.

This proof is based on methodology of [67], but with a few differences in system model and the update algorithms. The candidate Lyapunov function is chosen as

$$L(t) = \frac{1}{2\mu_{\rm H}} \tilde{w}_{\rm H_0}^{\rm T} \tilde{w}_{\rm H_0} + \frac{1}{2\mu_{\rm H}} \tilde{w}_{\rm H_1}^{\rm T} \tilde{w}_{\rm H_1} + c_{\rm H} L_s \mathbb{1}_{\|s_i\| \ge s_{\rm dest}}.$$
 (A.1)

Then, the corresponding derivative function is

$$\dot{L}(t) = \frac{1}{\mu_{\rm H}} \tilde{w}_{\rm H_0}^{\sf T} \dot{\bar{w}}_{\rm H_0} + \frac{1}{\mu_{\rm H}} \tilde{w}_{\rm H_1}^{\sf T} \dot{\bar{w}}_{\rm H_1} + c_{\rm H} [\nabla L_s]^{\sf T} \dot{s} \mathbb{1}_{\|s_i\| \ge s_{\rm dest}}, \qquad (A.2)$$

which can be rewritten in the following form

$$\begin{split} L(t) &= \tilde{w}_{H_{0}}^{\mathsf{T}}[(\nabla_{\hat{w}_{H_{0}}}e_{H})e_{H} + c_{H}\nabla_{\hat{w}_{H_{0}}}R_{i}] \\ &+ \tilde{w}_{H_{1}}^{\mathsf{T}}[(\nabla_{\hat{w}_{H_{1}}}e_{H})e_{H}] + c_{H}[\nabla L_{s}]^{\mathsf{T}}\dot{s}\mathbb{1}_{\|s_{i}\|\geq s_{\text{dest}}} \\ &= [\tilde{w}_{H_{0}}^{\mathsf{T}}(\nabla_{\hat{w}_{H_{0}}}e_{H}) + \tilde{w}_{H_{1}}^{\mathsf{T}}(\nabla_{\hat{w}_{H_{1}}}e_{H})]e_{H} \\ &+ c_{H}\tilde{w}_{H_{0}}^{\mathsf{T}}\nabla_{\hat{w}_{H_{0}}}R_{i} + c_{H}[\nabla L_{s}]^{\mathsf{T}}\dot{s}\mathbb{1}_{\|s_{i}\|\geq s_{\text{dest}}} \\ &= -[\tilde{w}_{H_{0}}^{\mathsf{T}}b_{H} + \tilde{w}_{H_{0}}^{\mathsf{T}}R_{H}\tilde{w}_{H_{0}} + \tilde{w}_{H_{1}}^{\mathsf{T}}\sigma_{H}] \\ &\times [\tilde{w}_{H_{0}}^{\mathsf{T}}b_{H} + \frac{1}{2}\tilde{w}_{H_{0}}^{\mathsf{T}}R_{H}\tilde{w}_{H_{0}} + \tilde{w}_{H_{1}}^{\mathsf{T}}\sigma_{H} + \varepsilon_{H}] \\ &+ c_{H}\tilde{w}_{H_{0}}^{\mathsf{T}}\nabla_{\hat{w}_{H_{0}}}R_{i} + c_{H}[\nabla L_{s}]^{\mathsf{T}}\dot{s}\mathbb{1}_{\|s_{i}\|\geq s_{\text{dest}}} \end{split} \tag{A.3}$$

where b_{H} and matrix R_{H} are defined as

$$\begin{split} b_{\mathsf{H}} &= [\nabla \sigma_{\mathsf{H}}] (f - \frac{1}{2c_3} B B^{\mathsf{T}} [\nabla \sigma_{\mathsf{H}}]^{\mathsf{T}} w_{\mathsf{H}_0} - \frac{1}{2c_3} B B^{\mathsf{T}} [\nabla \varepsilon_{\mathsf{H}_0}] \\ &+ \frac{1}{2c_3} [\nabla \sigma_{\mathsf{H}}] B B^{\mathsf{T}} [\nabla \varepsilon_{\mathsf{H}_0}]) + \frac{1}{2} \sum_{k=1}^{N} \operatorname{tr} (G G^T [\Delta \sigma_{\mathsf{H}}^{[k]}]) \mathbf{e}_k \\ &= [\nabla \sigma_{\mathsf{H}}] \bar{s} + \frac{1}{2} \sum_{k=1}^{N} \operatorname{tr} (G G^T [\Delta \sigma_{\mathsf{H}}^{[k]}]) \mathbf{e}_k + \frac{1}{2c_3} [\nabla \sigma_{\mathsf{H}}] B B^{\mathsf{T}} [\nabla \varepsilon_{\mathsf{H}_0}], \end{split}$$
(A.4)
$$R_{\mathsf{H}} &= \frac{1}{2c_3} [\nabla \sigma_{\mathsf{H}}] B B^{\mathsf{T}} [\nabla \sigma_{\mathsf{H}}]^{\mathsf{T}}, \tag{A.5}$$

and \overline{s} is the dynamics of nominal system defined as

$$\bar{s} = f - \frac{1}{2c_3} B B^{\mathsf{T}} [\nabla \sigma_{\mathsf{H}}]^{\mathsf{T}} w_{\mathsf{H}_0} - \frac{1}{2c_3} B B^{\mathsf{T}} [\nabla \varepsilon_{\mathsf{H}_0}].$$
(A.6)

Using the following relation as

$$ab = \frac{1}{2} \left(-(ha - \frac{b}{h})^2 + h^2 a^2 + \frac{b^2}{h^2} \right), \tag{A.7}$$

for scalars a and b, we obtain the upper-bounds for the terms on the right hand side of (A.3) as

$$-\frac{3}{2}(\tilde{w}_{\mathsf{H}_{0}}^{\mathsf{T}}b_{\mathsf{H}})(\tilde{w}_{\mathsf{H}_{0}}^{\mathsf{T}}R_{\mathsf{H}}\tilde{w}_{\mathsf{H}_{0}}) \leq \frac{3}{4}h_{1}^{2}\lambda_{1M}^{2}\|\tilde{w}_{\mathsf{H}_{0}}\|^{2} + \frac{3}{4h_{1}^{2}}\lambda_{2M}^{2}\|\tilde{w}_{\mathsf{H}_{0}}\|^{4}, \quad (A.8)$$

$$-(\tilde{w}_{\mathsf{H}_{0}}^{\mathsf{T}}b_{\mathsf{H}})(\varepsilon_{\mathsf{H}}) \leq \frac{1}{2}h_{2}^{2}\lambda_{1M}^{2}\|\tilde{w}_{\mathsf{H}_{0}}\|^{2} + \frac{1}{2h_{2}^{2}}\lambda_{3M}^{2}, \tag{A.9}$$

$$-(\tilde{w}_{\mathsf{H}_{0}}^{\mathsf{T}}R_{\mathsf{H}}\tilde{w}_{\mathsf{H}_{0}})(\varepsilon_{\mathsf{H}}) \leq \frac{1}{2}h_{3}^{2}\lambda_{2M}^{2}\|\tilde{w}_{\mathsf{H}_{0}}\|^{4} + \frac{1}{2h_{3}^{2}}\lambda_{3M}^{2}, \tag{A.10}$$

$$-\frac{3}{2}(\tilde{w}_{\mathsf{H}_{1}}^{\mathsf{T}}\sigma_{\mathsf{H}})(\tilde{w}_{\mathsf{H}_{0}}^{\mathsf{T}}R_{\mathsf{H}}\tilde{w}_{\mathsf{H}_{0}}) \leq \frac{3}{4}h_{4}^{2}\lambda_{4M}^{2}\|\tilde{w}_{\mathsf{H}_{1}}\|^{2} + \frac{3}{4h_{4}^{2}}\lambda_{2M}^{2}\|\tilde{w}_{\mathsf{H}_{0}}\|^{4}, (A.11)$$

$$-(\tilde{w}_{\mathsf{H}_{1}}^{\mathsf{T}}\sigma_{\mathsf{H}})(\varepsilon_{\mathsf{H}}) \leq \frac{1}{2}h_{5}^{2}\lambda_{4M}^{2}\|\tilde{w}_{\mathsf{H}_{1}}\|^{2} + \frac{1}{2h_{5}^{2}}\lambda_{3M}^{2}, \qquad (A.12)$$

$$-2(\tilde{w}_{\mathsf{H}_{0}}^{\mathsf{T}}b_{\mathsf{H}})(\tilde{w}_{\mathsf{H}_{1}}^{\mathsf{T}}\sigma_{\mathsf{H}}) \le h_{6}^{2}\lambda_{1M}^{2}\|\tilde{w}_{\mathsf{H}_{0}}\|^{2} + \frac{1}{h_{6}^{2}}\lambda_{4M}^{2}\|\tilde{w}_{\mathsf{H}_{1}}\|^{2}, \quad (A.13)$$

$$-(\tilde{w}_{\mathsf{H}_{0}}^{\mathsf{T}}b_{\mathsf{H}})^{2} \leq -\lambda_{1m}^{2} \|\tilde{w}_{\mathsf{H}_{0}}\|^{2}, \tag{A.14}$$

$$\left(\tilde{w}_{\mathsf{H}_{0}}^{\mathsf{T}}R_{\mathsf{H}}\tilde{w}_{\mathsf{H}_{0}}\right)^{2} \leq -\frac{1}{2}\lambda_{2m}^{2}\|\tilde{w}_{\mathsf{H}_{0}}\|^{4},\tag{A.15}$$

$$-(\tilde{w}_{\mathsf{H}_{1}}^{\mathsf{T}}\sigma_{\mathsf{H}})^{2} \leq -\lambda_{4m}^{2} \|\tilde{w}_{\mathsf{H}_{1}}\|^{2}, \tag{A.16}$$

where we assumed that

$$\lambda_{1m} \le \|b_{\mathsf{H}}\| \le \lambda_{1M},\tag{A.17}$$

$$\lambda_{2m} \le \|R_{\mathsf{H}}\| \le \lambda_{2M},\tag{A.18}$$

$$\begin{aligned} \|\varepsilon_{\mathsf{H}}\| &\leq \lambda_{3M}, \qquad (A.19)\\ \lambda_{4m} &\leq \|\sigma_{\mathsf{H}}\| \leq \lambda_{4M}. \end{aligned}$$
(A.20)

Therefore, the derivative of Lyapunov function is upperbounded as

$$\begin{split} \dot{L}(t) &\leq -\lambda_0 \|\tilde{w}_{\mathsf{H}_0}\|^4 + \lambda_1 \|\tilde{w}_{\mathsf{H}_0}\|^2 + \lambda_2^2 - \lambda_3 \|\tilde{w}_{\mathsf{H}_1}\|^2 \\ &+ c_{\mathsf{H}} \tilde{w}_{\mathsf{H}_0}^{\mathsf{T}} \nabla_{\hat{w}_{\mathsf{H}_0}} R_i + c_{\mathsf{H}} [\nabla L_s]^{\mathsf{T}} \dot{s} \mathbb{1}_{\|s_i\| \geq s_{\mathsf{dest}}}, \end{split}$$
(A.21)

where λ_0 , λ_1 , λ_2 , and λ_3 are defined as

$$\begin{split} \lambda_0 &= -\frac{3}{4h_1^2}\lambda_{2M}^2 - \frac{1}{2}h_3^2\lambda_{2M}^2 - \frac{3}{4h_4^2}\lambda_{2M}^2 + \frac{1}{2}\lambda_{2m}^2, \\ \lambda_1 &= \frac{3}{4}h_1^2\lambda_{1M}^2 + \frac{1}{2}h_2^2\lambda_{1M}^2 + h_6^2\lambda_{1M}^2 - \lambda_{1m}^2, \\ \lambda_3 &= -\frac{3}{4}h_4^2\lambda_{4M}^2 - \frac{1}{2}h_5^2\lambda_{4M}^2 - \frac{1}{h_6^2}\lambda_{4M}^2 + \lambda_{4m}^2, \\ \lambda_2^2 &= \frac{1}{2h_2^2}\lambda_{3M}^2 + \frac{1}{2h_3^2}\lambda_{3M}^2 + \frac{1}{2h_5^2}\lambda_{3M}^2. \end{split}$$
(A.22)

Depending on the state of the UAV, three cases can occur in (A.21) as

Case 1: $\mathbb{1}_{||s_i|| \ge s_{dest}} = 0$. With this condition, we can conclude that the UAVs are in destination, and we focus only on the weights of the models

$$\dot{L}(t) \le -\lambda_0 \|\tilde{w}_{\mathsf{H}_0}\|^4 + \lambda_1 \|\tilde{w}_{\mathsf{H}_0}\|^2 + \lambda_2^2 - \lambda_3 \|\tilde{w}_{\mathsf{H}_1}\|^2.$$
(A.23)

Then, when the following conditions hold, i.e.,

$$\|\tilde{w}_{\mathsf{H}_0}\| \ge \sqrt{\frac{\lambda_1 + \sqrt{\lambda_1^2 + 4\lambda_0\lambda_2^2}}{2\lambda_0}} \triangleq \omega_{0,1}, \tag{A.24}$$

$$\tilde{\omega}_{\mathsf{H}_{1}} \| \ge \sqrt{\frac{4\lambda_{2}^{2}\lambda_{0} + \lambda_{1}^{2}}{4\lambda_{3}\lambda_{0}}} \stackrel{\Delta}{=} \omega_{1,1}, \tag{A.25}$$

the stability condition $\dot{L}(t) < 0$ is satisfied.

case 2: $\mathbb{1}_{\|s_i\| \ge s_{\text{dest}}} = 1$ and $\dot{L}_s \le 0$. In this case, the regulizer term is inactive, and the upper-bound for derivative of Lyapunov is reduced to

$$\begin{split} \dot{L}(t) &\leq -\lambda_{0} \|\tilde{w}_{\mathsf{H}_{0}}\|^{4} + \lambda_{1} \|\tilde{w}_{\mathsf{H}_{0}}\|^{2} + \lambda_{2}^{2} - \lambda_{3} \|\tilde{w}_{\mathsf{H}_{1}}\|^{2} + c_{\mathsf{H}} [\nabla L_{s}]^{\mathsf{T}} \dot{s}_{1} \|_{s_{i}} \|_{\geq s_{\text{dest}}} \\ &\leq -\lambda_{0} \|\tilde{w}_{\mathsf{H}_{0}}\|^{4} + \lambda_{1} \|\tilde{w}_{\mathsf{H}_{0}}\|^{2} + \lambda_{2}^{2} - \lambda_{3} \|\tilde{w}_{\mathsf{H}_{1}}\|^{2} + c_{\mathsf{H}} \lambda_{4} \|\nabla L_{s}\| \tag{A.26}$$

where λ_4 is a number such that $0 < \lambda_4 ||\nabla L_s|| \le -[\nabla L_s]^{\mathsf{T}} \dot{s}$. Therefore, when the following inequalities as

$$\|\tilde{w}_{\mathsf{H}_0}\| \ge \sqrt{\frac{\lambda_1 + \sqrt{\lambda_1^2 + 4\lambda_0\lambda_2^2}}{2\lambda_0}} \stackrel{\Delta}{=} \omega_{0,2} \tag{A.27}$$

$$\|\tilde{w}_{\mathsf{H}_1}\| \ge \sqrt{\frac{4\lambda_2^2\lambda_0 + \lambda_1^2}{4\lambda_3\lambda_0}} \triangleq \omega_{1,2}$$
(A.28)

$$\|\nabla L_s(s_i(t))\| \ge \frac{4\lambda_2^2\lambda_0 + \lambda_1^2}{4\lambda_4\lambda_0} \stackrel{\Delta}{=} \gamma_2 \tag{A.29}$$

occur, the stability condition $\dot{L}(t) < 0$ holds.

case 3: $\mathbb{1}_{||s_i|| \ge s_{dest}} = 1$ and $\dot{L}_s \ge 0$. In this case, we find the upper-bound for the $\dot{L}(t)$ as

$$\begin{split} \dot{L}(t) &\leq -\lambda_{0} \|\tilde{w}_{\mathsf{H}_{0}}\|^{4} + \lambda_{1} \|\tilde{w}_{\mathsf{H}_{0}}\|^{2} + \lambda_{2}^{2} - \lambda_{3} \|\tilde{w}_{\mathsf{H}_{1}}\|^{2} \\ &+ c_{\mathsf{H}} \tilde{w}_{\mathsf{H}_{0}}^{\mathsf{T}} \nabla_{\hat{w}_{\mathsf{H}_{0}}} R_{i} + c_{\mathsf{H}} [\nabla L_{s}]^{\mathsf{T}} \dot{s}^{1} \|_{\|s_{i}\|| \geq s_{\mathsf{dest}}} \\ &= -\lambda_{0} \|\tilde{w}_{\mathsf{H}_{0}}\|^{4} + \lambda_{1} \|\tilde{w}_{\mathsf{H}_{0}}\|^{2} + \lambda_{2}^{2} - \lambda_{3} \|\tilde{w}_{\mathsf{H}_{1}}\|^{2} \\ &+ c_{\mathsf{H}} \tilde{w}_{\mathsf{H}_{0}}^{\mathsf{T}} \nabla_{\hat{w}_{\mathsf{H}_{0}}} [[\nabla L_{s}]^{\mathsf{T}} \dot{s}] + c_{\mathsf{H}} [\nabla L_{s}]^{\mathsf{T}} \dot{s} \\ \stackrel{(1)}{=} -\lambda_{0} \|\tilde{w}_{\mathsf{H}_{0}}\|^{4} + \lambda_{1} \|\tilde{w}_{\mathsf{H}_{0}}\|^{2} + \lambda_{2}^{2} - \lambda_{3} \|\tilde{w}_{\mathsf{H}_{1}}\|^{2} \\ &+ c_{\mathsf{H}} [\nabla L_{s}]^{\mathsf{T}} \ddot{s} + \frac{c_{\mathsf{H}}}{2c_{3}} [\nabla L_{s}]^{\mathsf{T}} B B^{\mathsf{T}} [\nabla \varepsilon_{\mathsf{H}_{0}}] \\ \stackrel{(2)}{\leq} -\lambda_{0} \|\tilde{w}_{\mathsf{H}_{0}}\|^{4} + \lambda_{1} \|\tilde{w}_{\mathsf{H}_{0}}\|^{2} + \lambda_{2}^{2} - \lambda_{3} \|\tilde{w}_{\mathsf{H}_{1}}\|^{2} \\ &- c_{\mathsf{H}} \lambda_{5m} \|\nabla L_{s}\|^{2} + \frac{c_{\mathsf{H}}}{2c_{3}} \lambda_{6M} \|\nabla L_{s}\|, \end{split}$$
(A.30)

where equality (1) is obtained by the calculations as

$$c_{\mathsf{H}}\tilde{w}_{\mathsf{H}_{0}}^{\mathsf{T}}\nabla_{\hat{w}_{\mathsf{H}_{0}}}[[\nabla L_{s}]^{\mathsf{T}}\dot{s}] = \frac{1}{2c_{3}}c_{\mathsf{H}}[\tilde{w}_{\mathsf{H}_{0}}]^{\mathsf{T}}[\nabla\sigma_{\mathsf{H}}]BB^{\mathsf{T}}[\nabla L_{s}], \quad (A.31)$$

$$c_{\mathsf{H}}[\nabla L_s]^{\mathsf{T}}\dot{s} = c_{\mathsf{H}}[\nabla L_s]^{\mathsf{T}}(f - \frac{1}{2c_3}BB^{\mathsf{T}}[\nabla\sigma_{\mathsf{H}}]^{\mathsf{T}}\hat{w}_{\mathsf{H}_0}), \tag{A.32}$$

and inequality (2) is based on the following assumptions,

$$\begin{split} [\nabla L_s]^{\mathsf{T}} \bar{s} &= -[\nabla L_s]^{\mathsf{T}} A[\nabla L_s], \\ &\leq -\lambda_{5m} \|\nabla L_s\|^2, \end{split} \tag{A.33}$$

$$BB^{\mathsf{T}}[\nabla \varepsilon_{\mathsf{H}_0}] \le \lambda_{6M}. \tag{A.35}$$

where λ_{5m} is the minimum eigenvalue of matrix A.

Therefore, when the following conditions occur, i.e.,

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication. Citation information: DOI 10.1109/TCOMM.2020.3017281, IEEE Transactions on Communications

14

$$\|\tilde{w}_{\mathsf{H}_{0}}\| \geq \sqrt{\frac{\lambda_{1} + \sqrt{\lambda_{1}^{2} + 4\lambda_{0}(\lambda_{2}^{2} + \frac{c_{\mathsf{H}}\lambda_{6M}^{2}}{16c_{3}^{2}\lambda_{5m}})}{2\lambda_{0}}} \triangleq \omega_{0,3} \quad (A.36)$$

$$\|\tilde{w}_{\mathsf{H}_1}\| \ge \sqrt{\frac{4(\lambda_2^2 + \frac{c_{\mathsf{H}}\lambda_{6M}^2}{16c_3^2\lambda_{5m}})\lambda_0 + \lambda_1^2}{4\lambda_3\lambda_0}} \triangleq \omega_{1,3} \tag{A.37}$$

$$\|\nabla L_s(s_i(t))\| \ge \frac{\frac{c_{\mathsf{H}}}{2c_3}\lambda_{6M} + \sqrt{\left(\frac{c_{\mathsf{H}}}{2c_3}\lambda_{6M}\right)^2 + 4c_{\mathsf{H}}\lambda_{5m}(\lambda_2^2 + \frac{\lambda_1^2}{4\lambda_0})}}{c_{\mathsf{H}}\lambda_{5m}} \stackrel{\text{(A.38)}}{=}$$

the Lyapunov stability condition holds, i.e., $\dot{L}(t) < 0$.

In summary, when $\|\tilde{w}_{H_0}\| \geq \omega_0 = \max\{\omega_{0,1}, \omega_{0,2}, \omega_{0,3}\}$, or $\|\tilde{w}_{H_1}\| \geq \omega_1 = \max\{\omega_{1,1}, \omega_{1,2}, \omega_{1,3}\}$, or $\|\nabla L_s(s_i(t))\| \geq \max\{\gamma_2, \gamma_3\}$ occurs, then the Lyapunov stability condition holds, i.e., $\dot{L}(t) < 0$. Considering all the cases 1-3, we can conclude that there exist s_{dest} , w_0 , and w_1 at time T such that $\|s(t)\| \leq s_{dest}$, $\|w_{H_0}(t) - \hat{w}_{H_0}(t)\| \leq w_0$, and $\|w_{H_1}(t) - \hat{w}_{H_1}(t)\| \leq w_1$ for all $t \geq T + T'$.

APPENDIX B PROOF OF PROPOSITION 2.

The candidate Lyapunov function is chosen as

$$L(t) = \frac{1}{2\mu_{\mathsf{F}}} \tilde{w}_{\mathsf{F}_0}^{\mathsf{T}} \tilde{w}_{\mathsf{F}_0} + \frac{1}{2\mu_{\mathsf{F}}} \tilde{w}_{\mathsf{F}_1}^{\mathsf{T}} \tilde{w}_{\mathsf{F}_1}.$$
 (B.1)

Then, the derivative of Lyapunov function is obtained as

$$\begin{split} L(t) &= \tilde{w}_{\mathsf{F}_{0}}^{\mathsf{I}}(\nabla_{\hat{w}_{\mathsf{F}_{0}}}e_{\mathsf{F}})e_{\mathsf{F}} + \tilde{w}_{\mathsf{F}_{1}}^{\mathsf{I}}(\nabla_{\hat{w}_{\mathsf{F}_{1}}}e_{\mathsf{F}})e_{\mathsf{F}}, \\ &= -[\tilde{w}_{\mathsf{F}_{0}}^{\mathsf{I}}[\nabla_{\hat{w}_{\mathsf{F}_{0}}}e_{\mathsf{F}}] + \tilde{w}_{\mathsf{F}_{1}}^{\mathsf{I}}[\nabla_{\hat{w}_{\mathsf{F}_{1}}}e_{\mathsf{F}}]], \\ &\times [\tilde{w}_{\mathsf{F}_{0}}^{\mathsf{I}}[\nabla_{\hat{w}_{\mathsf{F}_{0}}}e_{\mathsf{F}}] + \tilde{w}_{\mathsf{F}_{1}}^{\mathsf{I}}[\nabla_{\hat{w}_{\mathsf{F}_{1}}}e_{\mathsf{F}}] + \varepsilon_{\mathsf{F}}], \\ &= -[\tilde{w}_{\mathsf{F}}^{\mathsf{I}}[\nabla_{\hat{w}_{\mathsf{F}}}e_{\mathsf{F}}]][\tilde{w}_{\mathsf{F}}^{\mathsf{I}}[\nabla_{\hat{w}_{\mathsf{F}}}e_{\mathsf{F}}] + \varepsilon_{\mathsf{F}}]. \end{split}$$
(B.2)

Each term of the derivative of Lyapunov function has an upperbound obtained as

$$-(\tilde{w}_{\mathsf{F}}^{\mathsf{T}}\nabla_{\hat{w}_{\mathsf{F}}}e_{\mathsf{F}})^{2} \leq -\lambda_{7m}^{2}\|\tilde{w}_{\mathsf{F}}\|^{2}, \tag{B.3}$$

$$-(\tilde{w}_{\mathsf{F}}^{\mathsf{T}}[\nabla_{\hat{w}_{\mathsf{F}}}e_{\mathsf{F}}])(\varepsilon_{\mathsf{F}}) \le \lambda_{8M}\lambda_{7M}\|\tilde{w}_{\mathsf{F}}^{\mathsf{T}}\|,\tag{B.4}$$

where it is assumed that

$$\begin{split} \lambda_{7m} &\leq \|\nabla_{\hat{w}_{\mathsf{F}}} e_{\mathsf{F}}\| \leq \lambda_{7M}, \\ \|\varepsilon_{\mathsf{F}}\| &\leq \lambda_{8M}. \end{split} \tag{B.5}$$

Then, the derivative of Lyapunov is upper-bounded as

$$\dot{L}(t) \le -\lambda_{7m}^2 \|\tilde{w}_{\mathsf{F}}\|^2 + \lambda_{8M} \lambda_{7M} \|\tilde{w}_{\mathsf{F}}^{\mathsf{T}}\|.$$
(B.7)

Therefore, when the following condition occurs, i.e.,

$$\|\tilde{w}_{\mathsf{F}}^{\mathsf{T}}\| \ge \frac{\lambda_{8M}\lambda_{7M}}{\lambda_{7m}^2},\tag{B.8}$$

the stability condition holds, i.e., $\dot{L}(t) \leq 0$. However, (B.8) means that the model which makes the term $\|\varepsilon_{\mathsf{F}}\|$ small, can increase the stability of FPK learning algorithm.

APPENDIX C PROOF OF PROPOSITION 3.

Here we aim to show the convergence and bias of the FPK learning updates following the proof method of [69]. Let us first define extended vectors as

$$\hat{w}_{\mathsf{F}}(n) = [[\hat{w}_{\mathsf{F}_0}(n)]^{\mathsf{T}} \quad [\hat{w}_{\mathsf{F}_1}(n)]^{\mathsf{T}}]^{\mathsf{T}},$$
(C.1)
$$\tilde{w}_{\mathsf{F}}(n) = [[\tilde{w}_{\mathsf{F}_0}(n)]^{\mathsf{T}} \quad [\tilde{w}_{\mathsf{F}_1}(n)]^{\mathsf{T}}]^{\mathsf{T}},$$
(C.2)

$$\nabla_{\hat{w}_{\mathsf{F}}} e_{\mathsf{F}} = [[\nabla_{\hat{w}_{\mathsf{F}_0}} e_{\mathsf{F}}]^{\mathsf{T}} \qquad [\nabla_{\hat{w}_{\mathsf{F}_1}} e_{\mathsf{F}}]^{\mathsf{T}}]^{\mathsf{T}}. \tag{C.3}$$

Then, the FPK error vector update is obtained as

$$\tilde{w}_{\mathsf{F}}(n+1) = [I - \mu_{\mathsf{F}}[\nabla_{\hat{w}_{\mathsf{F}}} e_{\mathsf{F}}][\nabla_{\hat{w}_{\mathsf{F}}} e_{\mathsf{F}}]^{\mathsf{T}}]\tilde{w}_{\mathsf{F}}(n) - \mu_{\mathsf{F}}[\nabla_{\hat{w}_{\mathsf{F}}} e_{\mathsf{F}}]\varepsilon_{\mathsf{F}}.$$
 (C.4)

A. FPK Mean Stability

Taking the expectation of equation (C.4) yields

$$\mathbb{E}[\tilde{w}_{\mathsf{F}}(n+1)] = [I - \mu_{\mathsf{F}} \mathbb{E}[[\nabla_{\hat{w}_{\mathsf{F}}} e_{\mathsf{F}}] [\nabla_{\hat{w}_{\mathsf{F}}} e_{\mathsf{F}}]^{\mathsf{T}}]] \mathbb{E}[\tilde{w}_{\mathsf{F}}(n)] - \mu_{\mathsf{F}} \mathbb{E}[[\nabla_{\hat{w}_{\mathsf{F}}} e_{\mathsf{F}}] \varepsilon_{\mathsf{F}}], \qquad (C.5)$$

We can assume that the vector $[\nabla_{\hat{w}_F} e_F]$ which depends on the inputs, and ε_F which depends on the neural network design, are independent. Then we can write

$$\mathbb{E}[[\nabla_{\hat{w}_{\mathsf{F}}} e_{\mathsf{F}}]\varepsilon_{\mathsf{F}}] = 0. \tag{C.6}$$

Let us define the matrix R as

$$R \stackrel{\Delta}{=} \mathbb{E}[[\nabla_{\hat{w}_{\mathsf{F}}} e_{\mathsf{F}}] [\nabla_{\hat{w}_{\mathsf{F}}} e_{\mathsf{F}}]^{\mathsf{T}}]. \tag{C.7}$$

By substituting (C.6) and (C.7) in equation (C.5), it can be rewritten in the form

$$\mathbb{E}[\tilde{w}_{\mathsf{F}}(n+1)] = [I - \mu_{\mathsf{F}}R]\mathbb{E}[\tilde{w}_{\mathsf{F}}(n)]. \tag{C.8}$$

Then, the necessary condition for the convergence of this equation is

$$0 < \mu_{\mathsf{F}} < \frac{2}{\lambda_{max}},\tag{C.9}$$

where λ_{max} is the largest eigenvalue of the matrix R.

B. Biasness

Assuming small step-sizes and also the condition (C.9), the bias of estimation is calculated as

$$bias = \lim_{n \to \infty} -\mathbb{E}[\tilde{w}_{\mathsf{F}}(n)] = 0, \tag{C.10}$$

which means that if the step size is small enough such that the convergence condition holds, the parameters of the FPK equation tend to its optimal values.

C. Mean Square Convergence Analysis

The mean square deviation (MSD) of the estimation algorithm is defined as

$$MSD_{\mathsf{F}} = \lim_{n \to \infty} \mathbb{E}[\|\tilde{w}_{\mathsf{F}}(n)\|^2]. \tag{C.11}$$

In order to find the MSD, let us first define the weighted MSD of the algorithm as $\mathbb{E}[\|\tilde{w}_{\mathsf{F}}(n)\|_{\Sigma}^2]$, which can be obtained by the recursive equation

$$\mathbb{E}[\|\tilde{w}_{\mathsf{F}}(n+1)\|_{\Sigma}^{2}] = \mathbb{E}[\|\tilde{w}_{\mathsf{F}}(n)\|_{\Sigma'}^{2}] + \mu_{\mathsf{F}}^{2} \mathrm{tr}(R_{\Sigma})\|\varepsilon_{\mathsf{F}}\|^{2}, \qquad (C.12)$$

where Σ is a positive definite matrix, and

$$\Sigma' = (I - \mu_{\mathsf{F}} R)^{\mathsf{T}} \Sigma (I - \mu_{\mathsf{F}} R), \qquad (C.13)$$
$$R_{\Sigma} = \mathbb{E}[[\nabla_{\hat{w}_{\mathsf{F}}} e_{\mathsf{F}}]^{\mathsf{T}} \Sigma [\nabla_{\hat{w}_{\mathsf{F}}} e_{\mathsf{F}}]]. \qquad (C.14)$$

We know that $\operatorname{tr}(\Sigma X) = [\operatorname{vec}(X)]^{\mathsf{T}}\sigma$, and $\operatorname{vec}(U\Sigma V) = (V^{\mathsf{T}} \otimes U)\sigma$, where $\operatorname{vec}(\cdot)$ is a vectorazation operator, i.e., $\operatorname{vec}(\Sigma) = \sigma$. Using these equalities, we can obtain

$$(R_{\Sigma}) = [\operatorname{vec}(R)]^{\mathsf{T}}\sigma, \tag{C.15}$$

$$' = \mathcal{F}\sigma,$$
 (C.16)

$$\mathcal{F} = (I - \mu_{\mathsf{F}} R)^{\mathsf{T}} \otimes (I - \mu_{\mathsf{F}} R)^{\mathsf{T}}.$$
 (C.17)

At the convergence stage, the MSD is written as

tr

 σ

$$\lim_{n \to \infty} \mathbb{E}[\|\tilde{w}_{\mathsf{F}}(n)\|_{\Omega}^2] = \mu_{\mathsf{F}}^2 \|\varepsilon_{\mathsf{F}}\|^2 [\operatorname{vec}(R)]^{\mathsf{T}} (I - \mathcal{F})^{-1} \operatorname{vec}(\Omega), \quad (C.18)$$

where $\operatorname{vec}(\Omega) = (I - \mathcal{F})\sigma$ Therefore the steady state MSD is obtained as

$$\mathrm{MSD}_{\mathsf{F}} = \mu_{\mathsf{F}}^{2} \|\varepsilon_{\mathsf{F}}\|^{2} [\mathrm{vec}(R)]^{\mathsf{T}} (I - \mathcal{F})^{-1} \mathrm{vec}(I), \qquad (C.19)$$

The value of MSD can be very small by choosing a small value for step sizes, i.e., $\mu_{\rm F}$, and choosing a model which makes $\|\varepsilon_{\rm F}\|$ small.

0090-6778 (c) 2020 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications_standards/publications/rights/index.html for more information. Authorized licensed use limited to: Oulu University. Downloaded on August 21,2020 at 06:53:25 UTC from IEEE Xplore. Restrictions apply.

REFERENCES

- H. Kim, J. Park, M. Bennis, and S.-L. Kim, "Massive UAV-to-ground communication and its stable movement control: A mean-field approach," in *Proc. IEEE SPAWC, Kalamata, Greece*, Jun. 2018.
- [2] E. Ackerman and E. Strickland, "Medical delivery drones take flight in east africa," *IEEE Spectrum*, vol. 55, no. 1, pp. 34–35, Jan. 2018.
- [3] J. Tisdale, Z. Kim, and J. K. Hedrick, "Autonomous UAV path planning and estimation," *IEEE Robot. Autom. Mag.*, vol. 16, no. 2, pp. 35–42, Jun. 2009.
- [4] H. Shiri, J. Park, and M. Bennis, "Massive autonomous UAV path planning: A neural network based mean-field game theoretic approach," *arXiv preprint arXiv:1905.04152*, 2019.
- [5] M. Huang, R. P. Malhamé, P. E. Caines *et al.*, "Large population stochastic dynamic games: closed-loop McKean-Vlasov systems and the Nash certainty equivalence principle," *Commun. Inf. Syst.*, vol. 6, no. 3, pp. 221–252, 2006.
- [6] M. Huang, P. E. Caines, and R. P. Malhamé, "Large-population costcoupled LQG problems with nonuniform agents: individual-mass behavior and decentralized ε-Nash equilibria," *IEEE Trans. Autom. Control*, vol. 52, no. 9, pp. 1560–1571, 2007.
- [7] J.-M. Lasry and P.-L. Lions, "Mean field games," *Japanese J. math.* (*JJM*), vol. 2, no. 1, pp. 229–260, 2007.
- [8] H. B. McMahan, E. Moore, D. Ramage, S. Hampson *et al.*, "Communication-efficient learning of deep networks from decentralized data," *arXiv preprint arXiv:1602.05629*, 2016.
- [9] R. Shokri and V. Shmatikov, "Privacy-preserving deep learning," in Proc. 22nd ACM SIGSAC conf. Comput. comm. Secur., 2015, pp. 1310–1321.
- [10] D. Liu, D. Wang, F. Wang, H. Li, and X. Yang, "Neural-networkbased online HJB solution for optimal robust guaranteed cost control of continuous-time uncertain nonlinear systems," *IEEE Trans. Cybern.*, vol. 44, no. 12, pp. 2834–2847, Dec. 2014.
- [11] H. Shiri, J. Park, and M. Bennis, "Remote UAV online path planning via neural network based opportunistic control," *IEEE Wireless Commun. Lett.*, pp. 1–1, Feb. 2020.
- [12] K. P. Valavanis and G. J. Vachtsevanos, *Handbook of unmanned aerial vehicles*. Springer, 2015, vol. 1.
 [13] J. Lyu, Y. Zeng, and R. Zhang, "UAV-aided offloading for cellular
- [13] J. Lyu, Y. Zeng, and R. Zhang, "UAV-aided offloading for cellular hotspot," *IEEE Trans. Wireless Commun.*, vol. 17, no. 6, pp. 3988–4001, 2018.
- [14] A. A. Khuwaja, G. Zheng, Y. Chen, and W. Feng, "Optimum deployment of multiple UAVs for coverage area maximization in the presence of cochannel interference," *IEEE Access*, vol. 7, pp. 85 203–85 212, 2019.
- [15] C. W. Lim, C. K. Ryoo, K. Choi, and J. Cho, "Path generation algorithm for intelligence, surveillance and reconnaissance of an UAV," in *Proc. SICE Annu. Conf.*, 2010, pp. 1274–1277.
- [16] J.-H. Lee, J. Park, M. Bennis, and Y.-C. Ko, "Integrating LEO satellite and UAV relaying via reinforcement learning for non-terrestrial networks," 2020.
- [17] Y. Zeng and R. Zhang, "Energy-efficient UAV communication with trajectory optimization," *IEEE Trans. Wireless Commun.*, vol. 16, no. 6, pp. 3747–3760, 2017.
- [18] S. Zhang, Y. Zeng, and R. Zhang, "Cellular-enabled UAV communication: A connectivity-constrained trajectory optimization perspective," *IEEE Trans. Commun.*, vol. 67, no. 3, pp. 2580–2604, 2019.
 [19] U. Challita, W. Saad, and C. Bettstetter, "Interference management
- [19] U. Challita, W. Saad, and C. Bettstetter, "Interference management for cellular-connected UAVs: A deep reinforcement learning approach," *IEEE Trans. Wireless Commun.*, vol. 18, no. 4, pp. 2125–2140, 2019.
- [20] Q. Yang, J. Zhang, and G. Shi, "Path planning for unmanned aerial vehicle passive detection under the framework of partially observable markov decision process," in *Chinese Control Decision Conf. (CCDC)*, Jun. 2018, pp. 3896–3903.
- [21] J.-A. Delamer, Y. Watanabe, and C. P. C. Chanel, "Towards a MOMDP model for UAV safe path planning in urban environment," 9th Int. Micro Air Veh. Conf. Competition (IMAV), pp. 1–8, Sep. 2017. [Online]. Available: https://oatao.univ-toulouse.fr/19443/
- [22] A. Mardani, M. Chiaberge, and P. Giaccone, "Communication-aware UAV path planning," *IEEE Access*, vol. 7, pp. 52 609–52 621, Apr. 2019.
- [23] Y. Zeng, Q. Wu, and R. Zhang, "Accessing from the sky: A tutorial on UAV communications for 5G and beyond," arXiv preprint arXiv:1903.05289, 2019.
- [24] R. Austin, Unmanned aircraft systems: UAVS design, development and deployment. John Wiley & Sons, 2011, vol. 54.
- [25] "Enhanced LTE support for aerial vehicles." [Online]. Available: Available:ftp://www.3gpp.org/specs/archive/36_series/36.777
- [26] M. M. Azari, F. Rosas, A. Chiumento, and S. Pollin, "Coexistence of terrestrial and aerial users in cellular networks," in *IEEE GLOBECOM Workshops (GC Wkshps)*. IEEE, 2017, pp. 1–6.
- [27] L. Liu, S. Zhang, and R. Zhang, "Multi-beam UAV communication in cellular uplink: Cooperative interference cancellation and sumrate maximization," *Trans. Wireless. Comm.*, vol. 18, no. 10, p. 4679–4691, Oct. 2019. [Online]. Available: https://doi.org/10.1109/ TWC.2019.2926981

- [28] H. Hellaoui, A. Chelli, M. Bagaa, and T. Taleb, "Towards mitigating the impact of UAVs on cellular communications," in *IEEE Global Commun. Conf. (GLOBECOM)*, Dec. 2018.
- [29] H. Hellaoui, A. Chelli, M. Bagaa, T. Taleb, and M. Pätzold, "Towards efficient control of mobile network-enabled UAVs," in *IEEE Wireless Commun. Netw. Conf. (WCNC)*, Apr. 2019.
 [30] S. Zhang and R. Zhang, "Radio map based path planning for cellulartic descented UAV" in *IEEE Cluber Conf. (Cluber Conf. Cluber Conf. Cluber Clu*
- [30] S. Zhang and R. Zhang, "Radio map based path planning for cellularconnected UAV," in *IEEE Global Commun. Conf. (GLOBECOM)*, Dec. 2019, pp. 1–6.
- [31] C. Yan, L. Fu, J. Zhang, and J. Wang, "A comprehensive survey on UAV communication channel modeling," *IEEE Access*, vol. 7, pp. 107769– 107792, 2019.
- [32] B. Li, Z. Fei, and Y. Zhang, "Uav communications for 5G and beyond: Recent advances and future trends," *IEEE Internet of Things J.*, vol. 6, no. 2, pp. 2241–2263, 2019.
- [33] S. Zhang, H. Zhang, B. Di, and L. Song, "Cellular UAV-to-X communications: Design and optimization for multi-UAV networks," *IEEE Trans. Wireless Commun.*, vol. 18, no. 2, pp. 1346–1359, 2019.
- [34] H. M. La, "Multi-robot swarm for cooperative scalar field mapping," in *Handbook of Research on Design, Control, and Modeling of Swarm Robotics.* IGI Global, 2016, pp. 383–395.
- [35] H. M. La, W. Sheng, and J. Chen, "Cooperative and active sensing in mobile sensor networks for scalar field mapping," *IEEE Trans. Syst.*, *Man, Cybern. Syst.*, vol. 45, no. 1, pp. 1–12, 2015.
- [36] P. Sujit, S. Saripalli, and J. B. Sousa, "Unmanned aerial vehicle path following: A survey and analysis of algorithms for fixed-wing unmanned aerial vehicless," *IEEE Contr. Syst. Mag.*, vol. 34, no. 1, pp. 42–59, 2014.
- [37] J. M. Cano et al., "Quadrotor UAV for wind profile characterization," Proyectos Fin de Carrera, Universidad Carlos III de Madrid, 2013.
- [38] P. Abichandani, D. Lobo, G. Ford, D. Bucci, and M. Kam, "Wind measurement and simulation techniques in multi-rotor small unmanned aerial vehicles," *IEEE Access*, vol. 8, pp. 54910–54927, Mar. 2020.
- [39] H. Chen, K. Chang, and C. S. Agate, "UAV path planning with tangent-plus-Lyapunov vector field guidance and obstacle avoidance," *IEEE Trans. Aerosp. Electron. Syst.*, vol. 49, no. 2, pp. 840–856, 2013.
 [40] J.-h. Cui, R.-x. Wei, Z.-c. Liu, and K. Zhou, "UAV motion strategies
- [40] J.-h. Cui, R.-x. Wei, Z.-c. Liu, and K. Zhou, "UAV motion strategies in uncertain dynamic environments: A path planning method based on Q-learning strategy," *Appl. Sci.*, vol. 8, no. 11, p. 2169, 2018.
 [41] W. Li, M. Tan, L. Wang, and Q. Wang, "A cubic spline method
- [41] W. Li, M. Tan, L. Wang, and Q. Wang, "A cubic spline method combing improved particle swarm optimization for robot path planning in dynamic uncertain environment," *Int. J. Adv. Robot. Syst.*, vol. 17, no. 1, p. 1729881419891661, 2020. [Online]. Available: https://doi.org/10.1177/1729881419891661
- [42] H. X. Pham, H. M. La, D. Feil-Seifer, and L. V. Nguyen, "Autonomous UAV navigation using reinforcement learning," arXiv preprint arXiv:1801.05086, 2018.
- [43] W. Koch, R. Mancuso, R. West, and A. Bestavros, "Reinforcement learning for UAV attitude control," ACM Trans. Cyber-Phys. Syst., vol. 3, no. 2, Feb. 2019. [Online]. Available: https://doi.org/10.1145/3301273
- [44] C. Richter, W. Vega-Brown, and N. Roy, "Bayesian learning for safe high-speed navigation in unknown environments," in *Robotics Research*. Springer, 2018, pp. 325–341.
- [45] M. Nourian, P. E. Caines, and R. P. Malhamé, "Mean field analysis of controlled Cucker-Smale type flocking: Linear analysis and perturbation equations," *18th IFAC World Congress, Milan*, vol. 44, no. 1, pp. 4471– 4476, Aug. 2011.
- [46] C. Beck, S. Becker, P. Grohs, N. Jaafari, and A. Jentzen, "Solving stochastic differential equations and Kolmogorov equations by means of deep learning," *arXiv e-prints*, p. arXiv:1806.00421, Jun. 2018.
 [47] H. Xu, Z. Sun, and S. Xie, "An iterative algorithm for solving
- [47] H. Xu, Z. Sun, and S. Xie, "An iterative algorithm for solving a kind of discrete HJB equation with M-functions," *Appl. Math. Lett.*, vol. 24, no. 3, pp. 279 – 282, 2011. [Online]. Available: http://www.sciencedirect.com/science/article/pii/S0893965910003514
- [48] I. Kharroubi, N. Langrené, and H. Pham, "A numerical algorithm for fully nonlinear HJB equations: an approach by control randomization," *Monte Carlo Meth. Appl.*, vol. 20, no. 2, pp. 145–165, 2014.
- [49] K. W. Morton and D. F. Mayers, Numerical solution of partial differential equations: an introduction. Cambridge university press, 2005.
- [50] S. Khardi, "Aircraft flight path optimization. the Hamilton-Jacobi-Bellman considerations," *Appl. Math. Sci.*, vol. 6, no. 25, pp. pp–1221, 2012.
- [51] C. Greif, "Numerical methods for hamilton-jacobi-bellman equations," University of Wisconsin-Milwaukee, U.S.A, 2017.
- [52] Y. Tassa and T. Erez, "Least squares solutions of the HJB equation with neural network value-function approximators," *IEEE Trans. Neural Netw.*, vol. 18, no. 4, pp. 1031–1041, 2007.
- [53] B. Luo, H.-N. Wu, T. Huang, and D. Liu, "Reinforcement learning solution for HJB equation arising in constrained optimal control problem," *Neural Netw.*, vol. 71, pp. 150–158, 2015.
- [54] A. Faust, P. Ruymgaart, M. Salman, R. Fierro, and L. Tapia, "Continuous action reinforcement learning for control-affine systems with unknown dynamics," *IEEE/CAA J. Automatica Sinica*, vol. 1, no. 3, pp. 323–336, 2014.

- [55] J. W. Kim, B. J. Park, H. Yoo, J. H. Lee, and J. M. Lee, "Deep reinforcement learning based finite-horizon optimal tracking control for nonlinear system," *IFAC-PapersOnLine*, vol. 51, no. 25, pp. 257–262, 2018.
- [56] S. Samarakoon, M. Bennis, W. Saad, and M. Debbah, "Federated learning for ultra-reliable low-latency V2V communications," in *IEEE Global Commun. Conf. (GLOBECOM)*, Dec. 2018, pp. 1–7.
- [57] T. Zeng, O. Semiari, W. Saad, and M. Bennis, "Joint communication and control for wireless autonomous vehicular platoon systems," *IEEE Trans. Comm.*, vol. 67, no. 11, pp. 7907–7922, 2019.
 [58] Y. Y. Nazaruddin, A. Widyotriatmo, T. A. Tamba, M. S. Arifin, and R. A.
- [58] Y. Y. Nazaruddin, A. Widyotriatmo, T. A. Tamba, M. S. Arifin, and R. A. Santosa, "Communication-efficient optimal-based control of a quadrotor UAV by event-triggered mechanism," in 2018 5th Asian Conf. Def. Tech. (ACDT), 2018, pp. 96–101.
 [59] H. Shiri, M. A. Tinati, M. Codreanu, and G. Azarnia, "Distributed sparse
- [59] H. Shiri, M. A. Tinati, M. Codreanu, and G. Azarnia, "Distributed sparse diffusion estimation with reduced communication cost," *IET Signal Process.*, vol. 12, no. 8, pp. 1043–1052, 2018.
 [60] N. Goddemeier and C. Wietfeld, "Investigation of air-to-air channel of the statement of the stat
- [60] N. Goddemeier and C. Wietfeld, "Investigation of air-to-air channel characteristics and a UAV specific extension to the Rice model," in *IEEE GLOBECOM Workshops (GC Wkshps), San Diego, CA, USA*, Dec. 2015, pp. 1–5.
- [61] R. Zárate-Minano, F. M. Mele, and F. Milano, "SDE-based wind speed models with Weibull distribution and exponential autocorrelation," in *Proc. IEEE PESGM*, Boston, MA, USA, 2016.
- Proc. IEEE PESGM, Boston, MA, USA, 2016.
 [62] H. Inaltekin, M. Gorlatova, and M. Chiang, "Virtualized control over fog: interplay between reliability and latency," *IEEE Internet Things J.*, vol. 5, no. 6, pp. 5030–5045, 2018.
- vol. 5, no. 6, pp. 5030–5045, 2018.
 [63] Y. Zeng, J. Xu, and R. Zhang, "Energy minimization for wireless communication with rotary-wing UAV," *IEEE Trans. Wireless Commun.*, vol. 18, no. 4, pp. 2329–2345, Apr. 2019.
- vol. 18, no. 4, pp. 2329–2345, Apr. 2019.
 [64] F. Cucker and J.-G. Dong, "Avoiding collisions in flocks," *IEEE Trans. Autom. Control*, vol. 55, no. 5, pp. 1238–1243, 2010.
 [65] Z. Lin, M. Broucke, and B. Francis, "Local control strategies for groups
- [65] Z. Lin, M. Broucke, and B. Francis, "Local control strategies for groups of mobile autonomous agents," *IEEE Trans. Autom. Control*, vol. 49, no. 4, pp. 622–629, 2004.
- [66] G. Vásárhelyi, C. Virágh, G. Somorjai, T. Nepusz, A. E. Eiben, and T. Vicsek, "Optimized flocking of autonomous drones in confined environments," *Sci. Robot.*, vol. 3, no. 20, 2018. [Online]. Available: http://robotics.sciencemag.org/content/3/20/eaat3536
- http://robotics.sciencemag.org/content/3/20/eaat3536
 [67] D. Liu, D. Wang, F.-Y. Wang, H. Li, and X. Yang, "Neural-network-based online HJB solution for optimal robust guaranteed cost control of continuous-time uncertain nonlinear systems," *IEEE Trans. Cybern.*, vol. 44, no. 12, pp. 2834–2847, Dec. 2014.
- [68] Z. Shiller and H.-H. Lu, "Computation of path constrained time optimal motions with dynamic singularities," 1992.
- [69] H. Shiri, M. A. Tinati, M. Codreanu, and S. Daneshvar, "Distributed sparse diffusion estimation based on set membership and affine projection algorithm," *Digital Signal Processing*, vol. 73, pp. 47–61, 2018.



Hamid Shiri (S'19) is a researcher at the Center for Wireless Communications, University of Oulu, Finland. He has been with the Faculty of Electrical and Computer Engineering, University of Tabriz, Tabriz, Iran from 2005 to 2018. His recent research interests focus are machine learning, signal processing, distributed control and optimization in 5G networks and beyond, URLLC, and Mean-Field Game theory.



Jihong Park (S'09-M'16) is a Lecturer (assistant professor) at the School of IT, Deakin University, Australia. He received the B.S. and Ph.D. degrees from Yonsei University, Seoul, Korea, in 2009 and 2016, respectively. He was a Post-Doctoral Researcher with Aalborg University, Denmark, from 2016 to 2017; the University of Oulu, Finland, from 2018 to 2019. His recent research focus includes communication-efficient distributed machine learning, distributed control, and distributed ledger technology, as well as their applications for beyond 5G/6G communication systems. He served as a

16

Conference/Workshop Program Committee Member for IEEE GLOBECOM, ICC, and WCNC, as well as NeurIPS, ICML, and IJCAI. He received the IEEE GLOBECOM Student Travel Grant in 2014, the IEEE Seoul Section Student Paper Contest Bronze Prize in 2014, and the 6th IDIS-ETNEWS (The Electronic Times) Paper Contest Award sponsored by the Ministry of Science, ICT, and Future Planning of Korea. Currently, he is an Associate Editor of Frontiers in Aerial and Space Networks, and a Guest Editor of MDPI Telecom SI on "millimeter wave communiations and networking in 5G and beyond."



Mehdi Bennis is an Associate Professor at the Centre for Wireless Communications, University of Oulu, Finland, Academy of Finland Research Fellow and head of the intelligent connectivity and networks/systems group (ICON). His main research interests are in radio resource management, heterogeneous networks, game theory and distributed machine learning in 5G networks and beyond. He has published more than 200 research papers in international conferences, journals and book chapters. He has been the recipient of several prestigious awards including the 2015 Fred W. Ellersick Prize

teris, he has been ute recipient of several prestigious awards including the 2015 Fred W. Ellersick Prize from the IEEE Communications Society, the 2016 Best Tutorial Prize from the Journal of Wireless Communications and Networks, the all-University of Oulu award for research and the 2019 IEEE ComSoc Radio Communications Committee Early Achievement Award. Dr Bennis is an editor of IEEE TCOM and Specialty Chief Editor for Data Science for Communications in the Frontiers in Communications and Networks journal.