

Age of Information Minimization using Multi-agent UAVs based on AI-Enhanced Mean Field Resource Allocation

Yousef Emami, *Member, IEEE*, Hao Gao, *Student Member, IEEE*, Kai Li, *Senior Member, IEEE*, Luis Almeida, *Senior Member, IEEE*, Eduardo Tovar, *Member, IEEE*, and Zhu Han, *Fellow, IEEE*

Abstract—Unmanned Aerial Vehicle (UAV) swarms play an effective role in timely data collection from ground sensors in remote and hostile areas. Optimizing the collective behavior of swarms can improve data collection performance. This paper puts forth a new mean field flight resource allocation optimization to minimize age of information (AoI) of sensory data, where balancing the trade-off between the UAVs' movements and AoI is formulated as a mean field game (MFG). The MFG optimization yields an expansive solution space encompassing continuous state and action, resulting in significant computational complexity. To address practical situations, we propose, a new mean field hybrid proximal policy optimization (MF-HPPO) scheme to minimize the average AoI by optimizing the UAV's trajectories and data collection scheduling of the ground sensors given mixed continuous and discrete actions. Furthermore, a long short term memory (LSTM) is leveraged in MF-HPPO to predict the time-varying network state and stabilize the training. Numerical results demonstrate that the proposed MF-HPPO reduces the average AoI by up to 45% and 57% in the considered simulation setting, as compared to multi-agent deep Q-learning (MADQN) method and non-learning random algorithm, respectively.

Index Terms—UAV, Mean-field game, Age of information, Proximal policy optimization, Long short term memory.

I. INTRODUCTION

Thanks to high mobility, unmanned aerial vehicles (UAVs) are widely used in search and rescue [1], monitoring and surveillance [2], [3], aerial data relay [4], [5], construction [6] and parcel delivery [7]. A large number of ground sensors can be deployed in a target area to monitor the condition. Multiple UAVs can be employed to gather sensory data [8]. UAVs-assisted data collection offers several advantages for the data collection in remote and human-unfriendly environments. Specifically, UAVs have the ability to reach areas that are difficult to access by humans, making data collection more efficient and cost-effective. This reduces safety risks since the use of UAVs eliminates the need for human intervention in hazardous environments. Due to mobility, UAVs are able to cover large areas, which reduces the time and resources required for data collection [9], [10].

In UAVs-assisted sensor networks, Age of Information (AoI) is typically used to describe the freshness of the sensory data [11], which refers to the elapsed time between the generation of a piece of sensory data at a ground sensor and the receipt of that data at the UAV. Therefore, AoI takes into account not

only the time it takes for data to be transmitted, but also any delays that may occur in the network. When the UAV flight is not properly controlled, the UAV can move away from the ground sensor, causing the AoI of the sensor to be extended and resulting in data expiration. In addition, AoI of the ground sensors' data can be different from each other, since the sensory data generation is impacted by natural conditions monitored [12]. The instantaneous knowledge of data generation rate and channel conditions of the ground sensors are not available or can only be partially observed by the UAV in practical systems. Therefore, jointly optimizing the cruise control of the UAVs and communication schedules of the ground sensors to minimize AoI is non-trivial. Moreover, swift movements of the UAV results in poor channel conditions and fast signal attenuation, giving rise to frequent data retransmissions and a prolonged AoI. In contrast, slow movements of the UAV extend flight time, thereby extending AoI of the ground sensors. Meanwhile, the UAV should consider the movement of other UAVs to minimize the average AoI by coordinating their velocities. Therefore, finding an equilibrium solution among UAVs to optimize the velocity and reduce the average AoI is essential.

The decentralized approach is particularly relevant in situations where the UAVs have limited information about the actions of others, in terms of the trajectory, flight speed, and the scheduled ground sensors. Game theory can be leveraged to design decentralized control and determine the equilibrium in UAV networks [13]. However, traditional game theory approaches that focus on interactions between a finite number of UAVs (or players) can become intractable when dealing with a large number of UAVs, as the complexity of solving such games grows exponentially with the number of UAVs. Diverging from traditional game theory, mean field game (MFG) offers a scalable framework to addressing the proposed joint optimization of the cruise control and communication schedules, by approximating the interactive behavior of a large number of UAVs using a continuum or mean field, thereby considerably reducing computational complexity. Moreover, MFG allows the UAVs to make decisions based on the behavior of the overall swarms rather than on the actions of individual UAVs.

In this paper, we propose a cruise control optimization of the UAV swarms based on MFG to minimize AoI, while balancing the trade-off between the UAVs' movements and AoI. In our MFG, the optimal velocities of the UAVs are determined by a Fokker-Planck-Kolmogorov (FPK), which describes an evolution of the mean field for achieving an equilibrium of the

optimal velocities of the UAVs. However, the proposed MFG is difficult to be solved online in practical scenarios, since the instantaneous knowledge of the UAV's cruise control decision and AoI is hardly known by other UAVs. Furthermore, the MFG for the flight resource allocation optimization is formulated as a multi-agent Markov Decision Process (MMDP), where network states consist of AoI of the ground sensors and waypoints of the UAV swarm. The action space in the MMDP contains the waypoints and velocities in a continuous space and transmission schedule in a discrete space. Based on the formulated MMDP, we propose a mean field hybrid proximal policy optimization (MF-HPPO) to minimize AoI of the ground sensors, where the UAV swarms learn each other's decisions of cruise control and transmission schedules of the ground sensors.

In summary, UAV swarms are employed to collect time critical sensory data. Time-critical data collection is influenced by the velocity of the UAVs and their coordinated interactions in the swarms, which can be modeled using MFG. This raises the importance of an age-optimal cruise control based on MFG for UAVs. However, determining the equilibrium online is difficult in practical scenarios, and thus we propose MF-HPPO that highly reduces the complexity while minimizing the average AoI. The contributions of this paper are listed as follows:

- We novelly formulate the MFG optimization with a large number of UAVs to address the trade-off between the cruise control of the UAVs and AoI. Due to the high computational complexity of the MFG, MF-HPPO is proposed to minimize the average AoI, where the state dynamics are learned and the actions of the UAVs are optimized in a mixed discrete and continuous action space.
- To capture temporal dependencies of the cruise control while improving the learning convergence, a new long short term memory (LSTM) layer is developed with the proposed MF-HPPO to predict the time-varying network states, i.e., AoI and UAVs' waypoints.
- Numerical results demonstrate that MF-HPPO achieves fast convergence (less than 200 iterations). Meanwhile, our proposed MF-HPPO witnesses 45% and 57% reduction in AoI as compared to a multi-agent deep Q-learning (MADQN) method (which performs trajectory planning in the discrete space) and an existing non-learning random algorithm, respectively.

The rest of this paper is organized as follows: In Section II, we present the state-of-the-art in intelligent flight resource allocation, with a focus on MFG and time-critically. In Section III, we formulate the channel model as well as the AoI in the UAVs-assisted sensor network. In Section IV, we formulate the flight resource allocation of the UAV swarm as the MFG to minimize the AoI. Section V develops the proposed MF-HPPO, to jointly optimize the cruise control of multiple UAVs and data collection scheduling. Section VI presents the implementation of the proposed MF-HPPO in Pytorch as well as performance evaluation. Finally, Section VII concludes this paper.

A. Mean field flight resource allocation

In [14], the authors explore energy-efficient control strategies for UAVs that provide fair communication coverage for ground users. The UAV control problem is modeled as an MFG and a mean-field trusted region policy optimization (TRPO) algorithm is studied to design the UAVs' trajectories. In [15], the authors apply the MFG theory to the downlink power control problem in ultra-dense UAV networks to improve the network's energy efficiency. Due to the complexity of the MFG, a DRL-MFG algorithm is developed to learn the optimal power control strategy. [16] studies the task allocation in cooperative mobile edge computing and a mean field guided Q-function is formulated to reduce the network latency. MFG and deep reinforcement learning (DRL) are integrated to guide the learning process of DRL according to the equilibrium of MFG. In [17], the authors model the trajectory planning and power control for heterogeneous UAVs as an MFG, aiming to reduce energy consumption. A mean field Q-learning is studied to find the optimal solution. In [18], the authors study UAV-assisted ultra-dense networks, where each UAV can adjust its location to reduce the AoI. They formulate the problem as an MFG and apply a deep deterministic policy gradient (DDPG)-MFG algorithm to find the mean field equilibrium. In [15], downlink power control for a large number of UAVs is suggested to enhance the energy efficiency by learning the optimal power control policy. MFG is used to model the power control problem of the UAV network, where each UAV tries to enhance the energy efficiency by adjusting its transmit power. Then, due to the complexity of solving the formulated MFG, an effective DRL-MFG algorithm is suggested to learn the optimal power control strategy.

Although, DRL-based solutions mainly used, the following works adopt numerical solutions. In [19], the focus is on adaptive coverage problem in emergency communication system, where multiple UAV act as aerial base stations to serve randomly distributed users. The problem is formulated using discrete MFG, each UAV aim to reduce its flight energy consumption and increase the number of users it can serve. Finally, optimal control and state of each UAV is computed. In [20], a discrete MFG is formulated to address joint adjustment of power and velocity for a large number of UAVs that act as aerial base stations. Decentralized control laws are developed, and mean field equilibrium is analyzed. In [21], the authors present an energy-efficient velocity control algorithm for a large number of UAVs based on the MFG theory. The velocity control of the UAVs is modeled using a differential game in which energy and delay are balanced by using an original double mixed gradient method. In [22], a multi-UAV enabled mobile Internet of Vehicles (IoV) model is designed, to enable the UAVs to track the movements of the vehicles. A joint optimization problem is formulated to improve the total system throughput over the flight time by jointly adjusting vehicle communication scheduling, UAV power allocation, and UAV trajectory. In [23], a multi-UAV-enabled IoT is investigated

to improve the minimum energy efficiency of each UAV by jointly adjusting communication scheduling, power allocations, and trajectories of the UAVs.

B. Time-critical flight resource allocation

In [24], the authors consider ground sensors with limited energy and apply airborne base stations to collect sensory data. Each UAV's task is decomposed into energy transfer and fresh data collection. A centralized multi-agent DRL based on DDPG is developed to adjust the UAV trajectories in a continuous action space, to reduce the AoI of the ground sensors. In [25], the authors study UAV-assisted sensor networks where multiple UAVs cooperatively conduct the data collection to reduce the AoI. The trajectory planning is formulated as a decentralized partially observable Markov Decision Process (Dec-POMDP). A multi-agent DRL is studied to find the optimal strategy. In [26] and [27], the authors develop the trajectory planning for multiple UAVs that perform cooperative sensing and transmission, aiming to reduce the AoI. In [28], ground sensors sample and upload data in a UAV-assisted IoT network. PPO is used to explore the optimal scheduling policy and altitude control for the UAV to reduce the AoI. In [29], a data collection scheme characterized by AoI and energy consumption in a UAV-assisted IoT network is investigated. The average AoI, and energy consumption of propulsion and communication are reduced by adjusting the UAV flight speed, hovering waypoints, and bandwidth allocation for data collection using a TD3-based approach. In [23] a multi-UAV enabled IoT is investigated to maximize the minimum energy efficiency of each UAV by jointly optimizing communication scheduling, power allocations, and trajectories of the UAVs. In [30], a UAV-assisted IoT network is investigated in the presence of eavesdroppers. The communication UAV and the jamming UAV cooperate to collect data from IoT devices. A TD3-based solution is developed to reduce the AoI by adjusting the UAV trajectory, computations, and spectrum resource allocation strategies.

Although, DDPG and PPO used to adjust continuous and discrete actions to reduce AoI, the following works use DQN and QMIX to adjust discrete actions. In [31], the authors investigate UAV-assisted IoT networks where multiple UAVs relay data between sensors and base station. A DQN-based trajectory planning algorithm is presented to reduce the AoI. In [32], ground sensors with limited energy are used to observe various physical processes in the context of a UAV-assisted wireless network. The trajectory and scheduling policy are adjusted to reduce the weighted sum of AoI, and a DQN-based solution is applied to obtain the best strategy. In [33], trajectory planning of the UAV is performed to reduce the AoI in a UAV-assisted IoT network. The problem is formulated as an MDP, and a DQN-based algorithm is studied to find the optimal trajectories of the UAV. In [34], a UAV-assisted data collection for ground sensors is studied, where the UAV with limited energy is dispatched to collect sensory data. The UAV's trajectory is adjusted to reduce the average AoI and keep the packet loss rate low. The trajectory planning is formulated as

an MDP while DQN is applied to design the UAV's trajectory. In [35], a UAV-assisted wireless network with an energy supply is used, where the UAV performs wireless energy transmission to ground sensors, and the sensors transmit data to the UAV using the harvested energy. A DQN-based trajectory planning algorithm is presented to reduce the average AoI by adjusting the trajectory, transmission schedule, and harvested energy. In [36], a massive deployment of up to a hundred IoT devices is investigated, where multiple UAVs serve as mobile relay nodes to reduce the AoI and energy consumption. A DQN-based solution is developed to jointly reduce the AoI and energy consumption of the devices by adjusting the trajectory and scheduling policy. In [37], a UAV-assisted sensor network is investigated, where the UAV flies between the resource-constrained sensors and collects their status updates. A DQN-based solution is developed to jointly adjust the trajectory of the UAV and the scheduling of the sensors to reduce the Age of Synchronization (AoS), which considers both the freshness and the content of the information. In [38], a massive IoT communication scenario is investigated where a UAV swarm collects fresh information from IoT devices and provides better coverage and LoS for the IoT network. To mitigate high-dimensional problems with high complexity, a multi-agent DRL based on DQN is developed. In [39], the problem of optimal data collection in IoT networks with multiple cooperative UAVs is investigated. Kinematic, energy, trajectory, and collision avoidance constraints are considered. To achieve this goal, a QMIX-based algorithm is developed to reduce the total average AoI by jointly adjusting the trajectories of the UAVs and the scheduling of the sensors.

Most of the works, in Section II-A, formulate an MFG for multiple UAVs to address energy efficiency. Authors in [18] formulate an MFG to minimize AoI and suggest DDPG-MFG in continuous action space to find the optimal solution.

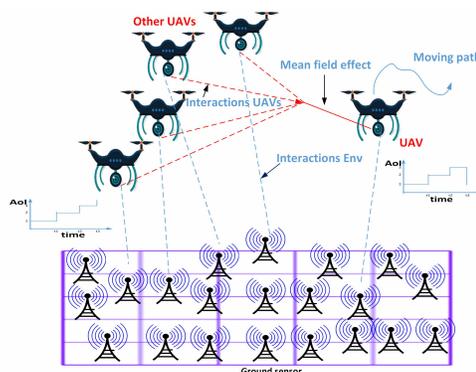


Fig. 1: Mean field representation of UAVs-assisted sensor networks.

The works in Section II-B investigate resource allocation to reduce AoI, however, actions are adjusted in continuous or discrete action spaces. In contrast, in this paper, we novelly formulate the MFG optimization with a large number of UAVs to address the trade-off between the cruise control of the UAVs

and AoI. Due to the high computational complexity of the MFG, MF-HPPO is proposed to minimize the average AoI, where the state dynamics are learned and the actions of the UAVs are optimized in a mixed discrete and continuous action space.

III. SYSTEM MODEL

In this section, we present the system model of the considered UAVs-assisted sensor network. Notations used in this paper are summarized in Table I. The system consists of I UAVs, $i \in [1, I]$ and J ground sensors, $j \in [1, J]$ in which the ground sensors are deployed in a target region. The UAVs are employed to patrol in the target zone while collecting the sensory data. Fig. 1 depicts an example of UAVs-assisted sensor network along with mean field representation. With the increase in the number of UAVs in Fig. 1 the interactions between them become complex and can dominate the overall behavior of the system. MFG designed to deal with the optimal control problem involving a large number of players. It has unique characteristics suitable for UAV swarm and modelling these interactions. Each UAV seeks to minimize the AoI according to the actions of other agents surrounded. As depicted, the UAV consider the mean field effect of the other UAVs, which represents the collective behavior of the UAVs in the system. The coordinates (x_i, y_i, z_i) and $(x_j, y_j, 0)$ represent the position of UAV i and ground sensor j , respectively. The UAVs fly to the ground sensors, collect sensory data, and then their operation is terminated. The UAVs fly at a constant altitude, represented by $\zeta_i(t) = (x_i, y_i, z)$. The distance between ground sensor j and UAV i is $\sqrt{(x_i - x_j)^2 + (y_i - y_j)^2 + z^2}$. For the safety of the UAV during flight by preventing it from exceeding the maximum safe speed or stalling, we denote the maximum and minimum velocity of the UAV as v_{max} and v_{min} , respectively. The lower bound of the speed, i.e., v_{min} , is set to 0. We consider that UAV i moves in low attitude for data collection, where the probability of LoS communication between UAV i and ground sensor j is given by [40]

$$\Pr_{LoS}(\varphi_j^i) = \frac{1}{1 + a \exp(-b[\varphi_j^i - a])} \quad (1)$$

where a and b are constants, and φ_j^i denotes the elevation angle between UAV i and ground sensor j . The elevation angle of the path loss model is given by

$$\varphi = \tan^{-1} \left(\frac{h}{d} \right), \quad (2)$$

where h is the height of the UAV above the ground and d is the horizontal distance between the UAV and the ground sensor. By changing d , the horizontal distance between the UAV and the ground sensor, we can influence the path loss of our UAV communication through the elevation angle. Moreover, path loss

TABLE I: Notation and Definition

Notation	Definition
J	number of ground sensors
I	number of UAVs
$h_j^i(t)$	channel gain between device j and UAV i
$\zeta_i(t)$	location of the UAV on its trajectory
$v_i(t)$	velocity of UAV i
v_{max}, v_{min}	the maximum and minimum velocity of UAV i
M	number of episodes
L	length of each episode
γ	discount factor
η	learning rate
D	buffer size
B	mini-batch size
a_i	action of UAV i
o_i	mean field of UAV i
a_i^c	continuous action of UAV i
a_i^d	discrete action of UAV i
$s_{\alpha, i}$	state of UAV i
$E[...]$	mathematical expectation
A	advantage function
θ	network parameter
π	policy
π^c	continuous policy
π^d	discrete policy
σ	diffusion coefficient
W	weiner process
H	entropy

of the channel between UAV i and device j can be modeled by [41]

$$\gamma_j^i = \Pr_{LoS}(\varphi_j^i)(\eta_{LoS} - \eta_{NLoS}) + 20 \log(r \sec(\varphi_j^i)) + 20 \log(\lambda) + 20 \log\left(\frac{4\pi}{v_c}\right) + \eta_{NLoS} \quad (3)$$

where r is the radius of the radio coverage of UAV i , λ is the carrier frequency, and v_c is the speed of light. η_{LoS} and η_{NLoS} are the excessive path losses of LoS or non-LoS, respectively.

To characterize the freshness of the collected sensory data at the UAV, AoI is defined as the time that has passed since ground sensor generates the latest information. The AoI of ground sensor j that generated a data packet at t_j and collected by UAV i at t_i is given by

$$AoI_j^i(t) = t_i - t_j. \quad (4)$$

According to (4), it can be also known that maintaining a low $AoI_j^i(t)$ is critical for improving the effectiveness and timeliness of the sensory data, reducing the response time, and providing real-time information for decision-making at the UAVs. In UAVs-assisted sensor networks, the AoI can be affected by the path loss condition: Decreased path loss attenuates signal deterioration during the conveyance from ground sensors to UAVs, culminating in enhanced signal strength, expedited data transfer, diminished error rates, and subsequently, a reduced necessity for retransmissions. Such efficacious communication markedly diminishes the AoI and guarantees the contemporaneity and pertinence of the data for real-time applications. Conversely, elevated path loss results in attenuated signals, escalated error frequencies, increased

retransmissions, and ultimately, an augmented AoI, thereby impacting the promptness and dependability of the data.

IV. PROBLEM FORMULATION

In this section, we formulate the MFG optimization with a large number of UAVs to address the trade-off between the cruise control of the UAVs and AoI. We also explore the FPK equation to determine the optimal velocities of the UAVs while characterizing the collective behavior of the UAVs. We begin with optimal control formulation in Section IV-A and then proceed with MFG formulation in Section IV-B.

A. Optimal Control Formulation

We derive the state dynamics and cost function, then we formulate the velocity control problem using the optimal control theory.

1) Time-varying Dynamics of Network States

Let $\zeta_i(t)$ denote the position of the UAV i at time t and $v_i(t)$ denotes the velocity. According to Newton's laws of motion [42], the location dynamics of UAV i can be expressed by

$$d\zeta_i(t) = v_i(t)dt + \sigma dW_i(t) \quad (5)$$

where $W_i(t)$ is a standard Wiener process [43] with a diffusion coefficient σ .

2) Cost Function

Each UAV intends to optimize its velocity to minimize the cost function. Our cost is defined as the average AoI of all ground sensors. The average AoI can be computed as:

$$c(t) = \frac{1}{IJ} \sum_{j=1}^J \sum_{i=1}^I AoI_j^i(t). \quad (6)$$

3) Velocity Control Problem Formulation

Given a period of time T regarding the data collection, the velocity of UAV i at t , denoted as $v_i^*(t)$, is optimally controlled to minimize $c(t)$, which gives:

$$v_i^*(t) = \arg \min_{v_i(t)} E \int_0^T c(t) dt, \quad (7)$$

s.t. (5).

Equation (7) is the integral of $c(t)$ on the given time range $(0, T)$, while $c(t)$ is the average AoI of all UAVs defined in (6). And the definition of AoI is defined in (4). In summary, physical meaning of (7) is to find the optimal velocity of UAVs to minimize the accumulated average AoI in the time range $(0, T)$. Constraint (5) is a differential constraint of a velocity control problem. To be more precise, the Newtonian motion function describes the change of locations of UAVs w.r.s.t their velocities. The Brownian motion, σ , represents the random effects, which might influence the locations of UAVs. The use of game theory in our model, as opposed to direct optimization, is crucial due to the inherent interdependencies between the decisions of multiple UAVs that are not explicitly captured in (7). While equation (7) represents a stand-alone minimization problem for the control strategy $v_i(t)$ of a single

UAV, the choice of $v_i(t)$ by one UAV in a real-world scenario with multiple UAVs implicitly affects the operational efficiency and AoI outcomes of the others. To determine $v_i^*(t)$ in (7), classical game theories, such as differential game, fails to capture the aggregate behavior of all the UAVs. Differential game assumes each agent's movement is independent of others. This assumption fails to capture the fact that a large number of UAVs' trajectories decisions are influenced by the aggregate behavior of all the UAVs, thus hardly minimizing the average AoI, $c(t)$.

We novelly extend MFG to capture the impact of the aggregate behavior of the UAVs, in terms of cruise control. The MFG models the aggregate decision of UAVs as a probability distribution, rather than focusing on the actions of individual UAVs. This recognizes that the cruise control of each UAV is influenced by the behavior of all other UAVs. Moreover, the formulated MFG is defined to minimize $c(t)$ given a large number of UAVs, which classical game theory struggles with due to the computational complexity of solving for the equilibrium.

B. MFG Problem Formulation

We reformulate the optimal cruise control problem in (7) into a cooperative MFG problem. The computational complexity of the system is greatly reduced by formulating an MFG, since a large number of interactions with other agents is converted into an interaction with the mass. The interaction between each UAV with the other UAVs is modeled as a mean-field term, which is denoted by $m(\zeta(t))$. The mean-field term is the distribution over agents' state space or control to model the overall state and control of them. We can measure the state and control of all agents in an MFG using the mean-field term.

Given dynamics, $\zeta_i(t)$, the mean-field term of $m(\zeta(t))$ can be denoted by

$$m(\zeta(t)) = \lim_{I \rightarrow \infty} \frac{1}{I} \sum_{i=1}^I \mathbb{1}\{\zeta_i(t) = \zeta(t)\}, \quad (8)$$

where $\mathbb{1}$ is an indicator function which returns 1 if the given condition is true, or 0, otherwise.

Given $m(\zeta(t))$, the state dynamics, cost function and FPK equation can be defined as:

- **State dynamics:** The state dynamics of each UAV can be expressed by

$$d\zeta(t) = v(t)dt + \sigma dW(t). \quad (9)$$

- **Cost function:** The mean-field term affects the running cost function of each UAV. The average AoI of the all UAVs is computed by

$$c(v(t), m(\zeta(t))) = \int c(v(t)) \cdot m(\zeta(t)) d\zeta. \quad (10)$$

Mathematically, the cost function can be written by

$$J(v(t), m(\zeta(t))) = \int_{t=0}^T c(v(t), m(\zeta(t))) dt. \quad (11)$$

If the UAV move quickly, lead to poor channel condition and retransmissions thereby AoI prolongs. In contrast, slow movement of the UAV, may prolong the AoI of the ground sensors because the data are not collected in time. The cost function addresses these trade-offs and find the optimal velocity to balance these objectives.

- **Focker-Planck equation:** Based on (9) we develop the FPK equation. The FPK equation governs the evolution of the mean field function of UAVs and given by:

$$\partial_t m(\zeta(t)) + \nabla_{\zeta} m(\zeta(t)) \cdot v(t) - \frac{\sigma^2}{2} \nabla_{\zeta}^2 m(\zeta(t)) = 0. \quad (12)$$

See *Appendix*.

After deriving the state dynamics, cost function, and FPK equation, we now proceed to present the MFG.

To summarize, the cooperative MFG problem is given by

$$\begin{aligned} \min_{v, m} J(v(t), m(\zeta(t))) \quad (13) \\ s.t. \quad (12). \end{aligned}$$

V. PROPOSED MF-HPPO

In this section, we present background on PPO in Section V-A and then describe the MFG as an MMDP in Section V-B so that the optimal actions of UAVs can be learned by the proposed MF-HPPO. MF-HPPO is presented in Section V-C, which employs onboard PPO to minimize the average AoI of the ground sensors. The trajectory and instantaneous speed of the UAVs, and the selection of the ground sensors are optimized in a mixed action space. In Section V-D, an LSTM layer is developed with MF-HPPO to capture the long-term dependency of data.

A. Background on PPO

In this work, we use policy-based DRLs because of their superior performance compared to state-of-the-art algorithms. A major issue in this category is update instability, which is rooted in the variability of the step-size parameter for policy optimization. Small steps slow the learning process, while large steps degrade policy performance. TRPO [44] addresses this issue by defining a trusted region for changes and using a complex second-order method. PPO follows the same approach by presenting a first-order method to overcome the high complexity of TRPO. PPO bounds the policy update within a range as an alternative to the hard constraint of TRPO. PPO has two primary variants: PPO-penalty and PPO-clip. PPO-penalty changes the hard constraint of TRPO to a penalty in objective function. PPO-clip does not have a constraint and use clipping

techniques to bound the changes of policy [45]. The PPO-clip objective function can be written as follows:

$$\begin{aligned} L^{clip}(\theta) = \min \left(\frac{\pi(a|s)}{\pi_{old}(a|s)} A_{\pi_{old}}(s, a), \right. \\ \left. g(\epsilon, A_{\pi_{old}}(s, a)) \right), \quad (14) \end{aligned}$$

where

$$g(\epsilon, A) = \begin{cases} (1 + \epsilon)A, & A \geq 0, \\ (1 - \epsilon)A, & A < 0. \end{cases} \quad (15)$$

where ϵ is used to control the clip range. PPO uses Actor-Critic framework in the implementation step. The overall objective function is given by:

$$L^{total}(\theta) = L^{clip}(\theta) - K_1 L^{VF}(\theta) + K_2 * H \quad (16)$$

where K_1 and K_2 are loss coefficients, H is entropy. L^{VF} is loss for Critic network. The proposed approach, MF-HPPO, leverages PPO-clip as model-free, on-policy and policy gradient DRL algorithm and is capable to optimize continuous and discrete actions. In the following, we formulate our problem using MMDP and then address it using MF-HPPO.

B. MMDP Formulation

We reformulate the MFG using MMDP framework to enable the application of PPO for optimizing the actions and minimizing average AoI. By adapting the MMDP framework to our problem, we define the relevant state space, action space, transition probabilities, policy and cost function, thus facilitating an effective solution approach based on MF-HPPO. We define our MMDP as follows.

- *Agents:* the number of agents, i.e., UAVs is denoted by I .
- *State:* A state s_{α} of the MMDP consists of the positions of UAV i , the AoI of ground sensors, i.e., $s_{\alpha} = \{\zeta_i(t), AoI_j^i(t) : i \in [1, I], j \in [1, J]\}$. All states of the MMDP constitute the state space.
- *Action:* Each UAV i takes an action a_i that schedules a ground sensor for data transmission and determines the flight trajectory and velocity, i.e., $a_i = \{k_j^i, v_i(t), \zeta_i(t)\}$
- *Policy:* Policy π_i is the probability of taking each action of agent i .
- *State Transition:* The current state s_{α} transit to a new state s_{β} according to probability $P(s_{\beta} | s_{\alpha}, a)$, where a indicates a joint action set that includes the actions of all the UAVs.
- *Cost:* The immediate cost of the UAVs is $\frac{1}{IJ} \sum_{j=1}^J \sum_{i=1}^I AoI_j^i(s_{\alpha}, a)$.

C. MF-HPPO

The proposed MF-HPPO operates onboard at the UAVs to determine their trajectories and sensor selection. The UAV chooses a sensor and moves to it, then sends out a short beacon message with the ID of the chosen sensor. Upon the receipt of the beacon message, the selected sensor transmits its

data packets to the UAV, along with the state information of $AoI_j^i(t)$ in the control segment of the data packet. After the UAV correctly receives the data, it sends an acknowledgement to the ground sensor.

The following equation highlights the mean field idea of MF-HPPO [46]:

$$Q_i(s_{\alpha,i}, a) = \frac{1}{N_i} \sum_{k \in N(i)} Q_i(s_{\alpha,i}, a_i, a_k) = Q_i(s_{\alpha,i}, a_i, o_i). \quad (17)$$

Here, Q_i is the Q value of agent i , a represents the joint action of all agents. The neighbor agents of agent i are characterized by N_i . o_i is an indicator of the mean field. In essence, in multi-agent systems the Q value of an agent is computed based on the current state and joint action, but when we have a large number of agents computing joint action is impractical, therefore (17) allow an agent to compute its Q value just based on the mean field of its neighbors.

Fig. 2 shows the proposed MF-HPPO with LSTM layer, where each UAV equipped with the MF-HPPO to minimize the average AoI by optimizing the trajectory and data collection schedule. The use of the LSTM layer, continuous and discrete actors, and the objective function of PPO, are the features of the MF-HPPO in this diagram. As shown, The decision-making component of each agent consists of two actors and a critic, which is preceded by the LSTM layer to draw conclusions based on experience. The actor for continuous action spaces outputs continuous values for cruise control, such as position and velocity, and the actor for discrete action spaces outputs a categorical value that can be used to select one of the ground sensors. Each agent samples the actions and performs in the environment. The rollout buffer is filled with data generated by these interactions such as, state, mean field, action, cost and policy. As can be seen, we use Generalized Advantage Estimate (GAE) [47] as a sample-efficient method to estimate the advantage function. As depicted, based on the RolloutBuffer, mini-batches are then formed to train the LSTM and the actors and critics so that the agent can continuously improve its policies. The definition of the objective function of PPO is the total of actor losses and critic loss subtracted by entropy, as depicted in the diagram. The actor loss is inputted by the ratio of old policy and current policy and the advantage value. The critic loss is inputted by the critic's output and the return value. The policy is designed to encourage the agent to take advantageous actions, while punishing actions that deviate from the current policy.

Algorithm 1 summarizes the MF-HPPO with the LSTM-based characterization layer. In the initialization step, Input and Output are characterized; the algorithm receives parameters like Clip threshold, discount factor and mini-batch size as input and specify its output as trajectory and scheduling policy of UAV i . Next, the actor π_i and critic w_i are initialized with random weights for each agent. The number of training episodes is M , where the length of each episode is L . Each agent is trained using a predetermined set of iterations throughout

the learning phase. Sampling and optimization constitutes the learning phase. In the beginning of learning, the state $s_{\alpha,i}$ and mean field o_i are randomly initialized for each agent. With the start of the sampling policy, UAV i samples its action based on the policy θ_{old}^i . The sampled action represents sensor selection, velocity and locations, and executed in the environment to obtain the cost, new state and new mean field. Consequently, trajectories (i.e., sequence of states, actions, policy, mean field, and costs) are gathered and stored in the RolloutBuffer. In addition, GAE is applied to calculate the advantage that is used in (18). In the optimization step, the policies are optimized. In the optimization step, the policy parameter is updated for each epoch. The PPO objective is computed in each epoch according to the following equation:

$$L^{clip}(\theta^i) = \min \left(\frac{\pi_{\theta^i}(a_i | s_{\alpha,i}, o_i)}{\pi_{\theta_{old}^i}(a_i | s_{\alpha,i}, o_i)} A_{\pi_{\theta_{old}^i}}(s_{\alpha,i}, o_i, a_i), \right. \\ \left. g(\epsilon, A_{\pi_{\theta_{old}^i}}(s_{\alpha,i}, o_i, a_i)) \right) \quad (18)$$

where

$$\pi_{\theta^i}(a_i | s_{\alpha,i}, o_i) = \pi_{\theta^i}^c(a_i^c | s_{\alpha,i}, o_i) \pi_{\theta^i}^d(a_i^d | s_{\alpha,i}, o_i). \quad (19)$$

Here a_i^c and a_i^d correspond to actions in continuous and discrete spaces. In (19), to obtain the hybrid policy $\pi_{\theta^i}(a_i | s_{\alpha,i}, o_i)$, we multiply the policies for continuous and discrete actions [48]. Meanwhile, we assume that wireless radio of the UAV can cover the whole field.

Continuous policy $\pi_{\theta^i}^c$ is modeled using multivariate normal distribution and discrete policy $\pi_{\theta^i}^d$ is modeled using categorical distribution. In the next step, the overall objective function is optimized according to the following equation:

$$L^{total}(\theta^i) = L^{clip}(\theta^i) - K_1 L^{VF}(\theta^i) + K_2 * H. \quad (20)$$

Here, $L^{VF}(\theta^i)$ is the critic loss and H acts as a regularizer encourages the agent to execute actions more unpredictably for exploration and guard against the policy being overly deterministic. The entropy for continuous and discrete actions is computed based on the actions' distribution. We obtain the entropy by multiplication of the entropy of continuous and discrete action spaces to enable enforcing consistent regularization to both continuous and discrete action spaces. K_1 balances the importance of the critic loss and the actor loss, and K_2 coefficient controls the amount of entropy in the policy.

Finally, the sampling policy $\pi_{\theta_{old}^i}$ is updated with the policy π_{θ^i} , and the stored data are dropped. The next iteration then begins. The proposed MF-HPPO model is driven by DRL and can improve data aggregation by learning and refining the actions of cruise control and communication schedules on the fly. By this means, the MF-HPPO model can account for generic collision or obstacle avoidance via the offline training and can adapt to the specific real-world application scenario via online refinement. The UAVs can also be equipped with event cameras or utilize vision-based techniques to avoid collisions and adjust the flight behavior [49], [50].

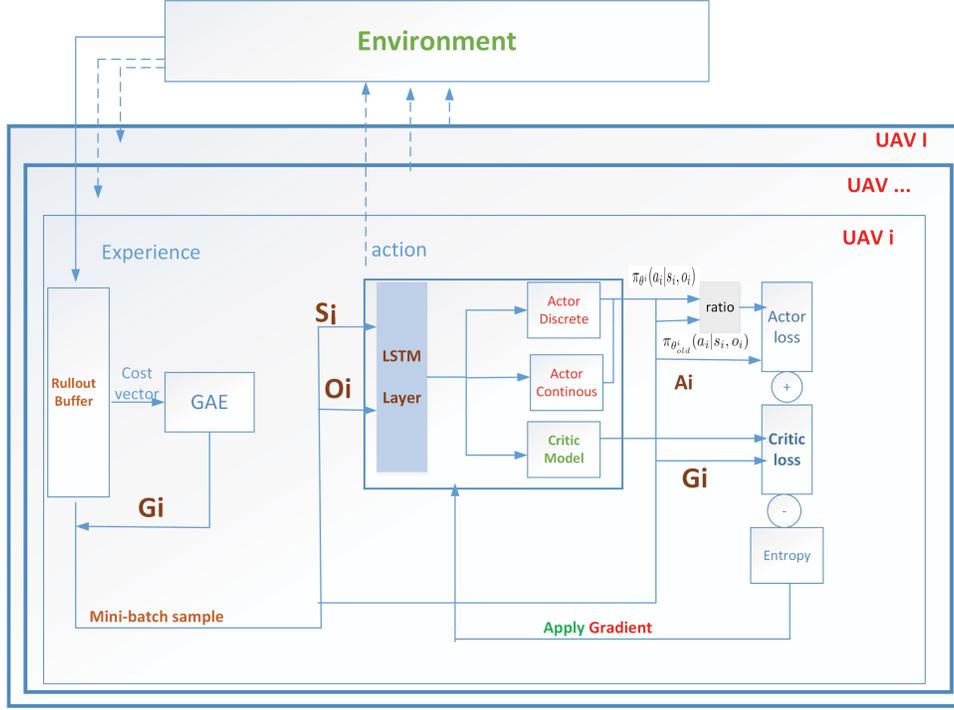


Fig. 2: An overview of MF-HPPO, where each UAV is equipped with the LSTM layer to optimize discrete and continuous actions using hybrid policy.

D. LSTM Layer

We further develop an LSTM layer in the proposed MF-HPPO, which captures long-term dependencies of time-varying network state s_α . Cell memory and the gating mechanism are main components of LSTM. Cell memory is responsible to store the summary of the past input data and the gating mechanism regulates the information flow between the input, output, and cell memory. The network states are fed into LSTM one by one (one at each step). The last hidden state κ_i^{hidd} is returned as the output of the state characterization layer. Each agent uses an LSTM layer to predict their respective hidden states. The hidden states κ_i^{hidd} are calculated by the following composite function:

$$\kappa_i^{hidd} = out_i \tanh(C_i), \quad (21)$$

$$out_i = \sigma(W_0 \cdot [C_i, \kappa_{i-1}^{hidd}, A_i] + e_i), \quad (22)$$

$$C_i = F_i C_{i-1} + p_i \tanh(W_c \cdot [\kappa_{i-1}^{hidd}, A_i] + e_c), \quad (23)$$

$$F_i = \sigma(W_f \cdot [\kappa_{i-1}^{hidd}, C_{i-1}, A_i] + e_f), \quad (24)$$

$$p_i = \sigma(W_p \cdot [\kappa_{i-1}^{hidd}, C_{i-1}, A_i] + e_p), \quad (25)$$

where the output gate, cell activation vectors, forget gate, and input gate of the LSTM layer are denoted by out_i , C_i , F_i , and p_i , respectively. σ and \tanh correspond to logistic sigmoid function and the hyperbolic tangent function, respectively. W_0, W_c, W_f, W_p are the weight matrix, and e_0, e_c, e_f, e_p are the bias matrix [51], [52].

E. Complexity and Convergence of MF-HPPO

The overall complexity of MF-HPPO is calculated as follows, $O(I \cdot ML \cdot (\sum_{g=1}^G n_{g-1} \cdot n_g))$ where n_g is the number of neural units in the g -th hidden layer. In this work, the PPO architecture is built with the same n_g in all hidden layers. Therefore, the PPO complexity can be reduced to $O(I \cdot ML \cdot (g-1) \cdot n_g^2) = O(I \cdot ML \cdot n_g^2)$. The convergence analysis is proved by simulation results (see Fig. 4).

VI. NUMERICAL RESULTS AND DISCUSSIONS

A. Implementation of MF-HPPO

MF-HPPO is implemented in Python 3.8 using Pytorch (the Python deep learning library). A Predator Workstation running 64-bit Ubuntu 20.04 LTS, with Intel Core i7-11370 H CPU @ 3.30 GHz 8 and 16 GB memory is used for the Pytorch setup. Table ?? clearly outlines the different considered simulation parameters. MF-HPPO algorithm is trained over 3000 episodes with 40 steps each. The discount factor and learning rate are set to 0.99 and $3e-4$, respectively. Each agent comprises the input layer, LSTM layer, the critic and actors with fully-connected hidden layers of size 256 and output layer. Each neuron uses Rectified Linear Unit (ReLU) as an activation function. In addition, Hyperbolic tangent (\tanh) and softmax are used as activation functions in the output layer of the continuous actor-network and discrete actor network. The input of each critic network is represented as a concatenation of states and mean field, and its output is a scalar that assesses the states according to the global policy. The total log probability of the hybrid

Algorithm 1: MF-HPPO Characterized by LSTM Layer

```

1 1.Initialize
  Input: Clip threshold  $\epsilon$ , discount factor  $\gamma$ , learning rate
            $\eta$ , buffer size  $D$ , mini-batch size  $B$ 
  Output: The scheduled ground sensor  $j$  and trajectory
             $\zeta_i$  of UAV  $i$ 
2 Randomly initialize the Actors  $\pi_i$  and Critics  $w_i$  with
   networks parameters  $\theta^i$ 
3 The LSTM layer with  $\{W_o, W_c, W_f, W_p\}$  and
    $\{e_o, e_c, e_f, e_p\}$ .
4 Initialize the sampling policy  $\pi_{\theta_{old}^i}$  with  $\theta_{old}^i \leftarrow \theta^i$ .
5  $\forall i \in (1, I)$ 
6 2.Learning
7 for  $episode=1$  to  $M$  do
8   Randomly obtain the initial state  $s_{\alpha,i}$ 
9   for  $t = 1$  to  $L$  do
10     *The sampling phase*
11     Sample: Sample action
12      $a_i \sim \pi_{\theta_{old}^i}(a_i | s_{\alpha,i}, o_i, \theta^i)$ ;
13     Execute the action  $a_i$  that specifies the
14     scheduled ground sensor  $j$  and trajectory  $\zeta_i$  of
15     UAV  $i$ .
16     Obtain the cost and new state  $s_{\beta,i}$  and new
17     mean field  $\phi_i(t+1)$ .
18     RolloutBuffer: store the trajectory
19      $(s_{\alpha,i}, a_i, c, o_i, \pi_{\theta_{old}^i}(a_i | s_{\alpha,i}, o_i, \theta^i))$ 
20      $s_{\alpha,i} = s_{\beta,i}$ 
21   end for
22   Compute the advantage using GAE
23   for  $epoch = 1$  to  $P$  do
24     *The optimization phase*
25     Sample the RolloutBuffer
26     Compute the PPO-Clip objective function using
27     (18)
28     Compute the critic loss.
29     Optimize the overall objective function using
30     (20)
31   end for
32   Synchronize the sampling policy  $\pi_{\theta_{old}^i} \leftarrow \pi_{\theta^i}$ 
33   Drop the stored data in RolloutBuffer.
34 end for

```

policy is the sum of the log probabilities of the continuous and discrete action spaces. This log probability would be used as part of the calculation of the objective function in MF-HPPO, along with the estimated cost and the entropy regularization term.

B. Baseline Description

The MF-HPPO characterized with LSTM layer is compared by single-agent PPO, random scheduling and trajectory design (RSTD), multi-agent DQN (MADQN) and MF-HPPO without

Table II: PyTorch Configuration

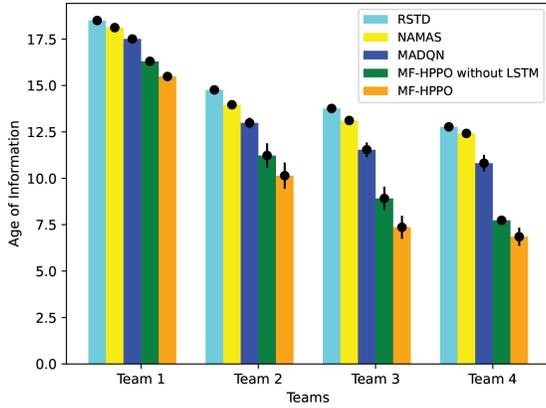
Parameters	Values
Number of ground sensors	100
Number of UAVs	30
Geographical area size [m]	1,000*1,000
Altitude of the UAVs	120 m
Critic Network Learning Rate	3e-4
Actor-Network Learning Rate	3e-4
Number of Hidden Layers for Networks	2
Number of Neurons	256
Loss Coefficients for K_1 and K_2	0.2 and 3
Optimizer Technique	Adam
Clip Fraction	0.2
Rollout Buffer size	40
Batch size	40
Mini Batch Size	4
PPO Epochs	8
Number of episodes	3,000
Discount Factor	0.99
Minimum speed	0 m/s
Maximum speed	15 m/s

LSTM Layer. A brief introduction of the four benchmarks is given below

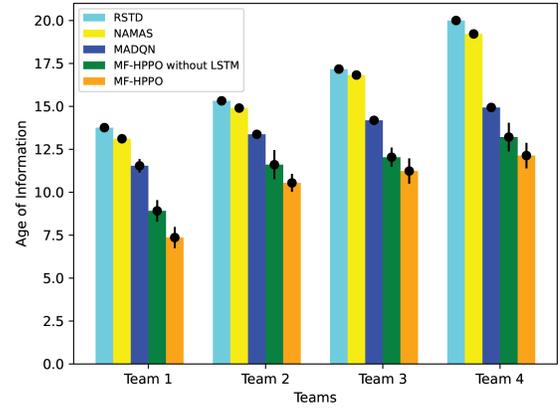
- 1) NAMAS [53], in this algorithm, each agent schedules the neighbors with maximum AoI to minimize the average AoI.
- 2) RSTD, in this algorithm transmission scheduling and trajectory design, are randomly designed.
- 3) MADQN, in this algorithm, each agent running DQN cooperates to reduce average AoI following circular trajectories.
- 4) MF-HPPO without LSTM Layer, the structure of this algorithm is the same as MF-HPPO but without LSTM layer.

C. Performance analysis of MF-HPPO

Fig. 3 depicts the performance evaluation of MF-HPPO in comparison to the baselines by changing the number of UAVs and ground sensors. Fig. 3a illustrates the influence of varying the quantity of UAVs on the AoI. It is observed that an increase of UAVs leads to a reduction in AoI, attributable to enhanced time efficiency and the capacity for quicker operation of more ground sensors. Specifically, augmenting the UAV count from 1 to 30 results in a 61% reduction in the average AoI for the MF-HPPO algorithm, in contrast to a 37% reduction observed for MADQN. This disparity is attributed to the fact that MF-HPPO executes optimization within a mixed action space, demonstrating greater training stability compared to MADQN, which utilizes circular trajectories. Furthermore, it is noteworthy that MF-HPPO significantly surpasses the performance of both random assignment and NAMAS strategies. Fig. 3b evaluates the average AoI given 20 UAVs and groups of 100, 200, 300, and 400 ground sensors. The MADQN, NAMAS, and the RSTD are used as baselines. Overall, increasing the number of ground sensors results in a uniform increase in the average



(a) Evaluation of MF-HPPO’s performance with a variable number of UAVs in comparison to RSTD, MADQN and MF-HPPO without LSTM



(b) Evaluation of MF-HPPO’s performance with a variable number of ground sensors in comparison to RSTD, MADQN and MF-HPPO without LSTM

Fig. 3: Performance evaluation of MF-HPPO by changing the number of UAVs and ground sensors

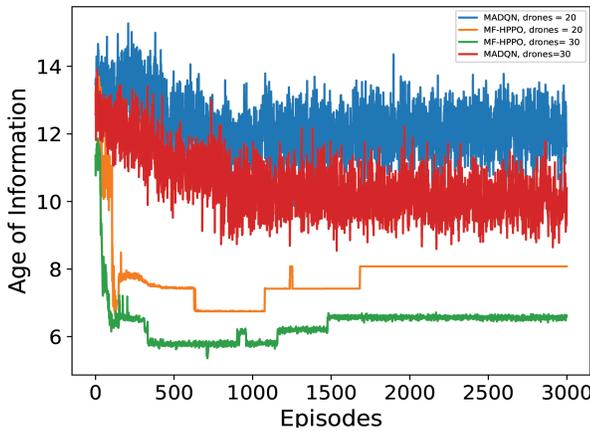


Fig. 4: The network cost for each episode of MF-HPPO with $I=30$ and benchmarks

AoI, since more sensor data should be collected. In particular, when the number of ground sensors is 400, the proposed MF-HPPO outperforms the RSTD by 38%, NAMAS by 33%, and the MADQN by 17%.

Fig. 4 captures the convergence trend of the MF-HPPO algorithm, which was assessed by deploying 20 UAVs to service 100 ground sensors. In this context, the MF-HPPO model with $I=30$ settings demonstrates a significantly lower Age of Information (AoI) when compared to the MADQN algorithm with $I=20$ and $I=30$, exhibiting improvements of 38% and 43%, respectively. This enhanced performance is attributed to the optimized trajectories and scheduling for data collection by multiple UAVs, resulting in superior time efficiency. Additionally, the integration of an LSTM layer in the MF-HPPO framework contributes to both an acceleration

and stabilization of convergence. Notably, the peak AoI for the proposed MF-HPPO model decreases dramatically from 14 seconds to 6 seconds within the initial 1,000 episodes. Subsequently, between episodes 1,500 and 3,000, the AoI stabilizes at around 7 seconds, with only minor fluctuations observed

MF-HPPO-generated trajectories for 20 UAVs are shown in Fig. 5, where the ground sensor distribution patterns are uniform, square, or normal ones. When designing trajectories for AoI minimization, the UAVs’ trajectories are impacted by the distribution of the ground sensors. The UAV needs to approach to the location of each scheduled sensor to collect the data and update its AoI. Fig. 5(a), refer to the normal distribution and shows trajectories for 20 UAVs, focusing on the center area of the ground sensors and less on the corners. The normal distribution of the ground sensors can affect the UAVs’ trajectories by determining which ground sensors are prioritized for data collection. For example, as can be seen, most ground sensors are centered and their data may become stale, in this case, the UAVs’ trajectories are designed to visit these ground sensors more frequently to minimize the average AoI. Figs. 5(b) is related to the square distribution. As can be seen, the ground sensors are less centered. This cause diverse set of ground sensors in wider range to be covered in comparison to normal distribution. Fig. 5(c) refer to the uniform distribution. As can be seen, the UAVs design wide-area trajectories due to the wider distribution of ground sensors covering the entire area and the AoI requirements of the scattered ground sensors.

Fig. 6 demonstrates the convergence figures for two variants of MF-HPPO by changing the clip threshold. PPO uses the clip threshold, commonly referred to as epsilon, to regulate the amount of policy updating. A larger clip threshold allows for more aggressive updating, while a smaller clip threshold restricts updating more severely, resulting in less policy change.

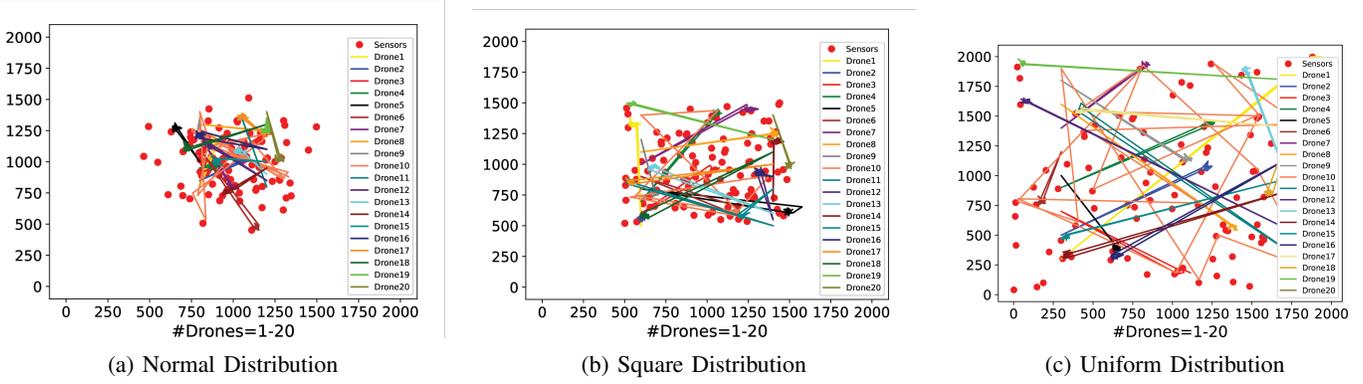


Fig. 5: MF-HPPO trajectory distributions for various UAV counts and ground sensor distributions.

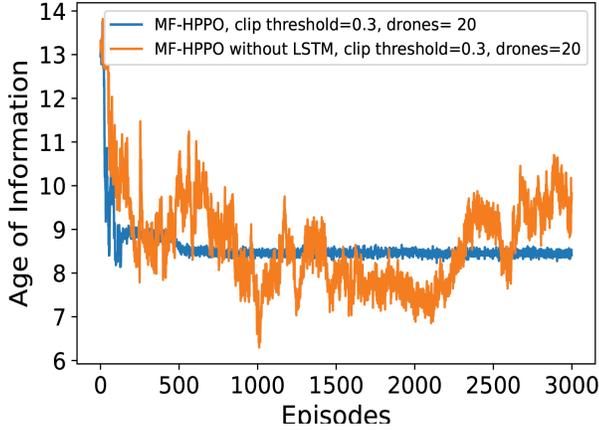


Fig. 6: Performance evaluation of MF-HPPO by changing clip threshold

The blue curve shows the MF-HPPO with LSTM layer and a clip threshold of 0.3 outperforming the MF-HPPO without LSTM layer clip threshold 0.3. The latter shows a deviating behavior due to the influence of the clip threshold, while the blue curve shows an absolutely stable trend despite the same value of the clip threshold thanks to the LSTM layer. Overall, adding the LSTM layer to MF-HPPO can stabilize the training and prevent divergence of the strategies.

VII. CONCLUSION

In this paper, we propose a mean field flight resource allocation to model velocity control for a swarm of UAVs, in which each UAV minimizes the average AoI by considering the collective behavior of others. Due to the high computational complexity of MFG, we leverage AI and propose MF-HPPO characterized with an LSTM layer to optimize the UAV trajectories and data collection scheduling in mixed action space. Simulation results based on PyTorch deep learning library show that the proposed MF-HPPO for UAVs-assisted sensor networks reduces average AoI by up to 57% and 45%, as compared to

existing non-learning random algorithm and MADQN method (which performs the action of trajectory planning in the discrete space), respectively. This confirms the AI-enhanced mean field resource allocation is a practical solution for minimizing AoI in UAV swarms.

APPENDIX

PROOF OF FPK EQUATION (12) FOR CRUISE CONTROL

We derive the mean field via an arbitrary test function $g(\zeta)$, which is a twice continuously differentiable compactly supported function of the state space. The integral of $m(\zeta)g(\zeta)d\zeta$ can be considered as the continuum limit of the sum $g(\zeta(t))$, where $\zeta(t)$ is the UAV's state at time t . It is known that,

$$\int m(\zeta(t))g(\zeta)d\zeta = \frac{1}{N} \sum_{i=1}^N g(\zeta(t)). \quad (26)$$

At time t , the first-order differential function with regard to time t is derived to check how this integral varies in time. By utilizing the chain rule, we can derive the heuristic formula as

$$\int \partial_t m(\zeta(t))g(\zeta)d\zeta = \frac{1}{N} \sum_{i=1}^N \partial_t \zeta(t) \nabla g(\zeta(t)) + \partial_t^2 \zeta(t) \nabla^2 g(\zeta(t)). \quad (27)$$

Taking the limit of the right side of the above equation when N tends to infinity, we get

$$\int [\partial_t m(\zeta(t)) + \nabla_{\zeta} m(\zeta(t)) \cdot \frac{\partial \zeta}{\partial t} - \frac{\eta^2}{2} \nabla_{\zeta}^2 m(\zeta(t))] g(\zeta(t)) d\zeta = 0, \quad (28)$$

for any test function g through integration by parts. Then the above equation leads to the following equation:

$$\partial_t m(\zeta(t)) + \nabla_{\zeta} m(\zeta(t)) \cdot v(t) - \frac{\sigma^2}{2} \nabla_{\zeta}^2 m(\zeta(t)) = 0. \quad (29)$$

which correspond to FPK equation defined in (12).

REFERENCES

- [1] A. Heidari, N. Jafari Navimipour, M. Unal, and G. Zhang, "Machine learning applications in internet-of-drones: Systematic review, recent deployments, and open issues," *ACM Comput. Surv.*, vol. 55, no. 12, Mar. 2023.
- [2] K. Li, R. C. Voicu, S. S. Kanhere, W. Ni, and E. Tovar, "Energy efficient legitimate wireless surveillance of uav communications," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 3, pp. 2283–2293, Jan. 2019.
- [3] K. Li, W. Ni, E. Tovar, and M. Guizani, "Joint flight cruise control and data collection in uav-aided internet of things: An onboard deep reinforcement learning approach," *IEEE Internet of Things Journal*, vol. 8, no. 12, pp. 9787–9799, Aug. 2020.
- [4] K. Li, W. Ni, X. Wang, R. P. Liu, S. S. Kanhere, and S. Jha, "Energy-efficient cooperative relaying for unmanned aerial vehicles," *IEEE Transactions on Mobile Computing*, vol. 15, no. 6, pp. 1377–1386, Aug. 2016.
- [5] K. Li, W. Ni, E. Tovar, and A. Jamalipour, "On-board deep q-network for uav-assisted online power transfer and data collection," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 12, pp. 12 215–12 226, Oct. 2019.
- [6] S. Guan, Z. Zhu, and G. Wang, "A review on uav-based remote sensing technologies for construction and civil applications," *Drones*, vol. 6, no. 5, p. 117, May 2022.
- [7] D. K. Villa, A. S. Brandao, and M. Sarcinelli-Filho, "A survey on load transportation using multicopter uavs," *Journal of Intelligent & Robotic Systems*, vol. 98, pp. 267–296, Oct. 2020.
- [8] K. Li, W. Ni, Y. Emami, and F. Dressler, "Data-driven flight control of internet-of-drones for sensor data aggregation using multi-agent deep reinforcement learning," *IEEE Wireless Communications*, vol. 29, no. 4, pp. 18–23, Aug. 2022.
- [9] K. Li, W. Ni, and F. Dressler, "Continuous maneuver control and data capture scheduling of autonomous drone in wireless sensor networks," *IEEE Transactions on Mobile Computing*, vol. 21, no. 8, pp. 2732–2744, Jan. 2021.
- [10] K. Li, W. Ni, X. Yuan, A. Noor, and A. Jamalipour, "Deep-graph-based reinforcement learning for joint cruise control and task offloading for aerial edge internet of things (edgeiot)," *IEEE Internet of Things Journal*, vol. 9, no. 21, pp. 21 676–21 686, Jun. 2022.
- [11] S. Kaul, R. Yates, and M. Gruteser, "Real-time status: How often should one update?" in *Proceedings IEEE INFOCOM*, Mar. 2012, pp. 2731–2735.
- [12] K. Li, W. Ni, A. Noor, and M. Guizani, "Employing intelligent aerial data aggregators for the internet of things: Challenges and solutions," *IEEE Internet of Things Magazine*, vol. 5, no. 1, pp. 136–141, Mar. 2022.
- [13] M. E. Mkiramweni, C. Yang, J. Li, and W. Zhang, "A survey of game theory in unmanned aerial vehicles communications," *IEEE Communications Surveys & Tutorials*, vol. 21, no. 4, pp. 3386–3416, May 2019.
- [14] D. Chen, Q. Qi, Z. Zhuang, J. Wang, J. Liao, and Z. Han, "Mean field deep reinforcement learning for fair and efficient uav control," *IEEE Internet of Things Journal*, vol. 8, no. 2, pp. 813–828, Jul. 2020.
- [15] L. Li, Q. Cheng, K. Xue, C. Yang, and Z. Han, "Downlink transmit power control in ultra-dense uav network based on mean field game and deep reinforcement learning," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 12, pp. 15 594–15 605, Dec. 2020.
- [16] D. Shi, H. Gao, L. Wang, M. Pan, Z. Han, and H. V. Poor, "Mean field game guided deep reinforcement learning for task placement in cooperative multiaccess edge computing," *IEEE Internet of Things Journal*, vol. 7, no. 10, pp. 9330–9340, Mar. 2020.
- [17] Y. Sun, L. Li, Q. Cheng, D. Wang, W. Liang, X. Li, and Z. Han, "Joint trajectory and power optimization in multi-type uavs network with mean field q-learning," in *IEEE International Conference on Communications Workshops (ICC Workshops)*, Dublin, Ireland, Jul. 2020.
- [18] M. Wang, L. Li, W. Lin, B. Wei, W. Chen, and Z. Han, "Uav position optimization based on information freshness: A mean field game approach," in *13th International Conference on Wireless Communications and Signal Processing (WCSP)*, Changsha, China, Dec. 2021.
- [19] K. Xue, Z. Zhang, L. Li, H. Zhang, X. Li, and A. Gao, "Adaptive coverage solution in multi-uavs emergency communication system: a discrete-time mean-field game," in *International Wireless Communications & Mobile Computing Conference (IWCMC)*, Limassol, Cyprus, Aug. 2018, pp. 1059–1064.
- [20] Y. Xu, L. Li, Z. Zhang, K. Xue, and Z. Han, "A discrete-time mean field game in multi-uav wireless communication systems," in *IEEE/CIC International Conference on Communications in China (ICCC)*, Beijing, China, Feb. 2018, pp. 714–718.
- [21] H. Gao, W. Lee, Y. Kang, W. Li, Z. Han, S. Osher, and H. V. Poor, "Energy-efficient velocity control for massive numbers of uavs: A mean field game approach," *IEEE Transactions on Vehicular Technology*, vol. 71, no. 6, pp. 6266–6278, Mar. 2022.
- [22] X. Liu, B. Lai, B. Lin, and V. C. Leung, "Joint communication and trajectory optimization for multi-uav enabled mobile internet of vehicles," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 9, pp. 15 354–15 366, Jan. 2022.
- [23] X. Liu, Z. Liu, B. Lai, B. Peng, and T. S. Durrani, "Fair energy-efficient resource optimization for multi-uav enabled internet of things," *IEEE Transactions on Vehicular Technology*, vol. 72, no. 3, pp. 3962–3972, Nov. 2022.
- [24] O. S. Oubbati, M. Atiquzzaman, H. Lim, A. Rachedi, and A. Lakas, "Synchronizing uav teams for timely data collection and energy transfer by deep reinforcement learning," *IEEE Transactions on Vehicular Technology*, vol. 71, no. 6, pp. 6682–6697, Apr. 2022.
- [25] K. Chi, F. Li, F. Zhang, M. Wu, and C. Xu, "Aoi optimal trajectory planning for cooperative uavs: A multi-agent deep reinforcement learning approach," in *IEEE International Conference on Electronic Information and Communication Technology (ICEICT)*, Hefei, China, Aug. 2022, pp. 57–62.
- [26] J. Hu, H. Zhang, K. Bian, L. Song, and Z. Han, "Distributed trajectory design for cooperative internet of uavs using deep reinforcement learning," in *IEEE Global Communications Conference (GLOBECOM)*, Waikoloa, HI, Dec. 2019.
- [27] J. Hu, H. Zhang, L. Song, R. Schober, and H. V. Poor, "Cooperative internet of uavs: Distributed trajectory design by multi-agent deep reinforcement learning," *IEEE Transactions on Communications*, vol. 68, no. 11, pp. 6807–6821, Aug. 2020.
- [28] M. Samir, C. Assi, S. Sharafeddine, and A. Ghrayeb, "Online altitude control and scheduling policy for minimizing aoi in uav-assisted iot wireless networks," *IEEE Transactions on Mobile Computing*, vol. 21, no. 7, pp. 2493–2505, Dec. 2022.
- [29] M. Sun, X. Xu, X. Qin, and P. Zhang, "Aoi-energy-aware uav-assisted data collection for iot networks: A deep reinforcement learning method," *IEEE Internet of Things Journal*, vol. 8, no. 24, pp. 17 275–17 289, May 2021.
- [30] X. Li, X. Xu, J. Huo, and W. Huangfu, "Aoi minimization in uav-assisted iot network: A reinforcement learning approach," in *International Conference on Ubiquitous Communication (Ucom)*, Jul. 2023, pp. 315–320.
- [31] E. Eldeeb, D. E. Pérez, J. Michel de Souza Sant'Ana, M. Shehab, N. H. Mahmood, H. Alves, and M. Latva-Aho, "A learning-based trajectory planning of multiple uavs for aoi minimization in iot networks," in *Joint European Conference on Networks and Communications & 6G Summit (EuCNC/6G Summit)*, Grenoble, France, Jun. 2022, pp. 172–177.
- [32] M. A. Abd-Elmagid, A. Ferdowsi, H. S. Dhillon, and W. Saad, "Deep reinforcement learning for minimizing age-of-information in uav-assisted networks," in *IEEE Global Communications Conference (GLOBECOM)*, Waikoloa, HI, Dec. 2019.
- [33] C. Zhou, H. He, P. Yang, F. Lyu, W. Wu, N. Cheng, and X. Shen, "Deep rl-based trajectory planning for aoi minimization in uav-assisted iot," in *International Conference on Wireless Communications and Signal Processing (WCSP)*, Xi'an, China, Oct. 2019.
- [34] P. Tong, J. Liu, X. Wang, B. Bai, and H. Dai, "Deep reinforcement learning for efficient data collection in uav-aided internet of things," in *IEEE International Conference on Communications Workshops (ICC Workshops)*, Dublin, Ireland, Jun. 2020.
- [35] L. Liu, K. Xiong, J. Cao, Y. Lu, P. Fan, and K. B. Letaief, "Average aoi minimization in uav-assisted data collection with rf wireless power transfer: A deep reinforcement learning scheme," *IEEE Internet of Things Journal*, vol. 9, no. 7, pp. 5216–5228, Sep. 2021.
- [36] E. Eldeeb, J. M. d. S. Sant'Ana, D. E. Pérez, M. Shehab, N. H. Mahmood, and H. Alves, "Multi-uav path learning for age and power optimization in iot with uav battery recharge," *IEEE Transactions on Vehicular Technology*, vol. 72, no. 4, pp. 5356–5360, Nov. 2023.
- [37] W. Liu, D. Li, T. Liang, T. Zhang, Z. Lin, and N. Al-Dhahir, "Joint trajectory and scheduling optimization for age of synchronization minimization

- in uav-assisted networks with random updates,” *IEEE Transactions on Communications*, vol. 71, no. 11, pp. 6633–6646, Jul. 2023.
- [38] E. Eldeeb, M. Shehab, and H. Alves, “Age minimization in massive iot via uav swarm: A multi-agent reinforcement learning approach,” in *IEEE 34th Annual International Symposium on Personal, Indoor and Mobile Radio Communications (PIMRC)*, Oct. 2023.
- [39] X. Wang, M. Yi, J. Liu, Y. Zhang, M. Wang, and B. Bai, “Cooperative data collection with multiple uavs for information freshness in the internet of things,” *IEEE Transactions on Communications*, vol. 71, no. 5, pp. 2740–2755, Mar. 2023.
- [40] A. Al-Hourani, S. Kandeepan, and S. Lardner, “Optimal lap altitude for maximum coverage,” *IEEE Wireless Communications Letters*, vol. 3, no. 6, pp. 569–572, Jul. 2014.
- [41] Y. Emami, B. Wei, K. Li, W. Ni, and E. Tovar, “Joint communication scheduling and velocity control in multi-uav-assisted sensor networks: A deep reinforcement learning approach,” *IEEE Transactions on Vehicular Technology*, vol. 70, no. 10, pp. 10986–10998, 2021.
- [42] B. Waldrip, V. Prain, and P. Sellings, “Explaining newton’s laws of motion: Using student reasoning through representations to develop conceptual understanding,” *Instructional Science*, vol. 41, pp. 165–189, Mar. 2013.
- [43] P. Mörters and Y. Peres, *Brownian motion*. Cambridge University Press, 2010, vol. 30.
- [44] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, “Trust region policy optimization,” in *International conference on machine learning*. PMLR, 2015, pp. 1889–1897.
- [45] M. Samir, C. Assi, S. Sharafeddine, and A. Ghayeb, “Online altitude control and scheduling policy for minimizing aoi in uav-assisted iot wireless networks,” *IEEE Transactions on Mobile Computing*, vol. 21, no. 7, pp. 2493–2505, Dec. 2020.
- [46] Y. Yang, R. Luo, M. Li, M. Zhou, W. Zhang, and J. Wang, “Mean field multi-agent reinforcement learning,” in *International Conference on Machine Learning*. PMLR, 2018, pp. 5571–5580.
- [47] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel, “High-dimensional continuous control using generalized advantage estimation,” *arXiv preprint arXiv:1506.02438*, Jun. 2015.
- [48] M. Neunert, A. Abdolmaleki, M. Wulfmeier, T. Lampe, T. Springenberg, R. Hafner, F. Romano, J. Buchli, N. Heess, and M. Riedmiller, “Continuous-discrete reinforcement learning for hybrid control in robotics,” in *Conference on Robot Learning*. PMLR, May 2020, pp. 735–751.
- [49] H. Huang, G. Zhu, Z. Fan, H. Zhai, Y. Cai, Z. Shi, Z. Dong, and Z. Hao, “Vision-based distributed multi-uav collision avoidance via deep reinforcement learning for navigation,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Oct. 2022, pp. 13745–13752.
- [50] J. N. Yasin, S. A. S. Mohamed, M.-H. Haghbayan, J. Heikkonen, H. Tenhunen, and J. Plosila, “Unmanned aerial vehicles (uavs): Collision avoidance systems and approaches,” *IEEE Access*, vol. 8, pp. 105139–105155, 2020.
- [51] K. Li, W. Ni, and F. Dressler, “Lstm-characterized deep reinforcement learning for continuous flight control and resource allocation in uav-assisted sensor network,” *IEEE Internet of Things Journal*, vol. 9, no. 6, pp. 4179–4189, Aug. 2022.
- [52] J. Zheng, K. Li, N. Mhaisen, W. Ni, E. Tovar, and M. Guizani, “Exploring deep-reinforcement-learning-assisted federated learning for online resource allocation in privacy-preserving edgeiot,” *IEEE Internet of Things Journal*, vol. 9, no. 21, pp. 21099–21110, May 2022.
- [53] M. A. Abd-Elmagid, A. Ferdowsi, H. S. Dhillon, and W. Saad, “Deep reinforcement learning for minimizing age-of-information in uav-assisted networks,” in *IEEE Global Communications Conference (GLOBECOM)*, Dec. 2019.