

# UAV Trajectory Planning in Wireless Sensor Networks for Energy Consumption Minimization by Deep Reinforcement Learning

Botao Zhu<sup>1</sup>, Ebrahim Bedeer<sup>2</sup>, *Member, IEEE*, Ha H. Nguyen<sup>1</sup>, *Senior Member, IEEE*, Robert Barton<sup>3</sup>, *Member, IEEE*, and Jerome Henry, *Senior Member, IEEE*

**Abstract**—Unmanned aerial vehicles (UAVs) have emerged as a promising candidate solution for data collection of large-scale wireless sensor networks (WSNs). In this paper, we investigate a UAV-aided WSN, where cluster heads (CHs) receive data from their member nodes, and a UAV is dispatched to collect data from CHs. We aim to minimize the total energy consumption of the UAV-WSN system in a complete round of data collection. Toward this end, we formulate the energy consumption minimization problem as a constrained combinatorial optimization problem by jointly selecting CHs from clusters and planning the UAV's visiting order to the selected CHs. The formulated energy consumption minimization problem is NP-hard, and hence, hard to solve optimally. To tackle this challenge, we propose a novel deep reinforcement learning (DRL) technique, pointer network-A\* (Ptr-A\*), which can efficiently learn the UAV trajectory policy for minimizing the energy consumption. The UAV's start point and the WSN with a set of pre-determined clusters are fed into the Ptr-A\*, and the Ptr-A\* outputs a group of CHs and the visiting order of CHs, i.e., the UAV's trajectory. The parameters of the Ptr-A\* are trained on small-scale clusters problem instances for faster training by using the actor-critic algorithm in an unsupervised manner. Simulation results show that the trained models based on 20-clusters and 40-clusters have a good generalization ability to solve the UAV's trajectory planning problem in WSNs with different numbers of clusters, without retraining the models. Furthermore, the results show that our proposed DRL algorithm outperforms two baseline techniques.

**Index Terms**—Combinatorial optimization, deep reinforcement learning, trajectory planning, UAV, WSN.

## I. INTRODUCTION

THE use of unmanned aerial vehicles (UAVs) has recently attracted a lot of attention from both the research community and industry. UAVs have been used for a variety of purposes [1], such as environmental monitoring, mobile cloud

computing, disaster management, security operations, and wireless power transfer, to name a few. The popularity and widespread applications of UAVs are due to their many advantages, such as cost-effectiveness, having line-of-sight (LoS) links with the ground devices, mobility, and reliable network access [2].

In order to overcome the limitations of the traditional wireless sensor networks (WSNs), extensive research has been done on the integration of UAVs and WSNs for long-distance mission communications where UAVs are considered as mobile sinks for receiving data from cluster heads (CHs), and then, they can transmit the collected data to terrestrial BSs for further processing. Using UAVs as mobile sinks can reduce the energy consumption of ground nodes in WSNs when compared to the traditional multi-hop WSNs which transmit data from each node to the sink node over a long distance or several hops [3]. Despite the advantages, the integration of UAVs and WSNs still grapples with many challenges. First of all, due to the limited energy source carried by UAVs, the service range of UAVs is constrained by the reality that they cannot travel very long distances or fly for long periods of time. Second, battery life of sensor nodes in WSNs is typically limited, and in many cases it is hard to regularly replace their batteries. As a result, frequent communication with UAVs can cause sensor nodes to exhaust their energy rapidly. Hence, it is important to study the energy saving problem in UAV-enabled WSNs.

## A. Motivation

Prior works on the energy consumption minimization for UAV-enabled WSNs can be classified into three categories depending on the objectives. The first category only considers minimizing the UAV's energy consumption, e.g., [4], [5]. In contrast, the second category considers only minimizing the energy consumption of the ground devices in the UAV-aided wireless networks, e.g., [6], [7]. In the third category, the energy of both the UAV and the ground devices are taken into account when minimizing the energy consumption of the UAV-enabled system, e.g., [8], [9]. However, most of the aforementioned studies assume the UAV directly communicates with each device of the ground wireless network. In this case, if the UAV flies over all devices in a large-scale WSN, it would lead to a long flight trajectory for the UAV which increase its energy consumption.

Manuscript received January 14, 2021; revised May 18, 2021; accepted July 25, 2021. Date of publication August 4, 2021; date of current version September 17, 2021. This work was supported by an NSERC/Cisco Industrial Research Chair in Low-Power Wireless Access for Sensor Networks. The review of this article was coordinated by Dr. Feng Ye. (Corresponding author: Botao Zhu.)

Botao Zhu, Ebrahim Bedeer, and Ha H. Nguyen are with the Department of Electrical and Computer Engineering, University of Saskatchewan, Saskatoon S7N 5A2, Canada (e-mail: botao.zhu@usask.ca; e.bedeer@usask.ca; ha.nguyen@usask.ca).

Robert Barton and Jerome Henry are with the Cisco Systems Inc, San Jose, California, CA 95134 USA (e-mail: robbarto@cisco.com; jerhenry@cisco.com).

Digital Object Identifier 10.1109/TVT.2021.3102161

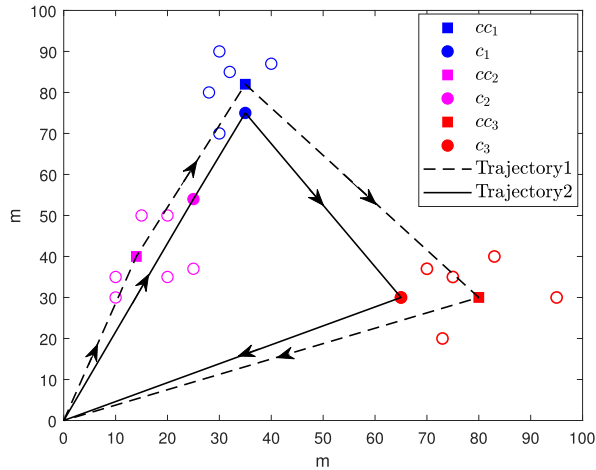


Fig. 1. Comparison of different UAV's trajectories.

As a result, the UAV may run out of its energy in flight or may need to recharge its battery frequently.

Motivated by the aforementioned works, in this paper we investigate the problem of minimizing the total energy consumption of the UAV and the ground devices in a clustered WSN, which has not been well researched in prior works. We assume that devices on the ground have been clustered according to some specific criterion, e.g., based on their geographical locations; hence, clustering techniques for WSNs will not be discussed in this paper. In each pre-determined cluster in our system model, one of the ground devices will be selected as the CH, which is responsible for collecting data from the non-CH devices in the same cluster. Hence, the UAV only needs to visit a set of CHs for gathering data along the planned trajectory that is determined by locations of the ground CHs. The selection of CHs affects both the energy consumption of the ground devices and the UAV, as shown in the following illustrative example.

Consider the case where nodes of three clusters are deployed in a given area as shown in Fig. 1. There are two candidate solutions of potential CHs. One possible solution, i.e., trajectory 1, selects the “center” node of each cluster as the CH, such as  $cc_1, cc_2, cc_3$ , and the other possible solution, i.e., trajectory 2, selects non-center nodes as the CHs, e.g.  $c_1, c_2, c_3$ . The start/end location of the UAV is (0, 0). If the UAV chooses to follow trajectory 1, the energy consumption between CHs and their member nodes will be minimal because the Euclidean distances between member nodes and their CHs are minimal (on average) per cluster [10]. However, this will lead to an increase in energy consumption of the UAV as its flight trajectory may be longer, and hence, may not be optimal from the start point to the end point, which can be seen from trajectory 1 in Fig. 1. On the contrary, if the UAV goes through trajectory 2, it will consume less energy in flight because trajectory 2 is shorter than trajectory 1. But, the communication energy consumption between CHs and their member nodes will be higher in this case. From the above simple example and discussion, it is clearly important and relevant to study the energy-efficient UAV's trajectory planning in clustered WSNs in order to minimize

the overall energy consumption of the UAV and the ground network.

## B. Related Works

1) *UAV Trajectory Planning*: Energy-efficient trajectory planning for UAVs has recently attracted significant research interest, and multiple solutions have been proposed for UAV-enabled wireless networks. In general, existing solutions for energy-efficient UAV trajectory planning can be loosely classified into two categories: non-machine learning-based methods and machine learning-based methods. In the first category, researchers mostly use mathematical programming or heuristic algorithms to solve the trajectory optimization problem. However, the computation time of mathematical programming algorithms may increase exponentially as the problem size increases, e.g., [11], [12]. Although some heuristic algorithms are applied to design the energy-efficient path in the UAV-enabled wireless networks, such as ant colony optimization [13] and cuckoo search [14], they usually cannot fully adapt to the increasing complexity of scalable wireless networks.

Regarding the machine learning-based category, deep reinforcement learning (DRL) and reinforcement learning (RL) are the most common techniques in solving the UAV's trajectory planning. In [15], the authors propose a DRL-based method which is composed of two deep neural networks (DNNs) and deep deterministic policy gradient (DDPG) to maximize the energy efficiency for a group of UAVs by jointly considering communications coverage, energy consumption, and connectivity. In order to minimize the UAV's transmission and hovering energy, the authors in [16] formulate the energy-efficient optimization problem as a Markov decision process. Then, they use two DNNs and the actor-critic-based RL algorithm to develop an online DRL algorithm that shows a good performance in terms of energy savings. In [17], the authors jointly optimize the UAV's 3D trajectory and the frequency band allocation of ground users by considering the UAV's energy consumption and the fairness of the ground users. A DDPG-based DRL algorithm is developed to generate the energy-efficient trajectory with fair communication service to ground users. In [18], with the aim of designing an energy-efficient UAV's route for long-distance sensing tasks, the authors propose a DRL-based framework where convolutional neural networks (CNNs) are used for extracting features and the deep Q-network (DQN) is utilized to make decisions. Towards realizing green UAV-enabled Internet of Things (IoT), the authors in [19] formulate the UAV's path planning problem as a dynamic decision optimization problem, which is solved by dueling DQN. The aforementioned machine learning-based methods show strong ability to handle complex wireless environments and effectively learn the UAV's trajectory policy from experiences; however, they are implemented by using some common neural network models, such as DNNs, DQN, and CNNs. Instead of common neural network models, in this paper we exploit the appealing concept of sequence-to-sequence learning that originally emerged in *Natural Language Processing* to design the DRL algorithm to solve the UAV's path planning problem in clustered WSNs.

2) *Sequence-to-Sequence Learning*: Sequence-to-sequence learning has shown great success in machine translation where sentences are mapped to correct translations [20]. Over the last several years, many neural networks based on sequence-to-sequence models have been proposed in different applications. For example, pointer network is one of the extensively studied models because of its excellent ability in solving sequence decision problems. In [21], pointer network is trained in a supervised fashion to solve the traveling salesman problem (TSP). The work of [22] uses a RL-based unsupervised method to train the pointer network and obtains better results when compared to the supervised learning in [21]. In [23], the authors propose a structural graph embedded pointer network to develop online vehicular routes in intelligent transportation systems. The authors in [24] propose a modified pointer network to solve the keyword recommendation problem in sponsored search advertising system. In [25], a simplified pointer network is introduced to solve Vehicle Routing Problem (VRP) in dynamic traffic environments. Different from the above-discussed research works, we extend the state-of-the-art pointer network-based DRL to solve the UAV's trajectory planning problem. Our contributions are elaborated in the next subsection.

### C. Contributions

We aim to minimize the overall UAV-WSN's energy consumption by designing an efficient UAV's trajectory in a clustered WSN. Since the UAV's visiting order to the ground CHs can be seen as a sequence decision problem, we propose a sequence-to-sequence learning-based DRL strategy with pointer network to deal with the challenging trajectory planning problem. The pointer network can capture the relation between a problem instance and its solution by using a sequence-to-sequence neural network. It has been demonstrated to be an effective method to solve some NP-hard problems, such as TSP [22] and VRP [25]. Hence, it is expected that the pointer network-based DRL algorithm is also promising for solving the problem of the UAV trajectory planning for the UAV-WSN system. The main contributions of this paper are summarized as follows:

- 1) We consider a UAV-enabled energy-efficient data collection framework for clustered WSNs. We formulate an optimization problem to minimize the energy consumption of the entire UAV-WSN system by jointly designing the UAV's trajectory and selecting CHs in pre-determined clusters of the ground WSN.
- 2) We show that the UAV's trajectory planning problem in the clustered WSN can be seen as a sequence of decisions. Hence, a sequence-to-sequence pointer network-A\* (Ptr-A\*) model is proposed to solve the formulated problem. Particularly, the pointer network is utilized to model the visiting order of all ground clusters, and A\* algorithm [26] is used to efficiently select CHs from clusters' ground nodes. The UAV's start point and all clusters, as the input, are fed into the Ptr-A\* model, and its output is a set of CHs and the visiting order to these CHs, i.e., the UAV's trajectory.

- 3) We use a self-driven learning mechanism that only needs the reward calculation to train the parameters of Ptr-A\* network on problem instances with small-scale clusters for faster training.
- 4) Our proposed DRL method has an excellent generalization capability with respect to the number of clusters used for training. In other words, given a new problem instance with any number of clusters, the trained model can automatically generate a trajectory for the UAV to visit clusters, without retraining the new model.
- 5) We perform extensive simulations to demonstrate that the proposed DRL method outperforms other baseline techniques when considering both the computation times and energy consumption results.

The rest of this paper is organized as follows. Section II presents the system model and the problem formulation. Section III describes the proposed DRL algorithm. Section IV provides simulation results. Finally, Section V concludes the paper.

## II. SYSTEM MODEL AND PROBLEM FORMULATION

As mentioned earlier, we assume that devices on the ground have been clustered according to some specific criterion, e.g., based on their geographical locations; hence, clustering techniques for the ground WSN are not discussed in this work. In particular, we consider  $K$  clusters of sensor nodes  $\{\mathbf{G}_1, \dots, \mathbf{G}_K\}$  located in the sensing (service) area for data collection. Each cluster contains  $N$  nodes, one of which is the CH, represented by  $b_k \in \mathbf{G}_k$ , that will be selected by our proposed algorithm. We assume that only one rotary-wing UAV is dispatched to visit CHs to collect data from the ground network. The UAV takes off from the start position  $b_0$  and then back to  $b_0$  after finishing the data collection task. The trajectory of the UAV should contain the start/end hovering position  $c_0$  corresponding to  $b_0$ , and the  $K$  target hovering positions  $\{c_1, \dots, c_K\}$ , which are vertically above the ground CHs. Hence, the trajectory planning problem of the UAV can be seen as a permutation of  $(K + 1)$  hovering positions. It is obvious that the locations of CHs determine the flight trajectory, and hence, the energy consumption of UAV and the ground nodes. We consider a three-dimensional (3D) Cartesian coordinates system to define the positions of ground nodes and the UAV. The position of the CH of the  $k$ -th cluster is  $b_k = (x_k, y_k, 0)$ . Correspondingly, the coordinate of the UAV's hovering position  $c_k$  can be represented by  $(x_k, y_k, H)$ , where  $H$  is the fixed flight height of the UAV. Similarly, the coordinate of the  $n$ -th member node of the  $k$ -th cluster  $b_k^{(n)}$  is denoted as  $(x_k^{(n)}, y_k^{(n)}, 0)$ ,  $n = 1, \dots, N - 1$ ,  $k = 1, \dots, K$ , and  $b_k^{(n)} \neq b_k$ .

### A. Channel Model

There is a number of channel models that have been developed for UAV communications, e.g., [27], [28]. In this work, we consider a simple air-to-ground channel model that is described as follows. For ground-to-air communication, there is a certain probability that each CH  $b_k$  has a LoS view towards the UAV



when it hovers at the hovering position  $c_k$ . This probability typically depends on the environment and elevation angle, and is given by [8]

$$P_{\text{LoS}} = \frac{1}{1 + \eta \exp(-\beta[\tau - \eta])}, \quad (1)$$

where  $\eta$  and  $\beta$  are constants determined by environment, and  $\tau = \frac{180}{\pi} \times \sin^{-1}(\frac{H}{d_k})$ , where  $d_k$  is the distance between  $b_k$  and  $c_k$ . Since it is assumed that each hovering position is directly above the corresponding CH, one has  $d_k = H$ . Obviously, the non-line-of-sight (NLoS) probability is given by  $P_{\text{NLoS}} = 1 - P_{\text{LoS}}$ . The average path loss between each CH and the UAV can be expressed as [8]

$$\bar{P}_{\text{loss}} = P_{\text{LoS}}(K_0 + \mu_{\text{LoS}}) + P_{\text{NLoS}}(K_0 + \mu_{\text{NLoS}}) \quad (2)$$

where  $\mu_{\text{LoS}}$  and  $\mu_{\text{NLoS}}$  are the mean values of the excessive path losses in LoS and NLoS links, respectively,  $K_0 = 10\alpha \log_{10}(\frac{4\pi f_c H}{c})$ ,  $\alpha$  is the path loss exponent,  $c$  is the speed of light, and  $f_c$  is the carrier frequency. Thus, the average data rate from each CH to the UAV can be computed as [8]

$$r_{\text{data}} = B_{\text{width}} \log_2 \left( 1 + \frac{P_{\text{CH}}}{\bar{P}_{\text{loss}} N_0} \right) \quad (3)$$

where  $B_{\text{width}}$  is the available bandwidth,  $N_0$  is the noise power spectral density, and  $P_{\text{CH}}$  is the transmit power of each CH.

### B. UAV's Energy and Trajectory Model

We assume that the UAV supports a flying-hovering mode without considering acceleration-deceleration patterns. After the UAV flies to the hovering position  $c_k$  with a fixed speed  $v_{\text{UAV}}$ , it hovers there and transfers a beacon frame to wake up the corresponding CH  $b_k$  from sleep mode to active model. Then,  $b_k$  starts to collect data from its member nodes by time-division multiple access (TDMA) and forwards the collected data to the UAV. At each hovering position, the energy consumption of the UAV includes two parts: communication-related energy and hover-related energy. The hovering power is given by [8], [29]

$$P_{\text{hover}} = \sqrt{\frac{(m_{\text{tot}}g)^3}{2\pi r_p^2 n_p \rho}} \quad (4)$$

where  $g$  is the earth gravity,  $\rho$  is the air density,  $n_p$  is the number of propellers,  $r_p$  is the propeller radius, and finally  $m_{\text{tot}}$  is the mass of the UAV. Thus, the energy consumed by the UAV at each hovering position  $c_k$  is given by

$$\begin{aligned} E_{c_k} &= T_k(P_{\text{hover}} + P_{\text{com}}) \\ &= \frac{D_k}{r_{\text{data}}}(P_{\text{hover}} + P_{\text{com}}) \end{aligned} \quad (5)$$

where  $T_k$  is the total hovering time of the UAV at  $c_k$ ,  $D_k$  is the amount of data that needs to be transferred from CH  $b_k$  to the UAV, and  $P_{\text{com}}$  is the communication power of the UAV. In order to simplify the analysis, we assume that the hovering time is equal to the data transmission time from  $b_k$  to the UAV.

The horizontal movement power is assumed as a linear function of the UAV's flight speed  $v_{\text{UAV}}$ , which is expressed

as [8], [29]

$$P_{\text{move}} = \frac{P_{\text{max}} - P_{\text{idle}}}{v_{\text{max}}} v_{\text{UAV}} + P_{\text{idle}} \quad (6)$$

where  $v_{\text{max}}$  is the maximum speed of the UAV,  $P_{\text{max}}$  and  $P_{\text{idle}}$  are the hardware power levels when the UAV is moving at full speed and when the UAV is in idle state, respectively. Because the UAV needs to start from the start hovering location  $c_0$ , goes through all target hovering positions  $c_1, \dots, c_K$ , and then back to  $c_0$ , the total energy consumption of the UAV in flight is given by [8], [29]

$$E_{\text{flight}} = T_{\text{flight}}(P_{\text{hover}} + P_{\text{move}}) \quad (7)$$

where  $T_{\text{flight}}$  is the total flight time, which can be expressed as

$$T_{\text{flight}} = \frac{1}{v_{\text{UAV}}} \sum_{i=0}^K \sum_{\substack{j=0 \\ j \neq i}}^K d_{c_i, c_j} L_{c_i, c_j}, \quad \forall c_i, c_j \in \mathcal{C} \quad (8)$$

where  $\mathcal{C} = \{c_0, c_1, \dots, c_K\}$ ,  $c_k$  is determined by  $b_k$ ,  $b_k \in \mathcal{G}_k$ , and  $L_{c_i, c_j}$  specifies whether the UAV travels from stop position  $c_i$  to  $c_j$ , which is defined as

$$L_{c_i, c_j} = \begin{cases} 1, & \text{if the path goes from } c_i \text{ to } c_j \\ 0, & \text{otherwise.} \end{cases} \quad (9)$$

The quantity  $d_{c_i, c_j}$  is the Euclidean distance between  $c_i$  and  $c_j$ , which is given by

$$d_{c_i, c_j} = \|c_i - c_j\| = \|b_i - b_j\|. \quad (10)$$

In order to meet the requirements of the UAV's trajectory, we need to consider the following constraints:

$$\sum_{\substack{i=0 \\ i \neq j}}^K L_{c_i, c_j} = 1, \quad \forall c_i, c_j \in \mathcal{C} \quad (11)$$

$$\sum_{\substack{j=0 \\ j \neq i}}^K L_{c_i, c_j} = 1, \quad \forall c_i, c_j \in \mathcal{C} \quad (12)$$

$$\sum_{c_i \in \mathbf{F}} \sum_{c_j \in \mathbf{F}} L_{c_i, c_j} \leq |\mathbf{F}| - 1, \quad \forall \mathbf{F} \subset \mathcal{C}; |\mathbf{F}| \geq 2. \quad (13)$$

The constraints (11) and (12) guarantee that there is only one UAV path entering and leaving a given node, which means that the UAV should visit each point in  $\mathcal{C}$  exactly once. Constraint (13) is the sub-trajectories elimination constraint and enforces that no partial loop exists where  $\mathbf{F}$  is the subset of  $\mathcal{C}$  [32], which means there is only one single trajectory covering all CHs.

According to the above analysis, the total energy consumption of the UAV is composed of the flying-related and the hovering-related energy consumption, which can be written as

$$E_{\text{UAV}} = E_{\text{flight}} + \sum_{k=1}^K E_{c_k}. \quad (14)$$

### C. Ground Network and Energy Model

We assume that all nodes have the same computation and transmission capabilities. In other words, all nodes are capable of acting as a CH. Nodes are static after being deployed. All member nodes transmit their sensing information to CHs periodically and CHs forward the collected data to the UAV. We also assume that the transmission energy of each node is sufficient to send messages to its CH. In addition, the UAV can simultaneously connect to at most one CH. Hence, there is no interference among neighboring CHs.

The energy consumption in the ground network includes two components. The first component is the communication energy consumption between CHs and their member nodes. The first-order radio model [30] is used to calculate the energy consumption of the ground network. The transmission energy is consumed by the transmitter's circuitry and power amplifier. If the distance between a member node and its CH is less than a given threshold, the power amplifier uses the free space model; otherwise, the multi-path model is used [31]. The energy consumed to transmit an  $l$ -bit message from a member node to its CH  $b_k$  is given by [30]

$$E_n^{b_k} = lE_{\text{elec}} + l(\chi \varepsilon_{\text{fs}} d_{n,b_k}^2 + (1 - \chi) \varepsilon_{\text{mp}} d_{n,b_k}^4) \quad (15)$$

where

$$\chi = \begin{cases} 1, & d_{n,b_k} \leq d_0 \\ 0, & d_{n,b_k} > d_0 \end{cases} \quad (16)$$

and

$$d_0 = \sqrt{\frac{\varepsilon_{\text{fs}}}{\varepsilon_{\text{mp}}}}. \quad (17)$$

In (15),  $E_{\text{elec}}$  is the dissipated energy per bit in the circuitry,  $d_{n,b_k}$  is the distance between the CH  $b_k$  and one of its member nodes  $n$ ,  $n = 1, \dots, N-1$ ,  $d_0$  is the distance threshold,  $\varepsilon_{\text{fs}}$  and  $\varepsilon_{\text{mp}}$  represent the radio amplifier's energy parameter of the free space and multi-path fading models, respectively. Moreover, the energy consumed to receive an  $l$ -bit message from member node  $n$  by  $b_k$  is given by [30]

$$E_{b_k}^{(n)} = lE_{\text{elec}}. \quad (18)$$

In addition, the second component of the energy consumption of the ground network is the energy consumed by each CH  $b_k$  to complete its data transmission to the UAV. This component can be written as

$$E_{b_k} = P_{\text{CH}} T_k = P_{\text{CH}} \frac{(N-1)l}{r_{\text{data}}} \quad (19)$$

where  $(N-1)l$  is the amount of data transferred by  $b_k$  to the UAV. Hence, the total energy consumption of all nodes in the ground network in a complete data collection task, where member nodes transmit the sensing data to their CHs and CHs forward data to the UAV, is

$$E_{\text{ground}} = \sum_{k=1}^K \sum_{n=1}^{N-1} (E_n^{b_k} + E_{b_k}^{(n)}) + \sum_{k=1}^K E_{b_k}. \quad (20)$$

### D. Problem Formulation for UAV's Trajectory

Based on (14) and (20), the total weighted energy consumption in the UAV-WSN system can be formulated as

$$E = \omega \left( \sum_{k=1}^K \sum_{n=1}^{N-1} (E_n^{b_k} + E_{b_k}^{(n)}) + \sum_{k=1}^K E_{b_k} \right) + (1 - \omega) \left( E_{\text{flight}} + \sum_{k=1}^K E_{c_k} \right), \quad 0 \leq \omega \leq 1 \quad (21)$$

where the first term corresponds to the total energy consumption of the ground network, while the second term is the energy consumption of the UAV, and  $\omega$  is the weighting coefficient that can be adjusted to achieve the trade-off between the two terms. With the aim of minimizing the total energy consumption of the ground network and the UAV, we jointly find a set of CHs from the ground cluster-based WSN and design the UAV's visiting order to these CHs. The optimization problem of interest is formulated as

$$\begin{aligned} & \min_{\substack{\{b_0, b_1, \dots, b_K, \dots, b_K\} \\ b_k \in \mathcal{G}_k}} E \\ & \text{s.t. (9), (11) - (13).} \end{aligned} \quad (22)$$

Obviously, the problem at hand is a constrained combinatorial optimization problem, which is NP-hard. Some promising approaches have been put forward to solve such combinatorial optimization problems, and their advantages and disadvantages are discussed below.

- 1) *Exact methods*: Exact methods often search for the optimal solution of the problem through systematic enumeration, integer programming, and constraint programming, etc. [33]. At least in theory, they can provide the optimal solution for the optimization problem. However, such algorithms cannot be applied to combinatorial optimization problems with large data scale because their computation complexity becomes prohibitive.
- 2) *Heuristics*: Heuristics are higher-level problem-independent algorithmic frameworks that provide a set of guidelines to develop optimization algorithms [34]. However, they generally cannot guarantee to find globally optimal solutions.
- 3) *RL*: Q-learning, one of the RL techniques, is demonstrated to be promising in solving NP-hard problems [35]. Specifically, it can deal with the path decision problem when provided with sufficient state space variables. However, if the number of ground nodes is high, Q-learning will need more storage space for action and state space variables [18].

As discussed before, DRL has recently shown to have important advantages in solving combinatorial optimization problems. A typical neural combinatorial framework is proposed in [22] that uses RL to optimize a policy modeled by the pointer network. In [25], the authors view a combinatorial optimization

problem as a sequence of decisions, and they use a sequence-to-sequence neural network model and the RL approach to obtain a near-optimal solution for the optimization problem. Inspired by these promising developments, we extend the application of sequence-to-sequence model to solve the UAV's trajectory planning problem described earlier.

### III. DEEP REINFORCEMENT LEARNING FOR UAV TRAJECTORY PLANNING

Because all clusters  $\{G_1, \dots, G_K\}$  must be visited by the UAV sequentially, we convert the visiting decisions problem into a sequence-to-sequence prediction problem. The problem can be simply formalized as follows. Given start position and all clusters, denoted by  $\mathcal{G} = \{b_0, G_1, \dots, G_K\}$ , we want to output a permutation of the items in  $\mathcal{G}$  that maximizes some measure of interest. The output sequence is denoted as  $\mathcal{T} = \{\pi_0, \pi_1, \dots, \pi_K\}$ , where each  $\pi_t$  is the index of any element in  $\mathcal{G}$  being placed at the  $t$ -th position of  $\mathcal{T}$ . In fact,  $\mathcal{T}$  is the UAV's visiting order to clusters in our problem. Thus, for a given input sequence  $\mathcal{G}$ , the probability of the output sequence  $\mathcal{T}$  can be factorized by a product of conditional probabilities according to the chain rule

$$P_\theta(\mathcal{T}|\mathcal{G}) = \prod_{t=0}^K P(\pi_t|\pi_0, \dots, \pi_{t-1}, \mathcal{G}) \quad (23)$$

where  $t$  is the time step,  $P_\theta(\mathcal{T}|\mathcal{G})$  parameterized by  $\theta$  is a *stochastic policy* for deciding the visiting order. The conditional probability  $P(\pi_t|\cdot)$  models the probability of any cluster being visited at the  $t$ -th time step according to the given  $\mathcal{G}$  and clusters already visited at previous time steps [36]. A trained  $\theta$  can assign high probabilities to good results and low probabilities to bad results. The reinforcement learning can be applied to train the optimal model policy  $\theta^*$  for producing the optimal visiting order  $\mathcal{T}^*$  with the highest probability.

#### A. Pointer Network-A\* Architecture for UAV's Trajectory Planning

With the rapid development of neural network techniques, the neural network-based frameworks have been applied to sequence-to-sequence learning [20]. The general sequence-to-sequence neural network [25] encodes the input sequence into a vector that includes information of the input by a recurrent neural network (RNN), called encoder, and decodes the vector to the target sequence by another RNN, called decoder [21]. In this work, we employ the pointer network to model the conditional probability  $P(\pi_t|\cdot)$ , which has been proved to be effective to solve the combinatorial optimization problems. The architecture of the pointer network is similar to sequence-to-sequence network model, but it uses attention mechanism as a pointer to choose items of its input sequence as the output.

The proposed Ptr-A\* model in this work is elaborated as follows.

1) *Encoder*: The traditional RNN shows poor performance in dealing with the problem of long-term dependencies, which makes it difficult to be trained in practice [38]. Hence, we use

Long Short-Term Memory (LSTM) cells which are capable of learning long-term dependencies to construct a RNN as the encoder. Each item (a start position or a cluster) in  $\mathcal{G}$  is converted into a high  $D$ -dimensional vector space, which enables the policy to extract useful features much more efficiently in the transformed space [37]. Then, the embedding vectors are fed into LSTM cells. At each encoding step, the LSTM cell reads one embedded item and outputs a latent memory state. Finally, the input sequence  $\mathcal{G}$  is transformed into a sequence of latent memory states  $\mathcal{E} = \{e_0, \dots, e_K\}$ , each  $e_k \in \mathbb{R}^D$ . In fact, the motive of the encoder network is to acquire the representation for each element in  $\mathcal{G}$ .

2) *Decoder*: We also adopt LSTM cells to construct the RNN of the decoder network. The output  $\{e_0, \dots, e_K\}$  of the encoder are given to the decoder network. At each decoding step  $t$ , the LSTM cell outputs the hidden state  $h_t \in \mathbb{R}^D$  that includes the knowledge of previous steps. And then, the decoder employs the attention mechanism to output the visiting decision  $\pi_t$  based on  $h_t$  and  $\{e_0, \dots, e_t\}$ . Attention mechanism can help the model to give different weights to different elements of the input and extract more critical information [39]. It tells us the relationship between each element in the input at current step  $t$  and the output  $\pi_{t-1}$  of the last decoding step. The most relevant element with the maximum conditional probability is selected as the access element at decoding step  $t$ . Thus, the calculation is given by

$$u_j^t = \begin{cases} \varphi \tanh(W_1 e_j + W_2 h_t), & \text{if } j \notin \{\pi_0, \dots, \pi_{t-1}\} \\ -\infty, & \text{otherwise} \end{cases} \quad (24)$$

where  $W_1, W_2 \in \mathbb{R}^{D \times D}$  are attention matrices,  $\varphi \in \mathbb{R}^{1 \times D}$  is the attention vector.  $W_1, W_2$ , and  $\varphi$  are denoted collectively by  $\theta$ , which is the learnable parameter in our pointer network. In essence,  $u_j^t$  is the score associated with item  $j$  ( $e_j$ ) in position  $t$ .

The conditional probability is calculated by a *softmax* function over the remaining items (not visited in the previous steps), as follows:

$$P(\pi_t = j|\pi_0, \dots, \pi_{t-1}, \mathcal{G}) = \text{softmax}(u_j^t) = \frac{\exp(u_j^t)}{\sum_{m \notin \{\pi_0, \dots, \pi_{t-1}\}} \exp(u_m^t)}, j \in m. \quad (25)$$

The probability  $P(\pi_t = j|\cdot)$  represents the degree to which the model points to item  $j$  at the decoding step  $t$  [36]. As shown in the example in Fig. 2, the start position  $b_0$  and clusters  $G_1, G_2, G_3$  are inputted into the encoder network.  $v_{go}$  is the start tag of the decoder, which is a learned vector. At each decoding step, the item of the input sequence with the highest probability is pointed by a thicker black arrow. The output of the 0-th decoding step points to  $b_0$  ( $\pi_0$ ), which will be visited at this step and given as the input of the next decoding step. Finally, we will obtain a visiting order sequence  $\mathcal{T} = \{\pi_0, \pi_1, \pi_2, \pi_3\}$  corresponding to the input sequence.

To help the reader familiarize with the attention mechanism, a numerical example is provided in Appendix A.

3) *A\* Search*: Once the output sequence of the decoder network is obtained, we can build a search graph for all clusters according to this sequence, where each layer is composed of nodes of one cluster. This is illustrated with an example in Fig. 2.

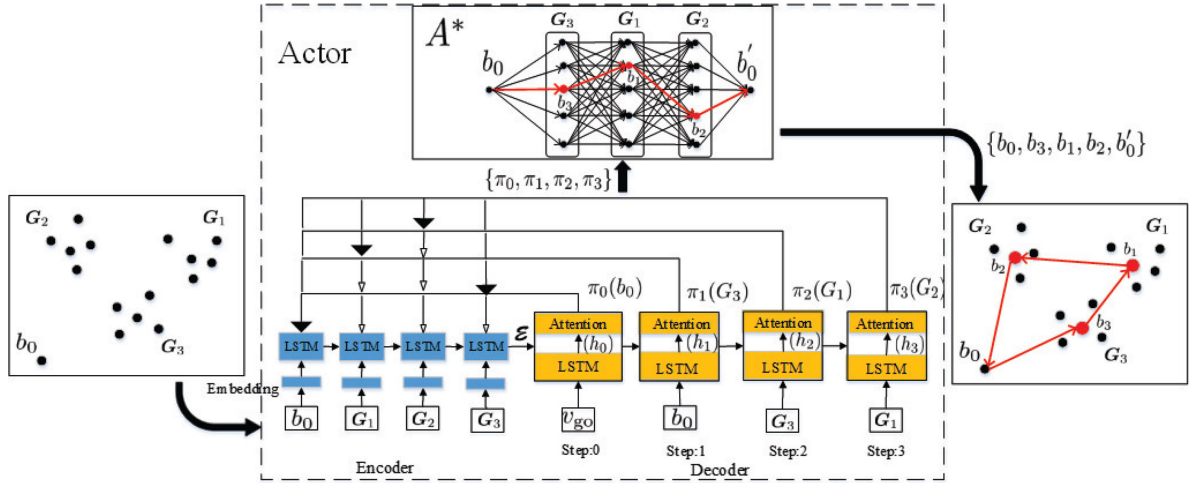


Fig. 2. Example of Ptr-A\* architecture for a 3-clusters network.

It is worth mentioning that the first layer is the start position  $b_0$  and the last layer is the end position  $b'_0$  which is the copy of  $b_0$ . Thus, the created graph has a total of  $(K+2)$  layers. We use the A\* search algorithm, one of best path-finding algorithms, to find the CH from each cluster to build a path having the smallest cost (total weighted energy consumption of the UAV-WSN system) from the start position to the end position. In each iteration, the A\* algorithm needs to calculate the cost of the traversed path and the estimated cost required to extend the path to the end to determine which of its partial paths to expand into one or more longer paths [26]. Any node  $m$  is chosen to be visited by the following function

$$f(m) = g(m) + h(m) \quad (26)$$

where  $g(m)$  represents the exact energy consumption of the UAV-WSN system when the UAV moves from the start node to a candidate node  $m$ , following the path generated to get there,  $h(m)$  is the estimated energy consumption of the UAV to travel from the candidate node  $m$  to the end. Then, the node with the lowest  $f(m)$  value is selected from candidate nodes as the next node to be traversed.

The main implementation of the A\* algorithm is to maintain two lists. The OPEN list contains those nodes that are candidates for checking. The CLOSED list contains those nodes that have been checked. The neighbor nodes of any node located in any layer are defined as all nodes in its previous and next layers. Also, each node keeps a pointer to its parent node so that we can determine how it was found, which is implemented by a map COME\_FROM. The pseudocode of using the A\* algorithm to find the path from the start position to the end position is described in Algorithm 1. As shown in the example in Fig. 2, the output of the A\* algorithm is the trajectory from  $b_0$  to  $b'_0$ , which can ensure the minimum energy consumption  $E$ . In addition, the CH of each cluster is found on this trajectory, given by  $\{b_0, b_3, b_1, b_2, b'_0\}$ . Finally, the Ptr-A\* model outputs the trajectory and the minimum energy consumption  $E$ .

### B. Parameters Optimization With Reinforcement Learning

In order to find a good trajectory for the UAV, we need to obtain the optimal model parameter  $\theta^*$  that can be trained from samples. If we adopt a supervised learning to train the model parameter, high-quality labeled data is needed because it decides the performance of the model. However, it is expensive to get the high-quality labeled data in practice for the proposed UAV's trajectory problem. Instead, we choose the well-known model-free policy-based RL, known as the actor-critic algorithm [40], to train the model because it is shown to be an appropriate paradigm for training neural networks for combinatorial optimization [22]. The UAV works as the agent to make a sequential action set in a given state of the environment. In the following, we describe the state, action, reward, and training of the proposed DRL algorithm.

1) *State*: The state includes coordinates for all clusters, the UAV's location, and the energy consumption of UAV-WSN at current step  $t$ .

2) *Action*: The action represents the choice of the next cluster to be selected at current step  $t$  and the CH in this cluster. Thus, we define the output of the right-hand side of (23) and the CH selection by A\* as the action at each step.

3) *Reward*: We design the reward as the negative of the total energy consumption in (21). This means that the DRL is set to get the maximal reward (minimal energy consumption).

4) *Training*: The actor-critic method includes the actor network and the critic network. The actor network is the proposed Ptr-A\* in this work. The critic network is used to provide an approximated baseline of the reward for any problem instance to reduce the variance of gradients during the training phase and increases the speed of learning [40]. Our critic network, parameterized by  $\psi$ , has the same architecture as that of the encoder of the Ptr-A\*. Then, its hidden states are decoded into a baseline prediction by two fully-connected ReLU layers [22].

Our training objective is the expected energy consumption, which is defined as

$$J(\theta|\mathcal{G}) = \mathbb{E}_{\mathcal{T} \sim p_{\theta}(\cdot|\mathcal{G})}[E]. \quad (27)$$



**Algorithm 1:** A\* Search Algorithm for the Trajectory Planning.

---

**Input:**  $\mathcal{T}$   
**Output:** Trajectory, minimum energy consumption  $E$

- 1: Build a search graph by  $\mathcal{T}$
- 2: Initialize OPEN, CLOSED, and COME\_FROM
- 3:  $f(b_0) = 0$ , OPEN.add( $b_0$ )
- 4: **while** OPEN is not empty **do**
- 5:   Find the node  $q$  with the lowest  $f(q)$  from OPEN
- 6:   **if**  $q = b'_0$  **then**
- 7:     Construct path from  $b_0$  to  $b'_0$  by COME\_FROM
- 8:     **return** Trajectory,  $E$
- 9:   **end if**
- 10: OPEN.remove( $q$ )
- 11: CLOSED.add( $q$ )
- 12: Obtain neighbor nodes of  $q$
- 13: **for** each neighbor node  $m$  of  $q$  **do**
- 14:    $cost = g(q) + \text{the energy consumption of UAV-WSN from } q \text{ to } m$
- 15:   **if**  $m$  in OPEN and  $cost < g(m)$  **then**
- 16:     OPEN.remove( $m$ )
- 17:   **end if**
- 18:   **if**  $m$  in CLOSED and  $cost < g(m)$  **then**
- 19:     CLOSED.remove( $m$ )
- 20:   **end if**
- 21:   **if**  $m$  not in OPEN and CLOSED **then**
- 22:      $g(m) = cost$
- 23:      $f(m) = g(m) + h(m)$
- 24:     OPEN.add( $m$ )
- 25:     COME\_FROM[ $m$ ] =  $q$  //set  $m$ 's parent
- 26:   **end if**
- 27: **end for**
- 28: **end while**

---

We use policy gradient method and stochastic gradient descent to optimize  $\theta$ . The gradient of (27) is formulated by REINFORCE [41] algorithm

$$\nabla_{\theta} J(\theta|\mathcal{G}) = \mathbb{E}_{\mathcal{T} \sim p_{\theta}(\cdot|\mathcal{G})} [(E - V_{\psi}(\mathcal{G})) \nabla_{\theta} \log p_{\theta}(\mathcal{T}|\mathcal{G})] \quad (28)$$

where  $V_{\psi}(\mathcal{G})$  is a baseline function for reducing the variance of the gradients, which is implemented by the critic network. Assume we have  $B$  *i.i.d* train samples, the gradient in (28) can be approximated with Monte Carlo sampling as follows

$$\nabla_{\theta} J(\theta) \approx \frac{1}{B} \sum_{i=1}^B (E_i - V_{\psi}(\mathcal{G}_i)) \nabla_{\theta} \log p_{\theta}(\mathcal{T}_i|\mathcal{G}_i). \quad (29)$$

We train the parameters of the critic with stochastic gradient descent on a mean squared error objective  $L(\psi)$  between its predictions  $V_{\psi}(\mathcal{G}_i)$  and the actual energy consumption.  $L(\psi)$  is formulated as

$$L(\psi) = \frac{1}{B} \sum_{i=1}^B (V_{\psi}(\mathcal{G}_i) - E_i)^2. \quad (30)$$

**Algorithm 2:** Training Ptr-A\* by Actor-Critic Algorithm.

---

**Input:** Training samples set  $D = \{\mathcal{G}_1, \mathcal{G}_2, \dots\}$ , batch size  $B$ , training steps  $S$

- 1: Initialize actor network  $\theta$  and critic network  $\psi$  with random weights
- 2: **for**  $s = 1$  to  $S$  **do**
- 3:   Sample  $\mathcal{G}_i$  from  $D$ ,  $\forall i \in \{1, \dots, B\}$
- 4:   Calculate  $E_i$  and  $\mathcal{T}_i$  with Ptr-A\* network,  $\forall i \in \{1, \dots, B\}$
- 5:   Calculate  $V_{\psi}(\mathcal{G}_i)$  with Critic network,  $\forall i \in \{1, \dots, B\}$
- 6:    $d\theta \leftarrow \frac{1}{B} \sum_{i=1}^B (E_i - V_{\psi}(\mathcal{G}_i)) \nabla_{\theta} \log p_{\theta}(\mathcal{T}_i|\mathcal{G}_i)$
- 7:    $L(\psi) \leftarrow \frac{1}{B} \sum_{i=1}^B (V_{\psi}(\mathcal{G}_i) - E_i)^2$
- 8:    $\theta \leftarrow \text{Adam}(\theta, d\theta)$
- 9:    $\psi \leftarrow \text{Adam}(\psi, \nabla_{\psi} L(\psi))$
- 10: **end for**
- 11: **return**  $\theta^* = \theta$

---

The training procedure is presented in Algorithm 2. Notice that Adam algorithm is used to update the parameters of the actor network and the critic network iteratively. Adam algorithm designs independent adaptive learning rates for different parameters via calculating the first and second moment estimates of the gradient instead of using a single learning rate to update all parameters by the traditional random gradient descent [42]. Given the initial learning rate, the learning rates in different steps adaptively change according to the learning results. Because of the generalization property of RNNs [43], our proposed models, including the Ptr-A\* and critic, have a very good generalization ability. In the training phase, we can use small problem instances to train the model, and then the trained model can be utilized to solve large problem instances.

## IV. NUMERICAL RESULTS

In this section, we first introduce detailed environment settings, and then describe the decoding search strategies at inference. Furthermore, we compare the performance of the proposed DRL algorithm with several baseline algorithms.

## A. Environmental Settings and Model Training

We consider the ground network size of 2 km  $\times$  2 km, and the start position of the UAV is located at (0 m, 0 m). Simulation parameters are listed in Table I. We use mini-batches of size 512 and LSTM cells with 128 hidden units in the encoder and the decoder. We implement the proposed model by using Pytorch 1.4 and Python 3.7 on a VM instance of Google Cloud Platform with 1 NVIDIA TESLA P100 GPU. The parameters of both the actor and critic networks are initialized by the Xavier initialization method and trained by the Adam optimizer with an initial learning rate of 0.0001 and decayed every 5000 steps by 0.96.

It is assumed that nodes in a given cluster  $\mathcal{G}_k$  are distributed according to the Gaussian distribution. Each cluster's nodes are sampled from a *torch.normal*( $\nu$ , std) where  $\nu$  is the mean and std



TABLE I  
SIMULATION PARAMETERS

Parameter	Description	Value
$\varepsilon_{fs}$	Amplifier's energy parameter of the free space fading	10 pJ/bit/m <sup>2</sup> [30]
$\varepsilon_{mp}$	Amplifier's energy parameter of the multi-path fading	0.0013 pJ/bit/m <sup>2</sup> [30]
$E_{elec}$	Energy consumption per bit in the circuitry	50 nJ/bit [30]
$P_{CH}$	Transmit power of each CH	21 dBm/Hz [8]
$N$	Number of nodes per cluster	20
$B_{width}$	Bandwidth	1 MHz
$N_0$	Noise power	-174 dBm/Hz [8]
$f_c$	Carrier frequency	2 GHz [8]
$\alpha$	Path loss exponent	3 [8]
$H$	UAV's flight height	50 m
$\mu_{LoS}, \mu_{NLoS}$	Mean values of the excessive path loss	1 dB, 20 dB [44]
$\beta, \eta$	Environmental parameters	0.03, 10 [8]
$v_{UAV} = v_{max}$	UAV's flight speed	15 m/s [8]
$m_{tot}$	UAV's mass	500 g [45]
$r_p$	Radius of UAV's propellers	20 cm [45]
$n_p$	Number of propellers	4 [45]
$P_{max}$	UAV's hardware power level at full speed	5 W [45]
$P_{idle}$	UAV's hardware power level when it hovers	0 W [45]
$P_{com}$	UAV's communication power	0.0126 W [8]

is the constant standard deviation. Each Gaussian distribution's  $\nu$  is randomly sampled from a *torch.rand()* function to determine the position of each cluster in the two-dimensional space. We train the model using instances of 20 clusters and 40 clusters, respectively. The 20-clusters model is trained for 100 000 steps, and the 40-clusters model is trained for 200 000 steps. We give a simple example on how to obtain the train data. At each training step of the 20-clusters model, we sample 20 means from *torch.rand()* for 20 Gaussian distributions, respectively. Then, we use these 20 *i.i.d.* distributions to generate a set (problem instance) of 20 clusters where the number of nodes per cluster is 20. The test data sets are also generated in the same way, only the number of clusters is different.

### B. Decoding Search Strategies At Inference

Given a new problem instance  $\mathcal{G}$  at inference, the decoder network of our trained Ptr-A\* architecture can easily output an access sequence for all clusters. The decoding process of the decoder at inference shows how solvers search over a large set of feasible access sequences. In this work, we consider the following three decoding search strategies.

1) *Greedy Search*: Greedy search strategy always select the cluster with the largest probability at each decoding step during inference, which is labeled as DRL-greedy in the simulation results.

2) *Sampling Search*: This strategy samples  $M$  candidate solutions from the *stochastic policy*  $P_\theta(\cdot|\mathcal{G})$  by running the trained Ptr-A\* on a single test input  $\mathcal{G}$  and selects the one with the minimum expected energy consumption from  $M$  candidate outputs. The more we sample, the more likely we will get the better output. In the simulation, we set  $M = 51200$ , and this strategy is labeled as DRL-sampling.

### Algorithm 3: Active Search.

---

**Input:** Test input  $\mathcal{G}$ , steps  $S, \zeta$

- 1: Randomly sample a solution  $\mathcal{T}$  for  $\mathcal{G}$
- 2: Calculate  $E$  according to  $\mathcal{T}$  by A\*
- 3:  $O \leftarrow E$
- 4: **for**  $s = 1$  to  $S$  **do**
- 5:    $\mathcal{T}_i \sim \text{Sample solutions } P_\theta(\cdot|\mathcal{G}), \forall i \in \{1, \dots, Q\}$
- 6:    $E(\mathcal{T}_i|\mathcal{G}) \leftarrow \text{argmin}(E(\mathcal{T}_1|\mathcal{G}), \dots, E(\mathcal{T}_Q|\mathcal{G}))$
- 7:   **if**  $E(\mathcal{T}_i|\mathcal{G}) < E$  **then**
- 8:      $\mathcal{T} \leftarrow \mathcal{T}_i$
- 9:      $E \leftarrow E(\mathcal{T}_i|\mathcal{G})$
- 10:   **end if**
- 11:    $d\theta \leftarrow \frac{1}{Q} \sum_{i=1}^Q (E(\mathcal{T}_i|\mathcal{G}) - O) \nabla_\theta \log p_\theta(\mathcal{T}_i|\mathcal{G})$
- 12:    $\theta \leftarrow \text{Adam}(\theta, d\theta)$
- 13:    $O \leftarrow \zeta O + (1 - \zeta) \frac{1}{Q} \sum_{i=1}^Q V_\psi(\mathcal{G})$
- 14: **end for**
- 15: **return**  $\mathcal{T}, E$

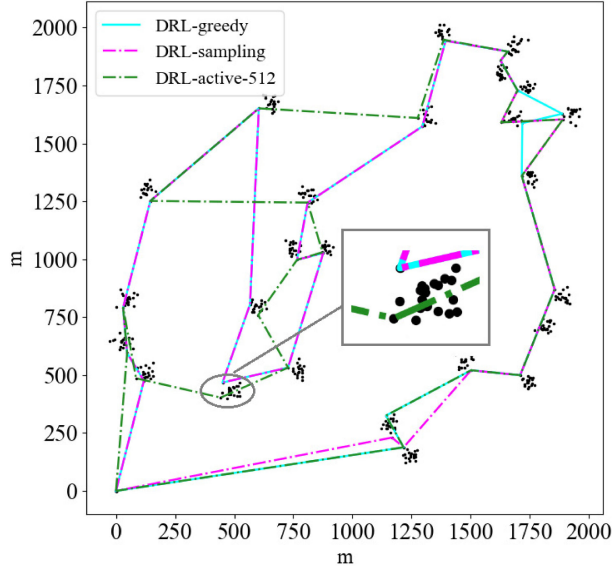
---

3) *Active Search*: Unlike the greedy search and the sampling search, this strategy can refine the parameter  $\theta$  of the Ptr-A\* during inference to minimize the expected energy consumption on a single test input  $\mathcal{G}$ . Active search samples multiple solutions  $\mathcal{T}_1, \dots, \mathcal{T}_Q$  from  $P_\theta(\cdot|\mathcal{G})$  for a single test input  $\mathcal{G}$  and uses policy gradients to refine  $\theta$  [22]. The process is presented in Algorithm 3. In the simulation, we sample three different sets of candidate solutions, {512, 5120, 10240}, which are labeled with DRL-active-512, DRL-active-5120, and DRL-active-10 240, respectively.

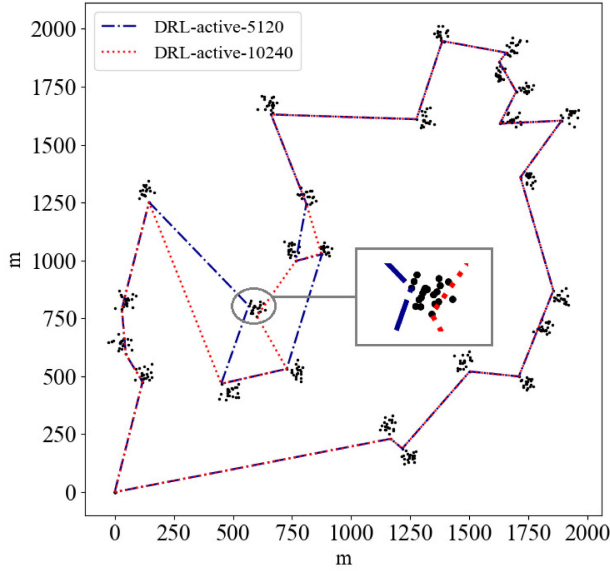
### C. Small-Scale Clusters

To thoroughly evaluate the performance of the proposed DRL algorithm, we first test the trained 20-clusters model on small-scale clusters. Since  $\omega$  in (21) is the weighting coefficient, its value does not impact the comparison results among algorithms. When  $\omega = 0$ , our optimization problem only considers the energy consumption of the UAV, which mainly depends on the flying distance of the UAV. Fig. 3 illustrates how the proposed DRL algorithm performs with different search strategies on the 25-clusters problem instance. As can be seen, the trajectory generated by DRL-greedy is the longest (9241 m), while DRL-active-10 240 produces the shortest trajectory (8693 m) among strategies.

Next, we compare our proposed DRL algorithm having different decoding search strategies with the nearest neighbor (NN) heuristic [46] and the genetic algorithm [47]. We first investigate the energy consumption comparison between our proposed DRL on the trained 20-clusters model and two baselines when  $\omega = 0.5$ . The genetic algorithm runs for 4000 generations, the chance of mutation is 0.5%, and the size of population is 150. In Fig. 4, we plot the average ratios of the energy consumption of our proposed DRL algorithm with different search strategies and two baselines to the energy consumption of DRL-active-10 240 versus different numbers of clusters  $K$ . Although the model is trained on 20-clusters problem instances, it still obtains good



(a)



(b)

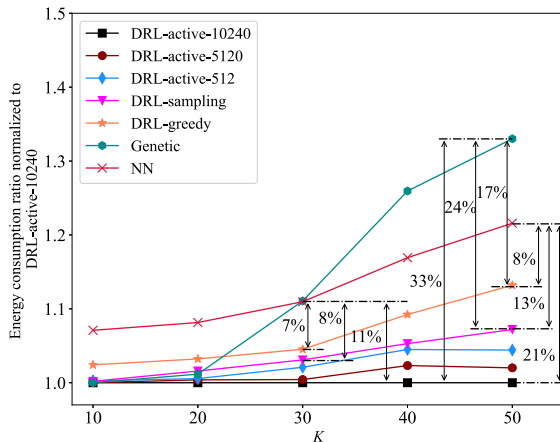
Fig. 3. Trajectories comparison on 25 clusters test instance when  $\omega = 0$ .

Fig. 4. Energy consumption comparison on small-scale clusters.

TABLE II  
RUNNING TIME COMPARISON ON SMALL-SCALE CLUSTERS

Algorithm	$K$	Time (s)				
		10	20	30	40	50
DRL-active-10240		17.1	57.15	121.62	209.27	319.64
DRL-active-5120		8.59	29.12	60.85	105.82	160.68
DRL-active-512		0.98	3.4	7.05	12.11	18.48
DRL-sampling		14.26	29.34	43.39	59.89	81.3
DRL-greedy		0.31	0.63	1.07	1.68	2.46
Genetic		60.37	64.85	73.35	80.21	90.91
NN		1.12	1.12	1.13	1.15	1.16

performance on the 10-clusters, 30-clusters, 40-clusters, and 50-clusters networks. This shows that the proposed DRL algorithm achieves an excellent generalization ability with respect to the number of clusters used for training. When  $K = 10$ , genetic, DRL-active-10 240, DRL-active-5120, DRL-active-512, and DRL-sampling obtain almost the same energy consumption result; however, the NN algorithm has higher energy consumption when compared with our proposed DRL with active search and sampling search strategies. As the number of clusters increases, the energy consumption savings of our proposed algorithm increase when compared to the NN and genetic algorithms. For example, when  $K = 30$ , the energy consumptions of UAV-WSN produced by NN and genetic algorithms are almost equal, which is about 11% more than that of DRL-active-10 240, 8% more than that of DRL-sampling, and 7% more than that of DRL-greedy. When the number of clusters increases to 50, the energy consumption of NN is around 21% more than that of DRL-active-10 240, 13% more than that of DRL-sampling, and 8% more than that of DRL-greedy. Likewise for  $K = 50$ , the energy consumption of the genetic algorithm is around 33% more than that of DRL-active-10 240, 24% more than that of DRL-sampling, and 17% more than that of DRL-greedy. It can be seen that our proposed DRL algorithm using any of three active search strategies can achieve better results than other search strategies and algorithms. This is because the active search strategy can refine the parameters of the Ptr-A\* model for producing the best solution while searching for candidate solutions on a single test instance at inference. From Fig. 4, we can see that DRL-sampling also obtains a relatively competitive result.

Table II compares the running time at inference. As the number of clusters increases, the running time of the proposed DRL algorithm with all strategies and two baseline techniques increases. Although DRL-active-10 240 obtains the best performance in reducing the energy consumption as can be seen from Fig. 4, it has the longest running time. This is because it needs more iterations to update the parameters. If the number of candidate solutions is relatively small, active search strategy tends to spend less time to produce the solution, like DRL-active-512. In addition, we can see that the running time of NN is always minimal among all algorithms for all values of  $K$ . DRL-greedy's running time is comparable with that of NN. Meanwhile, the

running times of DRL-sampling and DRL-active-512 are lower when compared with the genetic algorithm.

#### D. Large-Scale Clusters

In this subsection, we test the performance of the trained models on large-scale clusters test instances. The genetic algorithm runs for 10 000 generations. We first observe the results of the proposed DRL algorithm on the trained 40-clusters model, as shown in Fig. 5 (a). Clearly, our proposed DRL algorithm exhibits much better performances than the two baseline techniques in reducing the energy consumption of the UAV-WSN system. As the number of clusters increases, there is an increasing performance gap between the proposed DRL algorithm and the baseline techniques. For instance, when  $K = 80$ , the genetic algorithm consumes 50%, 58%, and 76% more energy than when compared to DRL-greedy, DRL-sampling, and DRL-active-10 240, respectively. As the number of clusters increases to 100, the energy consumption of the UAV-WSN when using the genetic algorithm is 72% more than that of DRL-greedy, 82% more than that of DRL-sampling, and 103% more than that of DRL-active-10 240. NN also shows a similar trend to that of the genetic algorithm. In particular, its energy consumption is 47% more than that of DRL-active-10 240 when  $K = 80$  and increases to 55% more than the energy consumption of DRL-active-10 240 when  $K = 100$ . When compared with DRL-greedy, the extra amount of energy consumed by the UAV-WSN when using NN increases from 25% to 31% as the number of clusters increasing from 80 to 100. However, NN exhibits an obviously superior performance than the genetic algorithm. As we can see, the active search strategies outperform the greedy strategy and sampling strategy for large-scale clusters problems. DRL-sampling shows a slightly better performance than DRL-greedy, which is reasonable.

Next, we use the trained 20-clusters model to evaluate the performance of the proposed DRL algorithm on large-scale clusters test instances. In Fig. 5 (b), our proposed DRL algorithm with different search strategies still obtains relatively good results as compared to the two baseline techniques on large-scale problem instances when the value of  $K$  varies, while the results generated by active search strategies are still the best among all search strategies. For example, when  $K = 80$ , the energy consumption of the genetic algorithm is 48% more than that of DRL-active-10 240, 40% more than that of DRL-sampling, and 31% more than that of DRL-greedy. Although NN exhibits a better performance than the genetic algorithm, its energy consumption is 23% more than that of DRL-active-10 240, 16% more than that of DRL-sampling, and 8% more than that of DRL-greedy when  $K = 80$ . However, compared with the results obtained on the trained 40-clusters model, the 20-clusters model clearly shows inferior performance. When  $K = 100$ , the energy consumption gap between the genetic algorithm and DRL-active-10 240 is 103% on the 40-clusters model, but this gap decreases to 65% on the 20-clusters model. Likewise, the

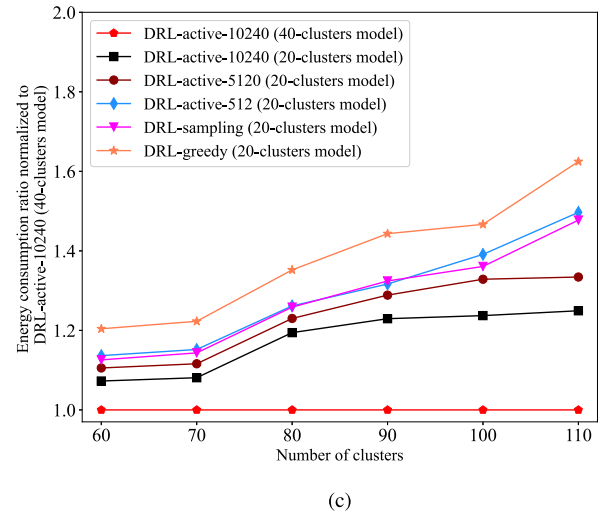
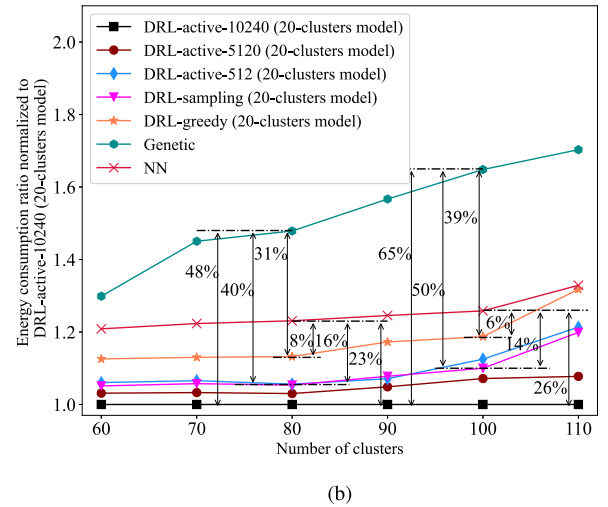
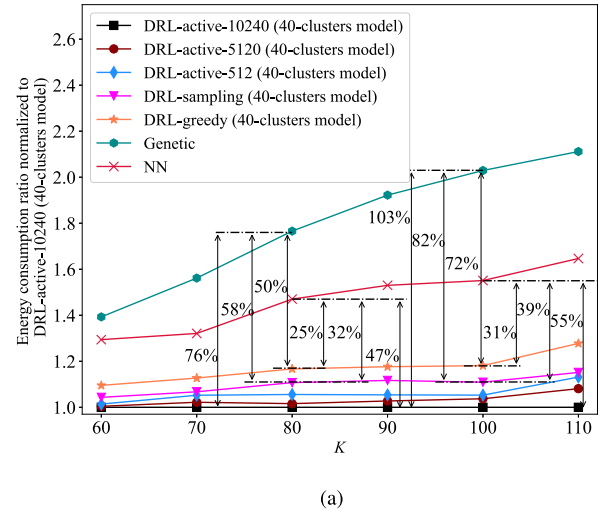


Fig. 5. Energy consumption comparison on large-scale clusters. (a) On 40-clusters model. (b) On 20-clusters model. (c) 40-clusters model vs. 20-clusters model.



TABLE III  
RUNNING TIME COMPARISON ON LARGE-SCALE CLUSTERS

Algorithm \ $K$		Time (s)					
		60	70	80	90	100	110
DRL-active-10240	20-clusters model	448.65	584.71	754.23	920.83	1149.36	1329.96
	40-clusters model	449.01	584.81	755.37	920.61	1150.05	1330.28
DRL-active-5120	20-clusters model	225.98	293.52	382.03	480.87	590.5	705.11
	40-clusters model	226.26	294.31	382.15	480.92	590.04	705.36
DRL-active-512	20-clusters model	24.99	33.08	44.08	55.54	67.21	79.31
	40-clusters model	25.75	33.57	44.1	55.78	67.54	79.88
DRL-sampling	20-clusters model	102.06	148.72	171.65	210.4	257.64	306.85
	40-clusters model	102.47	148.79	171.9	210.29	258.53	307.11
DRL-greedy	20-clusters model	3.11	4.23	5.84	7.83	9.96	12.55
	40-clusters model	3.39	4.54	5.81	7.48	10.2	12.89
Genetic		213.39	217.06	225.5	233.88	239.54	247.19
NN		1.17	1.17	1.18	1.19	1.2	1.22

gap in the energy consumption between the genetic algorithm and DRL-active-10 240 decreases from 76% on the 40-clusters model to 48% on the 20-clusters model when  $K = 80$ .

To investigate the difference between the two trained models, we compare the results obtained on the 20-clusters model with the results of DRL-active-10 240 on the 40-clusters model. As shown in Fig. 5 (c), DRL-active-10 240 (40-clusters model) clearly exhibits superior performance in reducing the energy consumption than three search strategies on the 20-clusters model. This performance merit constantly increases as the value of  $K$  increases. Thus, the trained 40-clusters model is more suitable for solving the trajectory planning problem for the UAV in cases of large-scale clusters.

The running time comparison of different strategies and algorithms at inference on large-scale problem instances is provided in Table III. We can observe that for a given number of clusters, the same strategy running on different models produces solution in almost the same amount of time. For example, DRL-active-10 240 takes 448.65 s on the 20-clusters model and 449.01 s on the 40-clusters model. DRL-active-10 240 takes the longest time among all the algorithms because it requires more iterations to refine the parameters of the Ptr-A\* to produce the best performance. The running time of DRL-sampling is acceptable in comparison with the genetic algorithm given its performance merits. Furthermore, the computation time spent by DRL-greedy is the least among all strategies, which is also significantly less than the running time of the genetic algorithm and slightly more than the time spent by NN. Although the NN algorithm spends the least amount of time, it produces worse results.

## V. CONCLUSION

In this paper, we investigated the problem of designing the UAV's trajectory for a clustered WSN to minimize the total energy consumption in the UAV-WSN system. Inspired by the recent developments of DRL, we propose a novel DRL-based method to solve the UAV's trajectory planning problem. Because the visiting order of the UAV to clusters can be regarded as

a sequential decision problem, we design a Ptr-A\* model to produce the trajectory of the UAV. The pointer network of the proposed Ptr-A\* model is used to determine the visiting order to clusters. Then, a search graph for all clusters is built according to the visiting order. The A\* search algorithm is utilized to quickly find the CH for each cluster from the search graph with the aim of minimizing the energy consumption of the UAV-WSN system. In order to obtain optimal parameters of the Ptr-A\*, we employ the model-free RL method to train the proposed Ptr-A\* model in an unsupervised manner. Lastly, we propose three search strategies at inference.

We conduct comprehensive experiments to evaluate the performance of the proposed DRL algorithm. The simulation results show that the proposed DRL algorithm with different search strategies can produce better trajectories for the UAV when compared with the baseline techniques. In particular, DRL-active-10 240 always produces the best results with different numbers of clusters of test instances. We also analyze the impact of different trained models on the results. The trained 40-clusters model is shown to be able to solve the trajectory planning problem of the UAV on large-scale clusters problems. The proposed DRL algorithm offers an appealing balance between performance and complexity. A key advantage of our proposed DRL algorithm is its generalization ability with respect to the number of clusters used for training. The model can be trained on small-scale clusters for faster training, and then can be used to solve larger-scale clusters problems. This makes it clearly more suitable for solving large-scale clusters problems as compared to the baseline techniques.

As for future work, we are interested in exploring other search strategies at inference to further improve the performance. It would also be interesting to develop a distributed DRL algorithm based on the Ptr-A\* model to solve multiple UAVs' trajectory planning problem jointly. We also plan to investigate and improve the generalization capability of the proposed DRL algorithm, i.e., when the number of nodes per cluster at inference is significantly larger than that used for training.

## APPENDIX A

### EXAMPLE OF ATTENTION MECHANISM

Here, we give a detailed numerical example of 3-clusters network to explain how the attention mechanism works. In Fig. 2, the input sequence  $\mathcal{G} = \{b_0, G_1, G_2, G_3\}$  is transformed into a sequence of latent memory states  $\mathcal{E} = \{e_0, e_1, e_2, e_3\}$ , which is the input of the decoder network. At decoding step 0, we calculate correlations between all elements in  $\mathcal{E}$  and the start tag  $v_{go}$  by (24) and (25), which can be expressed as

$$u_0^0 = \varphi \tanh(W_1 e_0 + W_2 h_0) \quad (\text{A.1})$$

$$u_1^0 = \varphi \tanh(W_1 e_1 + W_2 h_0) \quad (\text{A.2})$$

$$u_2^0 = \varphi \tanh(W_1 e_2 + W_2 h_0) \quad (\text{A.3})$$

$$u_3^0 = \varphi \tanh(W_1 e_3 + W_2 h_0). \quad (\text{A.4})$$

The learning parameters, namely  $\varphi$ ,  $W_1$ ,  $W_2$ , are initialized by the Xavier initialization method and trained by the Adam optimizer as explained in Section IV.A. Then, the *softmax* function is used to normalize the vector  $u^0 = \{u_0^0, u_1^0, u_2^0, u_3^0\}$ . For the sake of illustration, we assume that the four elements in  $u^0$  determine the following four conditional probability values:

$$\begin{aligned} P(\pi_0 = 0|\mathcal{G}) &= \frac{\exp(u_0^0)}{\exp(u_0^0) + \exp(u_1^0) + \exp(u_2^0) + \exp(u_3^0)} = 0.6 \end{aligned} \quad (\text{A.5})$$

$$\begin{aligned} P(\pi_0 = 1|\mathcal{G}) &= \frac{\exp(u_1^0)}{\exp(u_0^0) + \exp(u_1^0) + \exp(u_2^0) + \exp(u_3^0)} = 0.1 \end{aligned} \quad (\text{A.6})$$

$$\begin{aligned} P(\pi_0 = 2|\mathcal{G}) &= \frac{\exp(u_2^0)}{\exp(u_0^0) + \exp(u_1^0) + \exp(u_2^0) + \exp(u_3^0)} = 0.1 \end{aligned} \quad (\text{A.7})$$

$$\begin{aligned} P(\pi_0 = 3|\mathcal{G}) &= \frac{\exp(u_3^0)}{\exp(u_0^0) + \exp(u_1^0) + \exp(u_2^0) + \exp(u_3^0)} = 0.2. \end{aligned} \quad (\text{A.8})$$

Since  $e_0$  has the highest conditional probability value at decoding step 0, the output  $\pi_0$  of this step points to the first element of  $\mathcal{G}$ ,  $b_0$ . At decoding step 1, we use the same approach to calculate the correlations between  $b_0$  and the remaining elements in  $\mathcal{E}$  as follows

$$u_1^1 = \varphi \tanh(W_1 e_1 + W_2 h_1) \quad (\text{A.9})$$

$$u_2^1 = \varphi \tanh(W_1 e_2 + W_2 h_1) \quad (\text{A.10})$$

$$u_3^1 = \varphi \tanh(W_1 e_3 + W_2 h_1). \quad (\text{A.11})$$

Similarly, by using the *softmax* function, we calculate the following conditional probabilities based on the obtained elements

$$\{u_1^1, u_2^1, u_3^1\}:$$

$$\begin{aligned} P(\pi_1 = 1|\pi_0, \mathcal{G}) &= \frac{\exp(u_1^1)}{\exp(u_2^1) + \exp(u_2^1) + \exp(u_3^1)} = 0.2 \end{aligned} \quad (\text{A.12})$$

$$\begin{aligned} P(\pi_1 = 2|\pi_0, \mathcal{G}) &= \frac{\exp(u_2^1)}{\exp(u_2^1) + \exp(u_2^1) + \exp(u_3^1)} = 0.1 \end{aligned} \quad (\text{A.13})$$

$$\begin{aligned} P(\pi_1 = 3|\pi_0, \mathcal{G}) &= \frac{\exp(u_3^1)}{\exp(u_2^1) + \exp(u_2^1) + \exp(u_3^1)} = 0.7. \end{aligned} \quad (\text{A.14})$$

The output  $\pi_1$  of this step points to  $G_3$  because the conditional probability of  $e_3$  is maximum. Then, the process repeats until we obtain the full output sequence of the decoder network as  $\{b_0, G_3, G_1, G_2\}$ , i.e.,  $\{\pi_0, \pi_1, \pi_2, \pi_3\}$ .

## REFERENCES

- [1] M. Mozaffari *et al.*, "A tutorial on UAVs for wireless networks: Applications, challenges, and open problems," *IEEE Commun. Surv. Tut.*, vol. 21, no. 3, pp. 2334–60, Jul.–Sep. 2019.
- [2] Y. Zeng, R. Zhang, and T. J. Lim, "Wireless communications with unmanned aerial vehicles: Opportunities and challenges," *IEEE Commun. Mag.*, vol. 54, no. 5, pp. 36–42, May 2016.
- [3] D. A. Hedges, J. P. Coon, and G. Chen, "A continuum model for route optimization in large-scale inhomogeneous multi-hop wireless networks," *IEEE Trans. Commun.*, vol. 68, no. 2, pp. 1058–1070, Feb. 2020.
- [4] D. H. Tran, T. X. Vu, S. Chatzinotas, S. ShahbazPanahi, and B. Ottersten, "Coarse trajectory design for energy minimization in UAV-enabled," *IEEE Trans. Veh. Technol.*, vol. 69, no. 9, pp. 9483–9496, Sep. 2020.
- [5] 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.
- [6] C. Zhan and H. Lai, "Energy minimization in internet-of-things system based on rotary-wing UAV," *IEEE Wireless Commun. Lett.*, vol. 8, no. 5, pp. 1341–1344, Oct. 2019.
- [7] J. Baek, S. I. Han, and Y. Han, "Energy-efficient UAV routing for wireless sensor networks," *IEEE Trans. Veh. Technol.*, vol. 69, no. 2, pp. 1741–1750, Feb. 2020.
- [8] M. B. Ghorbel *et al.*, "Joint position and travel path optimization for energy efficient wireless data gathering using unmanned aerial vehicles," *IEEE Trans. Veh. Technol.*, vol. 68, no. 3, pp. 2165–2175, Mar. 2019.
- [9] D. Yang, Q. Wu, Y. Zeng, and R. Zhang, "Energy trade-off in ground-to-UAV communication via trajectory design," *IEEE Trans. Veh. Technol.*, vol. 67, no. 7, pp. 6721–6726, Jul. 2018.
- [10] B. Zhu, E. Bedeer, H. H. Nguyen, R. Barton, and J. Henry, "Improved soft- $k$ -means clustering algorithm for balancing energy consumption in wireless sensor networks," *IEEE Internet Things J.*, vol. 8, no. 6, pp. 4868–4881, Mar. 2021.
- [11] M. Samir *et al.*, "UAV trajectory planning for data collection from time-constrained IoT devices," *IEEE Trans. Wireless Commun.*, vol. 19, no. 1, pp. 34–46, Jan. 2020.
- [12] S. Zhang, S. Shi, S. Gu, and X. Gu, "Power control and trajectory planning based interference management for UAV-assisted wireless sensor networks," *IEEE Access*, vol. 8, pp. 3453–3464, 2019.
- [13] A. A. Al-Habob, O. A. Dobre, S. Muhaidat, and H. V. Poor, "Energy efficient data dissemination using a UAV: An ant colony approach," *IEEE Wireless Commun. Lett.*, vol. 10, no. 1, pp. 16–20, Jan. 2021.
- [14] K. Zhu, X. Xu, and S. Han, "Energy-efficient UAV trajectory planning for data collection and computation in mMTC networks," in *Proc. IEEE Globecom Workshops*, 2018, pp. 1–6.
- [15] C. H. Liu, Z. Chen, J. Tang, J. Xu, and C. Piao, "Energy-efficient UAV control for effective and fair communication coverage: A deep reinforcement learning approach," *IEEE J. Sel. Areas Commun.*, vol. 36, no. 9, pp. 2059–2070, Sep. 2018.

- [16] Y. Yuan, L. Lei, T. X. Vu, S. Chatzinotas, and B. Ottersten, "Actor-critic deep reinforcement learning for energy minimization in UAV-aided networks," in *Proc. Eur. Conf. Netw. Commun.*, 2020, pp. 348–352.
- [17] R. Ding, F. Gao, and X. S. Shen, "3D UAV trajectory design and frequency band allocation for energy-efficient and fair communication: A deep reinforcement learning approach," *IEEE Trans. Wireless Commun.*, vol. 19, no. 12, pp. 7796–7809, Dec. 2020.
- [18] B. Zhang *et al.*, "Learning-based energy-efficient data collection by unmanned vehicles in smart cities," *IEEE Trans. Ind. Inform.*, vol. 14, no. 4, pp. 1666–1676, Apr. 2018.
- [19] W. Liu, P. Si, E. Sun, M. Li, C. Fang, and Y. Zhang, "Green mobility management in UAV-assisted IoT based on dueling DQN," in *Proc. IEEE Int. Conf. Commun.*, 2019, pp. 1–6.
- [20] H. Zhang, J. Li, Y. Li, and H. Yue, "Understanding subtitles by character-level sequence-to-sequence learning," *IEEE Trans. Ind. Inform.*, vol. 13, no. 2, pp. 616–624, Apr. 2017.
- [21] O. Vinyals, M. Fortunato and N. Jaitly, "Pointer networks," in *Proc. Adv. Neural Inf. Process. Syst.*, 2015, pp. 2692–2700.
- [22] I. Bello *et al.*, "Neural combinatorial optimization with reinforcement learning," 2016, *arXiv:1611.09940*.
- [23] J. Q. Yu, W. Yu, and J. Gu, "Online vehicle routing with neural combinatorial optimization and deep reinforcement learning," *IEEE Trans. Intell. Transp. Syst.*, vol. 20, no. 10, pp. 3806–3817, Oct. 2019.
- [24] Z. Li, J. Wu, L. Sun, and T. Rong, "Combinatorial keyword recommendations for sponsored search with deep reinforcement learning," in *Proc. AdKDD*, 2019, pp. 1–6.
- [25] M. Nazari, A. Oroojlooy, L. Snyder, and M. Takác, "Reinforcement learning for solving the vehicle routing problem," in *Proc. Adv. Neural Inf. Process. Syst.*, 2018, pp. 9839–9849.
- [26] V. Razo and H. Jacobsen, "Smart charging schedules for highway travel with electric vehicles," *IEEE Trans. Transp. Electrification*, vol. 2, no. 2, pp. 160–173, Jun. 2016.
- [27] Z. Ma, B. Ai, R. He, G. Wang, Y. Niu, and Z. Zhong, "A wideband non-stationary air-to-air channel model for UAV communications," *IEEE Trans. Veh. Technol.*, vol. 69, no. 2, pp. 1214–1226, Feb. 2020.
- [28] Z. Lian, L. Jiang, C. He, and D. He, "A non-stationary 3-D wideband GBSM for HAP-MIMO communication systems," *IEEE Trans. Veh. Technol.*, vol. 68, no. 2, pp. 1128–1139, Feb. 2019.
- [29] D. Hulens, J. Verbeke, and T. Goedemé, "Choosing the best embedded processing platform for on-board UAV image processing," in *Proc. VISAPP*, 2015, pp. 455–472.
- [30] W. Heinzelman *et al.*, "An application-specific protocol architecture for wireless microsensor networks," *IEEE Trans. Wireless Commun.*, vol. 1, no. 4, pp. 660–670, Oct. 2002.
- [31] T. Rappaport, *Wireless Communications: Principles and Practice*. Englewood Cliffs, NJ, USA: Prentice-Hall, 1996.
- [32] R. Roberti and P. Toth, "Models and algorithms for the asymmetric traveling salesman problem: An experimental comparison," *EURO J. Transp. Logistics*, vol. 1, no. 1, pp. 113–133, Jun. 2012.
- [33] J. Puchinger and G. R. Raidl, "Combining metaheuristics and exact algorithms in combinatorial optimization: A survey and classification," *Artif. Intell. Knowl. Eng. Appl.: A BioInspired Approach*, vol. 3562, pp. 41–53, Jun. 2005.
- [34] D. Gong and Y. Yang, "Low-latency SINR-based data gathering in wireless sensor networks," *IEEE Trans. Wireless Commun.*, vol. 13, no. 6, pp. 3207–3221, Jun. 2014.
- [35] L. Liu and U. Mitra, "On sampled reinforcement learning in wireless networks: Exploitation of policy structures," *IEEE Trans. Commun.*, vol. 68, no. 5, pp. 2823–2837, May 2020.
- [36] I. Bello *et al.*, "Seq2Slate: Re-ranking and slate optimization with RNNs," in *Proc. Int. Conf. Learn. Representations*, 2019, pp. 1–12.
- [37] W. Koehrsen, "Neural network embeddings explained," *Towards Data Science*, 2018. [Online]. Available: <https://towardsdatascience.com/neural-network-embeddings-explained-4d028e6f0526>
- [38] Y. Zhu, X. Dong, and T. Lu, "An adaptive and parameter-free recurrent neural structure for wireless channel prediction," *IEEE Trans. Commun.*, vol. 67, no. 11, pp. 8086–8096, Nov. 2019.
- [39] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, "Show and tell: A neural image caption generator," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2015, pp. 3156–3164.
- [40] V. R. Konda and J. N. Tsitsiklis, "Actor-critic algorithms," in *Proc. Adv. Neural Inf. Process. Syst.*, 2000, pp. 1008–1014.
- [41] R. J. Williams, "Simple statistical gradient following algorithms for connectionist reinforcement learning," *Mach. Learn.*, vol. 8, pp. 229–256, 1992.
- [42] D. P. Kingma and J. L. Ba, "Adam: A method for stochastic optimization," in *Proc. Int. Conf. Learn. Represent.*, 2015, pp. 1–13.
- [43] Z. Tu, F. He, and D. Tao, "Understanding generalization in recurrent neural networks," in *Proc. Int. Conf. Learn. Representations*, 2020, pp. 1–16.
- [44] A. Al-Hourani, S. Kandeepan, and S. Lardner, "Optimal LAP altitude for maximum coverage," *IEEE Wireless Commun. Lett.*, vol. 3, no. 6, pp. 569–572, Dec. 2014.
- [45] H. Ghazzai, M. Ben Ghorbel, A. Kadri, M. J. Hossain, and H. Menouar, "Energy-efficient management of unmanned aerial vehicles for underlay cognitive radio systems," *IEEE Trans. Green Commun. Netw.*, vol. 1, no. 4, pp. 434–443, Dec. 2017.
- [46] B. Hu and G. Raidl, "Effective neighborhood structures for the generalized traveling salesman problem," in *Proc. Eur. Conf. Evol. Computation Combinatorial Optim.*, 2008, pp. 36–47.
- [47] J. Li, Q. Sun, M. Zhou, and X. Dai, "A new multiple traveling salesman problem and its genetic algorithm-based solution," in *Proc. IEEE Int. Conf. Syst. Man Cybern.*, 2013, pp. 627–632.



**Botao Zhu** is currently working toward the Ph.D. degree with the Department of Electrical and Computer Engineering, University of Saskatchewan, Saskatoon, SK, Canada. His current research interests include wireless sensor networks, Internet-of-Things, machine learning, and unmanned aerial vehicles.



**Ebrahim Bedeer** (Member, IEEE) received the B.Sc. (Hons.) and M.Sc. degrees from Tanta University, Tanta, Egypt and the Ph.D. degree from Memorial University, St. Johns, NL, Canada, all in electrical engineering.

In 2019, he joined the Department of Electrical and Computer Engineering, University of Saskatchewan, Saskatoon, SK, Canada, as an Assistant Professor. Before that, he was an Assistant Professor (Lecturer in the U.K.) with Ulster University, London, U.K., and a Postdoctoral Fellow with Carleton University, Ottawa, ON, Canada, and the University of British Columbia, Kelowna, BC, Canada. His current research interests include applications of optimization techniques in signal processing and wireless communications, spectral efficient communication systems, and Internet-of-Things (IoT).

He is an Associate Editor of the IEEE COMMUNICATIONS LETTERS. He has served on the Technical Program Committees of numerous main international communication conferences, such as IEEE GLOBECOM, IEEE ICC, and IEEE VTC. Dr. Bedeer was the recipient of numerous awards, including the Exemplary Reviewer of IEEE COMMUNICATIONS LETTERS and IEEE WIRELESS COMMUNICATIONS LETTERS.



**Ha H. Nguyen** (Senior Member, IEEE) received the B.Eng. degree from the Hanoi University of Technology, Hanoi, Vietnam, in 1995, the M.Eng. degree from the Asian Institute of Technology, Bangkok, Thailand, in 1997, and the Ph.D. degree from the University of Manitoba, Winnipeg, MB, Canada, in 2001, all in electrical engineering. In 2001, he joined the Department of Electrical and Computer Engineering, University of Saskatchewan, Saskatoon, SK, Canada, and became a Full Professor in 2007. He currently holds the position of NSERC/Cisco Industrial Research Chair in *Low-Power Wireless Access for Sensor Networks*. He is the Co-Author, with Ed Shwedyk, of the textbook *A First Course in Digital Communications* (published by Cambridge University Press). His research interests include communication theory, wireless communications, and statistical signal processing. He was an Associate Editor for the IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS and IEEE WIRELESS COMMUNICATIONS LETTERS during 2007–2011 and 2011–2016, respectively. He is currently an Associate Editor for the IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY. He served as a Technical Program Chair for numerous IEEE events. Dr. Nguyen is a Fellow of the Engineering Institute of Canada (EIC) and a registered member of the Association of Professional Engineers and Geoscientists of Saskatchewan (APEGS).





**Robert Barton** (Member, IEEE) is currently a Distinguished Engineer with Cisco's Digital Transformation and Innovation Group, where he has the role of the Chief Architect for Cisco Canada and Cisco's global IoT sales organization. He has worked in the IT industry for more than 23 years, the last 20 of which have been at Cisco. He Graduated from the University of British Columbia, Vancouver, BC, Canada, with a degree in engineering physics. He is a published author, with titles on the subjects of Network QoS, Wireless, IoT, Machine Learning and Data Analytics.

He has also contributed to many academic papers, and leads Cisco's university research partnership program. He also holds many patents in the areas of wireless communications, segment routing, and machine learning. His current research interests include wireless communications, industrial networking, IoT, and AI/ML in networking systems.



**Jerome Henry** (Senior Member, IEEE) received the master's degree from Universite Sorbonne Nouvelle, Paris, France, in 1994. He is currently a Principal Engineer with the office of the Wireless CTO, Cisco Systems, Toronto, ON, Canada. He has more than 20 years of experience in the wireless industry, focusing on optimization of short and long range performances. He has authored multiple titles on Wireless technologies, location, 802.11 security, IoT, QoS, and Machine Learning, and has also contributed to numerous academic papers on the same topics. His current

research interests include wireless communications and location with a particular focus on accuracy improvement through machine learning. He is currently the Vice-Chair in IEEE 802.11bi Task Group, the Chair of FiRa Localization Group, and is FiRa Board Member.