

# Differentiated Federated Reinforcement Learning Based Traffic Offloading on Space-Air-Ground Integrated Networks

Yeguang Qin, *Student Member, IEEE*, Yilin Yang, *Student Member, IEEE*, Fengxiao Tang, *Member, IEEE*, Xin Yao, *Member, IEEE*, Ming Zhao, *Member, IEEE*, and Nei Kato, *Fellow, IEEE*

**Abstract**—The Space-Air-Ground Integrated Network (SAGIN) plays a pivotal role as a comprehensive foundational network communication infrastructure, presenting opportunities for highly efficient global data transmission. Nonetheless, given SAGIN’s unique characteristics as a dynamically heterogeneous network, conventional network optimization methodologies encounter challenges in satisfying the stringent requirements for network latency and stability inherent to data transmission within this network environment. Therefore, this paper proposes the use of differentiated federated reinforcement learning (DFRL) to solve the traffic offloading problem in SAGIN, i.e., using multiple agents to generate differentiated traffic offloading policies. Considering the differentiated characteristics of each region of SAGIN, DFRL models the traffic offloading policy optimization process as the process of solving the Decentralized Partially Observable Markov Decision Process (DEC-POMDP) problem. The paper proposes a novel Differentiated Federated Soft Actor-Critic (DFSAC) algorithm to solve the problem. The DFSAC algorithm takes the network packet delay as the joint reward value and introduces the global trend model as the joint target action-value function of each agent to guide the update of each agent’s policy. The simulation results demonstrate that the traffic offloading policy based on the DFSAC algorithm achieves better performance in terms of network throughput, packet loss rate, and packet delay compared to the traditional federated reinforcement learning approach and other baseline approaches.

**Index Terms**—Space-air-ground Integrated Network (SAGIN), federated reinforcement learning (FRL), heterogeneous network, network optimization, traffic offloading.

## I. INTRODUCTION

THE optimized decision-making is a critical challenge in the Cyber-physical system, especially in the next-generation network of B5G/6G with a highly dynamic and large-scale environment. Reinforcement learning (RL) is one kind of machine learning technology to optimizes decision-making by maximizing the cumulative reward by continuously

exploring and exploiting the environment of an agent. Using RL to obtain network optimization strategies is a current research hotspot in the field of modern networks, e.g., network resources allocation [1], [2] and traffic control [3].

However, traditional RL applications still suffer from critical problems in scenarios with complex environments. For example, the local strategy trap when the network environment becomes dynamic and heterogeneous. The local strategy trap arises when decisions are derived from incomplete local information, overlooking the broader implications of global environmental change. This leads to ill-informed outcomes due to a lack of holistic consideration. Due to the features of high dynamic and heterogeneous complex networks, the local strategy trap occurs frequently in network scenarios such as Internet of Things (IoT) [4], [5], vehicular network [6], [7] and Space-Air-Ground Integrated Network (SAGIN) [8], [9].

Federated learning (FL) is an advanced learning technology to train a shared learning model without raw training data by considering the privacy of distributed devices in complex environments. By combining both RL and FL, a privacy-preserving multi-agent collaboration approach referred to as federated reinforcement learning (FRL) is proposed. In FRL, each agent trains the data separately, aggregates it into a uniform global policy model, and then distributes it to the agents, ensuring the user’s data privacy as much as possible while integrating local learning results from each agent.

The local strategy traps are solved to some extent by the shared training manner of FRL as the global information is implicitly shared during the integrated learning process. Unfortunately, the global sharing-based FRL employing a global learning model for local inference may still fall into the local strategy trap. The accuracy of using the global learning model for local inference is ensured by the assumption that the environments of diverse devices are homologous with the same state transition probability, which is not practical as the considered state transition probability is always differentiation in the heterogeneous environment.

For example, a dynamic and heterogeneous SAGIN consists of satellites, unmanned aerial vehicles (UAVs), ground-based stations in different regions, and user devices. In such a network environment, there is significant dynamism and heterogeneity among regions due to the high-speed movement of nodes, the extensive range of data transmission, and the convergence between different types of networks [10]. The differentiation between regions is generally due to network

This work was supported by the National Key R&D Program of China (Grant no.2021ZD0140301), National Key R&D Program of China (Project no.2020AAA0109602), Changsha Municipal Natural Science Foundation (Grant no.kq2208284), Hunan Provincial Natural Science Foundation (Grant no.2023jj40774), National Natural Science Foundation of China (Grant no.62302527), and the High Performance Computing Center of Central South University. (Corresponding author: Fengxiao Tang.)

Y. Qin, Y. Yang, F. Tang, X. Yao, and M. Zhao are with the School of Computer Science and Engineering, Central South University, ChangSha, China. Emails:{qinyeguang, yangyilin, tangfengxiao, xinyao, meanzhao}@csu.edu.cn

N. Kato is with Graduate School of Information Sciences (GSIS), Tohoku University, Sendai, Japan. Email:kato@it.is.tohoku.ac.jp

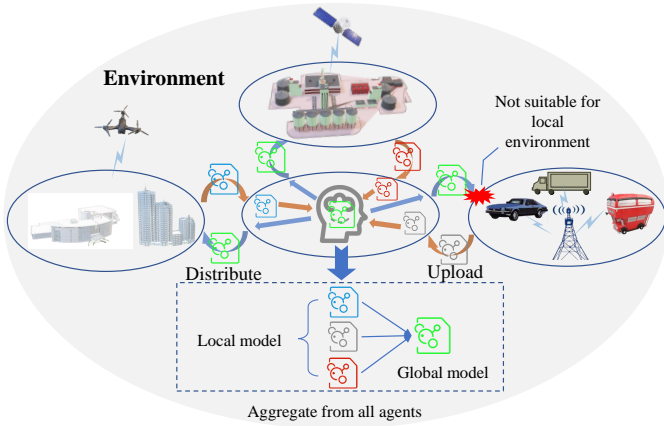


Fig. 1. The global learning model does not work well because of environmental differentiations.

state differences, sample distribution differences, and dynamic characteristic differences. These discrepancies lead to the global policy model obtained by the aggregation of agents in each region can only be suboptimal compared to the optimal policy in each region. It is not easy to use the traditional FRL approach to obtain a global policy model to address issues like traffic offloading in SAGIN that performs well across all differentiated regions.

As shown in Fig. 1, traditional FRL distributes a uniform global learning model to each region. When there are differentiations in the environment of each region, the decisions inferred by the global learning model may be misjudged. Therefore, considering the differentiations between regions in complex network environments, we propose a novel differentiated federated reinforcement learning (DFRL) approach. Instead of seeking a uniform global policy model, DFRL enables agents across regions to cooperate in training their respective policy models.

Precisely, the cooperation among the agents in each region is accomplished through the global trend model. The global trend model is generated by aggregating the trend model developed by all the agents involved in the cooperation, reflecting the state of the environment and its changing trend in each region. Then, under the guidance of the global trend model, the agents in each region get their local policy models by training with local data. The introduction of the global trend model allows agents to consider each region's differentiations based on the information from multiple regions to obtain local strategies that are more applicable to their regions.

The main work of the paper can be summarized as follows:

- In this paper, we study the traffic offloading problem in a dynamic heterogeneous network like SAGIN and model it as a DEC-POMDP decision problem with the objective of optimizing the network delay.
- Then, we propose a novel concept of differentiated federal reinforcement learning by introducing a trend model, and based on this concept, an algorithm named Differentiated Federated Soft Actor-Critic (DFSAC) is proposed to solve the joint policy in DEC-POMDP.
- We design a SAGIN struct and propose a dynamic

traffic offloading method based on the DFSAC algorithm to solve the traffic offloading problem under this complex network structure. Simulation results show that our method achieves better results regarding network throughput, packet loss rate, and delay than other methods.

The rest of the paper is organized as follows. In section II, we discuss related work. In section III, we detail the system model and the formulation of the problem. Section IV elaborates on the concept of DFRL and proposes the DFSAC algorithm. An empirical study of traffic offloading in SAGIN is presented in Section V and summarized in Section VI.

## II. RELATED WORK

### A. Federated Learning

FL is a distributed machine learning approach in privacy-preserving scenarios, where the critical point is that information is passed between collaborators by sharing models rather than data [11]. In recent years, researchers have proposed several solutions to the challenges faced in FL. For example, to obtain higher communication efficiency, McMahan et al. proposed FedAVG [12], which is now widely used as the baseline for FL research. FedAVG is a widely adopted federated learning algorithm that aggregates model parameters using weighted averaging by uploading local model parameters to a central server, computing the average of all model parameters, and then broadcasting this average to all local devices. To address the heterogeneity problem in FL, Yuan, Dinh, and Ruan et al. proposed [13], [14], and [15], respectively, to solve data heterogeneity, model heterogeneity, and device heterogeneity. Moreover, Qu et al. further analyzed the convergence of FL methods in [16], comprehensively studying how the convergence of FedAVG varies with the number of participating devices in the FL setting. They demonstrate the convergence of the FedAVG method in several differentiated scenarios and perform a comprehensive study of its convergence rate. This work provides a solid foundation for federated learning in further research.

### B. Federated Reinforcement Learning

RL is a branch of machine learning (ML). Compared with other machine learning methods such as supervised or unsupervised learning methods, RL generates samples and learns through these samples through constant interaction between the agent and the environment [17]. The RL is widely used for communication and networking optimization in various networks. Wang et al. [18] proposed a semantic communication framework for efficient textual data transmission in resource-constrained wireless networks. This framework utilizes knowledge graphs and a proximal-policy-optimization-based RL algorithm integrated with an attention network to optimize the partial transmission of semantic information and enhance semantic similarity metrics. Luan et al. [19] employed deep reinforcement learning (DRL) to ascertain the optimal traffic allocation ratio among multiple controllable paths for source-destination pairs. Gao et al. [20] proposed a deep reinforcement learning-based framework for joint optimization

of computing, pushing, and caching in mobile edge computing networks, effectively reducing transmission bandwidth and computing cost through dynamic orchestration and proactive content delivery. However, when the environment gets bigger, it becomes difficult for a single agent to complete the complex task in these scenarios due to constraints such as information processing capabilities or learning efficiency. The natural idea is to set up multiple agents in the environment and complete the task by cooperating among them. Hence, researchers have introduced distributed RL and parallel RL [21] with the aim of accelerating the learning of optimal strategies for single-agent RL problems through parallel hardware utilization. However, these approaches also raise concerns about potential agent privacy breaches.

In this context, the integration of FL, a paradigm centered on preserving privacy, into the realm of reinforcement RL represents a notable development. FRL is a novel, distributed, and collaborative methodology that effectively amalgamates the principles of FL and RL. In this paradigm, each agent trains data locally and builds a shared model, which protects agent privacy and accelerates agent learning efficiency [22].

### C. Federated Learning in Heterogeneous Network

Environmental heterogeneity is a significant challenge for FL in real application scenarios. For this challenge, Yuan et al. proposed an improved Federated Deep AUC (area under the ROC curve) Maximization algorithm in [13] to solve the data heterogeneity problem, while Hanzely et al. proposed a hybrid global model and local model to achieve balanced training in [23] to solve the conflicts caused by data heterogeneity. Shen et al. proposed a novel distributed learning scheme [24] that combines FL with a model-splitting mechanism to accommodate customer heterogeneity. In addition, researchers have also tried to use personalized federated learning methods [25]–[27] to address the challenges caused by environment heterogeneity.

In communication networks, Liu et al. proposed a method for Radio Access Network (RAN) slicing using the FRL approach in [28], which improves the communication efficiency and throughput of the network. Kwon et al. proposed an FRL-based resource allocation method for Internet of Things Underwater (IoUT) devices in [29], which has significant advantages over the single-agent DRL approach. Zhu et al. proposed a fast convergent federated-based dynamic task offloading method in the power grid Internet of Things [30] but needed to sufficiently consider the differentiation of environments. Wang et al. [31] combined mobile edge computing with fiber-wireless networks and applied a reputation-based FRL strategy, effectively optimizing network performance and resource allocation in IoT deployments while protecting user privacy.

The current methods for addressing environmental heterogeneity primarily involve various forms of aggregation or training of local models. However, these methods lack the capability to dynamically adapt to environmental changes. As large-scale network environments become increasingly dynamic and heterogeneous, these methods struggle to effectively adapt to variations in regional network conditions and network

types. This leads to reduced learning efficiency and potential convergence issues, posing a significant challenge.

### D. Offloading Solutions in SAGIN

Guo et al. [32] have devised a comprehensive energy-efficient optimization strategy that integrates UAV-assisted communication with Mobile Edge Computing (MEC), effectively reducing overall energy consumption through strategic task offloading, transmission bit allocation, and UAV trajectory planning. Li et al. [33] proposed an integrated satellite/terrestrial collaborative transmission scheme, incorporating cache-enabled Low Earth Orbit (LEO) satellites as part of the RAN, which offloads traffic from base stations through satellite broadcast transmission, thereby achieving energy-efficient RAN operations.

The existing research predominantly addresses scenarios involving either satellites or UAVs in isolation; however, given the inherent complexity, heterogeneity, and dynamic nature of SAGIN, these systems present considerable challenges in terms of accurate modeling. As a result, the model-free methodology of reinforcement learning has been extensively investigated and applied to effectively tackle the offloading issues within the SAGIN framework. Cheng et al. [34] have developed an innovative SAGIN edge/cloud computing architecture, tailored for offloading computation-intensive applications while considering remote energy and computational constraints. Tang et al. [8] have presented a reinforcement learning-based traffic offloading scheme, taking into account the substantial node mobility within SAGIN, as well as the frequent variations in network traffic and link statuses. Zhang et al. [35] proposed a Learning-based Orbital Edge Offloading (LOEF) method using multi-agent learning, enabling UAVs to coordinate and learn optimal offloading strategies for computational task scheduling in the Internet of Remote Things (IoRT) within SAGIN.

Although existing solutions have made significant progress in SAGIN, they often overlook the critical aspect of privacy protection. In this context, FRL offers a potential mechanism for privacy preservation. However, the application of FRL faces inherent limitations, mainly because it relies on a unified global policy model to guide the reasoning process of local agents. This approach can lead to suboptimal results in local reasoning, thereby affecting the quality of overall decision-making. Therefore, applying FRL in dynamic and heterogeneous SAGIN environments to effectively solve the traffic offloading problem remains a significant challenge.

## III. SYSTEM MODEL

### A. Network Model

In this paper, we consider traffic offloading in SAGIN and establish a multi-dimensional heterogeneous network model: The ground network is composed of mobile user equipment (UE), base stations (BS) and edge base stations. The air network includes a series of dynamically deployable UAVs as an extension of the ground network. The space network is a double-layer communication satellite network composed

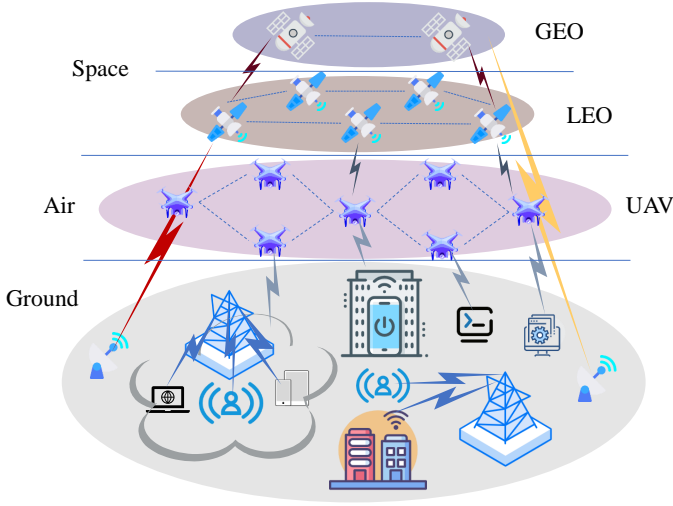


Fig. 2. The four-layer SAGIN architecture.

of LEO and Geostationary Earth Orbit satellites (GEO). The network structure is shown in Fig. 2.

Due to the uneven distribution of UE in the ground network, the traffic may be over the load of the BS in service-intensive areas. Therefore, we use UAV as a mobile flying base station to dynamically relay the offloaded traffic to other regions to avoid network congestion. We set the trajectory of the UAV as circular motion or random movement. When the UAV flies to an area, it can be used as the traffic relay node of the area. In this paper, the trajectory optimization of UAV in SAGIN is not our focus, so we do not introduce it additionally. However, the signal coverage of UAV is limited. In the process of moving, UAV may fly out of the target area, resulting in signal loss and a sharp increase in the packet loss rate. On the other hand, the loadable traffic of UAV is relatively limited, and a large number of dropped packets exceeding the capacity will also cause an increase in packet loss rate. Therefore, the communication satellite with large coverage is introduced as the available traffic relay node, and the ground traffic can be offloaded to available UAVs or communication satellites. Our network model further considers the double-layer satellite network composed of LEO and GEO based on the existing research on traffic offloading in SAGIN. The two-layer satellite network expands the available resource pool of SAGIN, and the use of GEO also improves the stability of the network and the overall performance.

In this paper, we model a four-layer SAGIN as an undirected graph  $G = (V, E)$ . Among them,  $V = \{V^D, V^B, V^U, V^L, V^G\}$  represents different types of nodes in SAGIN.  $V^D = \{v_1^D, v_2^D, \dots, v_n^D\}$  is a collection of UE nodes,  $n$  is the number of this type nodes. Similarly,  $V^B$  is the collection of BS nodes,  $V^U$  is the collection of UAV nodes,  $V^L$  and  $V^G$  are collections of LEO nodes and GEO nodes. Each node in  $V$  contains its attribute parameters and network state parameters. When the node encounters buffer overflow or connection breakage during packet transmission, it will discard these packets. And the collection of edges  $E = \{e_{11}^{DB}, e_{12}^{DB}, \dots, e_{ij}^{xy}\}$  indicates the

link between any two nodes in the network. For example,  $e_{ij}^{xy}$  is the link between node  $v_i^x$  and  $v_j^y$ , among them  $v_i^x, v_j^y \in V$ . The model will dynamically change the value of  $E$  according to the location and connection state of each node in the network. Further,  $V^D$  is divided into two types: source node and target node. The source node is responsible for sending data packets to the target node according to the preset generation mode. Other type nodes in  $V$  are relay nodes used to transfer data packets from the source node to the target node. Among them, base station nodes  $V^B$  will make the offloading decision as the traffic offloading node.

### B. Transmission Model

Since SAGIN consists of several different types of nodes, the transmission model between the nodes is not the same. For the transmission between UAV and BS, according to [36], its path loss can be found using the following equation:

$$PL(l_{UB}, \omega) = 10\varphi \log(l_{UB}) + \eta(\omega - \omega_0) e^{\frac{\omega_0 - \omega}{\gamma}} + k_0 \quad (1)$$

where  $l_{UB}$  denotes the horizontal distance between UAV and BS,  $\omega$  is vertical angle between UAV and BS,  $\omega_0$  is the angle offset. And  $\varphi$ ,  $\eta$  and  $k_0$  are the terrestrial path loss exponent, excess path loss and excess path loss offset,  $\gamma$  is the angle scalar. Due to the dynamic nature of the UAV nodes, the distance between UAV and BS changes over time. We use the model in [37] to calculate the transmission rate between UAV and BS during time slot  $t$ , with the transmission rate as follows:

$$v_t^{UB} = B_{UB} \log_2 \left( 1 + \frac{P_{UB} \cdot 10^{-\frac{PL}{10}}}{\sigma_{UB}^2} \right) \quad (2)$$

where  $B_{UB}$  is the channel bandwidth between UAV and BS,  $P_{UB}$  represents the transmission power,  $PL$  represents the path loss, and  $\sigma_{UB}^2$  represents the power of background noise.

In our transmission model, the distance between the satellite and UAV is much greater than the movement range of UAV and UAV height. Thus, UAVs can be considered ground devices along with the base station during satellite transmission. We consider the rain attenuation during satellite-to-ground communication as a Weibull-based stochastic process [38]. Therefore, the channel gain from the UAV to the satellite can be calculated by the following equation:

$$I_{US} = \frac{G_{US} G_{SU} \lambda_{US}^2}{(4\pi l_{US})^2} 10^{-\frac{F_{rain}}{10}} \quad (3)$$

where  $G_{US}$  and  $G_{SU}$  are the antenna gains of the UAV and the satellite, respectively,  $\lambda_{US}$  is the wavelength, and  $l_{US}$  is the distance between the satellite and the UAV. Rain attenuation  $F_{rain}$  is modeled by Weibull distribution [38]. Similarly, the channel gain from the BS to the satellite can be calculated as follows:

$$I_{BS} = \frac{G_{BS} G_{SB} \lambda_{BS}^2}{(4\pi l_{BS})^2} 10^{-\frac{F_{rain}}{10}} \quad (4)$$

where  $G_{BS}$  and  $G_{SB}$  are the antenna gains of the BS and the satellite respectively,  $\lambda_{BS}$  is the wavelength, and  $l_{BS}$  is the distance between the BS and the satellite. After getting the

channel gain, we can calculate the transmission rate between the UAV and the satellite as follows:

$$\nu_t^{US} = B_{US} \log_2 \left( 1 + \frac{P_{US} \cdot |I_{US}|}{\sigma_{US}^2} \right) \quad (5)$$

where  $B_{US}$  represents the channel bandwidth of the UAV and the satellite communication link,  $P_{US}$  represents the transmission power between the UAV and the satellite,  $I_{US}$  represents the channel gain between the UAV and the satellite, and  $\sigma_{US}^2$  is the power of the background noise. Additionally, the transmission rate between the BS and the satellite can be obtained as follows:

$$\nu_t^{BS} = B_{BS} \log_2 \left( 1 + \frac{P_{BS} \cdot |I_{BS}|}{\sigma_{BS}^2} \right) \quad (6)$$

where  $B_{BS}$  denotes the channel bandwidth between the BS and satellite communication link,  $P_{BS}$  denotes the transmission power between the BS and satellite,  $I_{BS}$  is the channel gain between the BS and satellite, and  $\sigma_{BS}^2$  is the power of the background noise.

In addition, the transmission between GEO and LEO adopts the free space propagation model. We can calculate the channel gain of inter-satellite communication as follows:

$$I_{HL} = \frac{G_{HL} G_{LH} \lambda_{HL}^2}{(4\pi d_{HL})^2} \quad (7)$$

where  $G_{HL}$  and  $G_{LH}$  are the antenna gains of GEO and LEO,  $\lambda_{HL}$  is the wavelength, and  $d_{HL}$  is the distance between the GEO and LEO. The transmission rate of GEO and LEO is as follows:

$$\nu_t^{HL} = B_{HL} \log_2 \left( 1 + \frac{P_{HL} \cdot |I_{HL}|}{\sigma_{HL}^2} \right) \quad (8)$$

where  $B_{HL}$  is the channel bandwidth between the GEO and LEO communication link,  $P_{HL}$  is the transmission power between the GEO and LEO,  $I_{HL}$  denotes the channel gain between the GEO and LEO, and  $\sigma_{HL}^2$  denotes the power of the background noise.

### C. Problem Formulation

Denoted by  $X_t$ , the total amount of data packets in  $t$  th time slot.  $X_t$  follows a normal distribution and these packets are transmitted to node  $v_d \in V$  through the set routing path  $L_t = \{v_x, v_y, \dots, v_d\}$ . This process can be expressed as  $M_{v_s \rightarrow v_d} = \{X_t, L_t, T, r\}$ ,  $X_t$  is the amount of data,  $L_t$  is the defined routing path,  $T$  is the transmission delay of the packet, and  $r$  is the transmission rate of the packet. The system will generate a new routing path when using the traffic offloading algorithm  $A(m)$  during data packet transmission.  $m$  represents the process of transmitting this data packet. At this time, there will be UAV or satellite nodes in the path. In this paper, our goal is to minimize the packet delay of the entire network over the designed time  $D$  by the usage of the traffic offloading method to generate a new routing path at time  $t$  dependent on specific parameters of:

$$Z_\mu(m, t) = \begin{cases} m_{X_t, L_t, T, r}, & a = 0 \\ A(m_{X_t, L_t, T, r}), & a \neq 0 \end{cases} \quad (9)$$

where  $a = 0$  means the offloading method is not used, and  $a \neq 0$  means that the offloading method is used.

$$\min_\mu \sum_D \sum_m O(Z_\mu(m, t)) \quad (10)$$

where  $O(\cdot)$  calculates the delay of all packets in the entire network. We are committed to minimizing the delay of data packets by minimizing the parameters of the traffic offloading method.

## IV. TRAFFIC OFFLOADING METHOD

To address the above issues, we define SAGIN as a differentiated environment consisting of a series of dynamically heterogeneous local environments. Since the steady-state Markov process is not applicable to this environment, this study will describe the process using the framework of the DEC-POMDP [39], which is an extended form of the partially observable Markov decision process. ‘‘Distributed’’ implies that the training of each intelligence in the process is decentralized. ‘‘Partially observable’’ implies that each agent can only observe a part of the environment. Therefore, this study formalizes the local policy trap problem in SAGIN as an environment discretization problem under DEC-POMDP to obtain the joint optimal traffic offloading policy in such a differentiated environment. Among them, the BS located in each region are considered agents. During packet transmission, BS will use the traffic offload algorithm  $A$  to make traffic offload decisions. And, the traffic offloading problem in SAGIN is transformed into DEC-POMDP and thus trained to learn.

In DEC-POMDP, multiple agents need to work together to maximize the overall reward in a partially observable environment, where each agent can only partially observe the environment state, and cooperation among the agents is required to solve the task. DEC-POMDP represents the connections among the agents through an interrelationship model between the agents and takes action based on this. The DEC-POMDP model can be defined as:

$$M = \langle S, A_1, \dots, A_N, T, R, \Omega_1, \dots, \Omega_N, O, \gamma \rangle \quad (11)$$

among them:

- $S$  is the state space, representing all possible states of the environment. In the traffic offloading problem discussed in this study, this state space can be considered as the network state information of all network nodes which includes both dynamically changing data and covers static characteristics such as generalized mobility model, service generation rate, and queue maximum size. This is not complete for a single agent to make observations.
- $A_i = \{0, op_j\}$  denotes the action space of agent  $i$ . Here, the first element signifies that the system forwards the packet to the next node following the pre-set routing path. The second element, the offloading path  $op_j$ , involves the selection of the next relay node, such as UAV, LEO, or GEO, to determine the new transmission path for the packet.
- $T : S \times A_1 \times \dots \times A_N \times S \rightarrow [0, 1]$  is a state transfer function that represents the probability distribution of

transfer to a new state given the actions of all agents and the current state of the environment. Since it is difficult to build a corresponding mathematical model for the state transfer probability of the environment, a model-free approach will be used to represent it. i.e., the training is guided by neural network models to predict future reward sums.

- $R : S \times A_1 \times \dots \times A_N \rightarrow \mathbb{R}$  is reward function. Our optimization goal is to minimize the total delay of the network, so set the reward function to

$$R(s_t, a_t) = \begin{cases} \frac{1}{D_t}, & \text{arrive} \\ -(T_{drop} - T_{born}), & \text{drop} \end{cases} \quad (12)$$

where  $D_t$  denotes packet delay,  $T_{drop}$  and  $T_{born}$  are the packet drop time and born time. When the packet arrives at the destination node, a positive reward is given according to the delay, and the lower the delay, the higher the reward value. If the packet is dropped, a negative reward is given according to the difference between the drop and the birth time.

- $\Omega_i$  denotes the observation space of agent  $i$ , which is defined as a multi-dimensional array encompassing all relevant information of the one-hop and two-hop neighboring nodes in its vicinity. The observational outcomes of the agent are constituted by the network information of these one-hop and two-hop neighbors, meticulously gathered through the use of the Hello protocol. For any given time point  $t$ , the local observation space  $\Omega_i$  for the agent can be defined as:

$$\Omega_i(t) = \left\{ \mathcal{Q}_t^{adj(i)}, \mathcal{Q}_t^{adj(adj(i))} \right\} \quad (13)$$

where  $\mathcal{Q}_t^{adj(i)}$  represents the collection of relevant information at time  $t$  for the immediate neighbors of the agent  $i$ .  $\mathcal{Q}_t^{adj(adj(i))}$  represents the collection of relevant information at time  $t$  for the neighbors of the neighbors of agent  $i$ .  $adj(i)$  denotes the set of adjacent nodes to agent  $i$ .

- $O : S \times A_1 \times \dots \times A_N \times \Omega_1 \times \dots \times \Omega_N \rightarrow [0, 1]$  is an observation probability function that represents the probability distribution of its local observations observed by agent  $i$ , given the actions of all agents and the current state of the environment.
- $\gamma \in (0, 1)$  denotes the discount factor used to measure the importance of immediate and future rewards.

In DEC-POMDP, the goal of the agents is to maximize the joint reward of all agents, i.e., the total expected discounted reward. So, we can define the policy  $\pi_i$  of agent  $i$  as a mapping from the observation space  $\omega_i$  to the action space  $A_i$ , indicating which action agent  $i$  should take given its local observations. The joint policy  $\Pi$  is then defined as the combination of all the agents' policies. Thus, this paper aims to solve the joint traffic offloading policy  $P_i$  for each agent in SAGIN.

### A. Differentiated Federated Reinforcement Learning

In recent years, many scholars have proposed various algorithms to solve the joint policy in DEC-POMDP, including

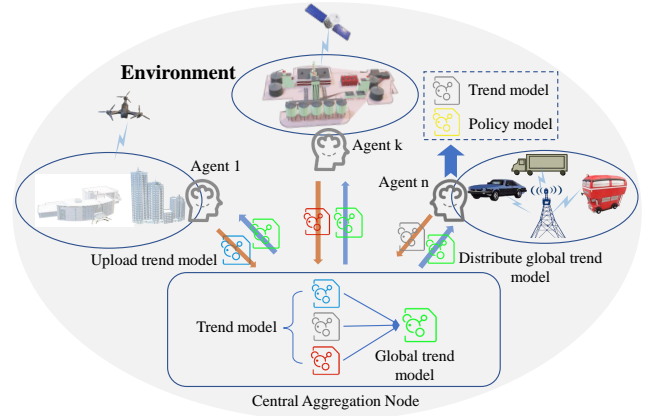


Fig. 3. The agents communicate through the trend model.

a value-based iterative approach, policy gradient-based approach, Monte Carlo tree search-based approach, and so on. In this paper, the concept of DFRL is proposed to solve this joint policy  $\Pi$  as shown in Fig. 3.

The main improvement between DFRL and FRL is that DFRL isolates local learning from the global learning process. We first divide the local learning model of the agent into two parts: the trend model and the policy model. The policy model generates policies based on the local environment state, which is isolated and never shared with other agents. The trend model is the window to communicate with the center node during the global learning process and guide the update of the policy model. And The global trend model is generated by the local trend model and shared among all agents, which is responsible for describing the state of the global environment, changing trends, and transferring the knowledge learned by all the agents.

The introduction of the trend model separates the learning of policy models, thus allowing agents to cooperate in obtaining policy models that satisfy their preferences. More importantly, the training process of the policy models of each agent is still distributed in the training nodes themselves, thus alleviating the environmental variability problem while satisfying the DEC-POMDP requirements, i.e., the non-IID data and local strategy trap problem in dynamically heterogeneous networks. Otherwise, separating the policy models ensures the privacy security of the policy models to meet the requirements of FL.

Theoretically, the DFRL architecture can be applied to many RL algorithms, such as Deep-Q network (DQN) [40] and Actor-critic methods [41]. In the following section, we describe the concept of DFRL in the context of specific algorithms.

### B. Algorithms: Differentiated Federated Soft Actor-Critic

In this section, we propose a novel DFRL algorithm called Differentiated Federated Soft Actor-Critic (DFSAC), which introduces the concept of DFRL based on the Soft Actor-Critic (SAC) algorithm [42]. Compared with the traditional FRL algorithm, this algorithm can better adapt to the differentiation between the environments in which the agents are located.

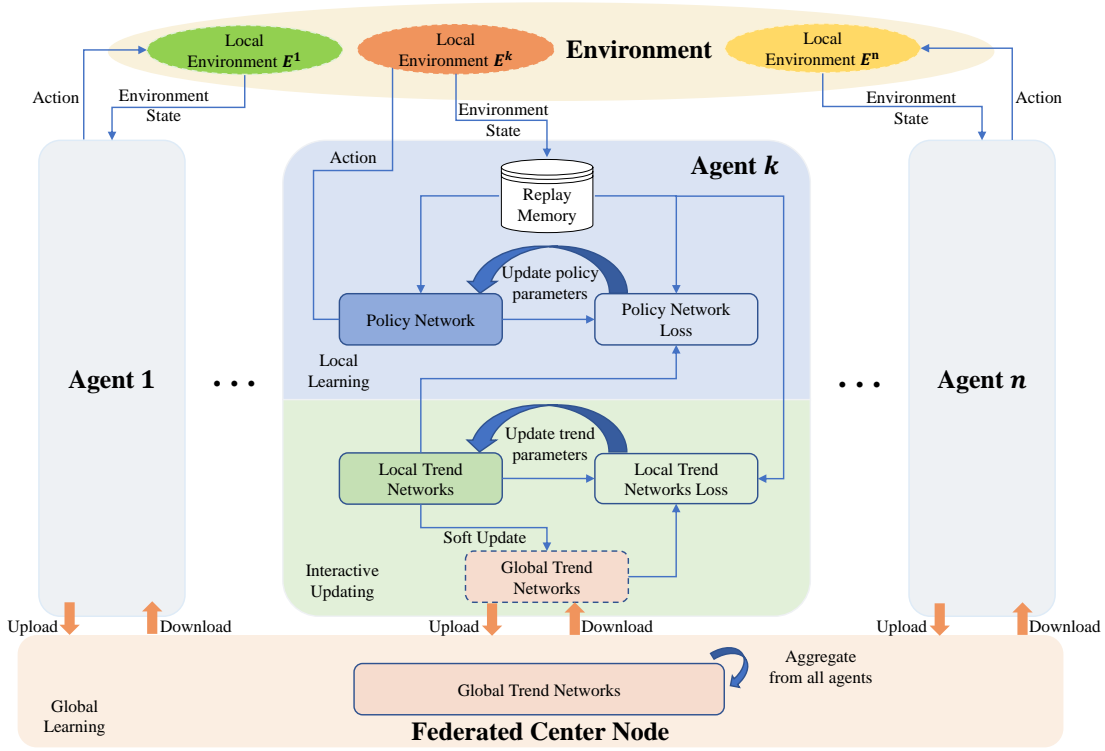


Fig. 4. The learning process of DFSAC algorithm.

Due to data heterogeneity and node dynamics, the traditional FRL algorithm is not capable to handle the cooperated learning tasks in differentiated environments. Therefore, we propose the DFSAC, which evolves SAC based on the concept of DFRL. The SAC algorithm is a DRL method consisting of Actor and Critic networks to optimize random policies in an off-policy manner. The Actor network is used to output the actions, and the Critic network is used to evaluate the states and actions. Its core feature is entropy regularization, where policy training trades off maximizing expected reward and entropy. Increasing entropy makes the policy explore more, speeding up the subsequent learning process and preventing the policy from prematurely converging to a local optimum. This makes it well-suited for exploring optimal strategies in several differentiated environments. Therefore, we choose this algorithm as our infrastructure. In the DFRL framework, the Actor network is employed as the policy model, which is responsible for determining the actions to be taken based on the current state. On the other hand, the Critic network serves as the trend model, providing feedback on the actions chosen by the Actor network by estimating the value function. Within the Critic network, the target-critic component is specifically utilized as the global trend model, which plays a crucial role in the aggregation process by integrating local and global trend information to enhance the decision-making strategy in DFRL. The DFSAC algorithm is processed by multiple agents and a federated center node, where the agents generate local trend networks and policy networks through the local environment, and the federated center node is responsible for collecting the local trend networks of each agent and aggregating them into global trend networks then distributing them to the agents.

Fig. 4 shows the whole learning process including global learning, interactive updating, and local learning, and the global trend model participates in all the processes in any agent  $k$ . We define a global environment  $E = \{E^1, \dots, E^k, \dots, E^n\}$  consisting of a set of  $N$  differentiated local environments. A corresponding agent exists for each local environment, i.e. environment  $E^k$  corresponds to agent  $k$ . Unlike the traditional FRL approach, the goal of the DFSAC algorithm is to obtain a set of policies  $\Pi = \{\tilde{\pi}^1, \dots, \tilde{\pi}^k, \dots, \tilde{\pi}^n\}$ , and the target policy  $\tilde{\pi}^i$  in each local environment  $E^i$  is:

$$\tilde{\pi}^i = \operatorname{argmax}_{\pi^i} \sum_{t=0}^T E_{(s_t^i, a_t^i) \sim \tau_{\pi^i}} [\gamma r(s_t^i, a_t^i) + \alpha^i \mathcal{H}(\pi^i(\cdot | s_t^i))] \quad (14)$$

where  $\tilde{\pi}^i$  is the target policy,  $\pi^i$  is a policy of agent  $i$  in local environment  $E^i$ ,  $\gamma$  is the discount rate and  $r$  is a reward from the environment,  $s_t^i \in S$  and  $a_t^i \in A$  denote state and action in local environment  $E^i$  with timestamp  $t$ , and  $S, A$  are global state space and global action space, respectively.  $\tau_{\pi^i}$  is the distribution of trajectories generated from policy  $\pi^i$ ,  $\alpha^i$  is the temperature parameter to control the positivity of the policy exploration in the local environment and  $\mathcal{H}(\cdot)$  indicates the entropy. And, in differentiated environments, the transfer probability varies among local environments, i.e.  $P(s_{t+1}^i | s_t^i, a) \neq P(s_{t+1}^j | s_t^j, a), i \neq j$ .

To obtain  $\Pi$ , the DFSAC algorithm uses an approach like the soft policy iteration [43], that alternately performs the two steps of policy evaluation and policy improvement to converge to the optimal value function and optimal policy. An example

is that in agent  $k$ , the policy network takes the state  $s_t^k$  of the local environment  $E^k$  at moment  $t$  as input to obtain the action  $a_t^k$ . To evaluate the impact of the policy on the local environment as well as on the global environment, the soft state value is defined as:

$$V(s_t^k) := \pi^k(s_t^k)^T [\chi(s_t^k) - \alpha^k \log(\pi^k(s_t^k))] \quad (15)$$

among them,  $\chi$  denote the trend networks. And, we use the following loss function to update trend networks:

$$J_{\chi^k}(\theta) = \mathbb{E}_{(s_t^k, a_t^k) \sim \mathcal{D}} \left[ \frac{1}{2} \left( \chi_{\theta}^k(s_t^k, a_t^k) - \left( r(s_t^k, a_t^k) + \gamma \mathbb{E}_{s_{t+1}^k \sim p} \left[ V_{\bar{\theta}}(s_{t+1}^k) \right] \right) \right)^2 \right] \quad (16)$$

where  $\mathcal{D}$  denotes a trajectory stored in replay memory and  $V_{\bar{\theta}}$  denotes the use of global trend networks to calculate the soft state value. This can be seen as a collaborative learning process using the global information shared by the collaborators. And the global trend networks are obtained by a soft update of each agent and then aggregated:

$$\bar{\chi} \leftarrow \varepsilon \chi^k + (1 - \varepsilon) \bar{\chi}, k \in \{1, 2, \dots, N\} \quad (17)$$

where  $\bar{\chi}$  is global trend networks,  $\chi^k$  is local trend networks in agent  $k$ ,  $\varepsilon$  is the aggregation factor. Then, the updated trend networks are used to guide the policy improvement:

$$J_{\pi^k}(\phi) = E_{s_t^k \sim \mathcal{D}} [\pi^k(s_t^k)^T [\alpha^k \log(\pi_{\phi}^k(s_t)) - \chi_{\theta}^k(s_t^k)]] \quad (18)$$

And  $\alpha$  is the temperature parameter, but the appropriate value of  $\alpha$  is different at different stages of training. Therefore, the selection of the  $\alpha$  is formulated as a constrained optimization problem [44] that maximizes the expected return while keeping the entropy of the policy greater than a threshold as follows:

$$\max_{\pi_0, \dots, \pi_T} \mathbb{E} \left[ \sum_{t=0}^T r(s_t, a_t) \right] \quad \text{s.t. } \forall t, \mathcal{H}(\pi_t) \geq \mathcal{H}_0 \quad (19)$$

we also dynamically train with  $\alpha$  as a parameter of model, the loss function of  $\alpha^k$  in agent  $k$  is

$$J(\alpha^k) = \pi_t^k(s_t^k)^T [-\alpha^k (\log(\pi_t^k(s_t^k)) + \bar{H})] \quad (20)$$

Through the continuous iteration of the above two steps, agent  $K$  is guided by global trend networks to obtain the target policy  $\tilde{\pi}^k$  applicable to the local environment  $E^k$  by sharing knowledge with other agents eventually.

The complete DFSAC algorithm we proposed is shown in Algorithm 1. Firstly, the algorithm initializes the local network parameters of local networks (two local trend networks and one policy network). Then initialize the global network (two global trend networks), and equalize the global trend networks and local trend networks' parameters. The agents first initialize an empty replay memory and store a backup of the global trend networks. And the agent will collect information about the local environment as the input of the policy network and output an action. Then, the agent gets the reward value and the next state from the local environment calculates the cumulative discount function at each moment and stores the transition in the replay memory. When the number of transitions in the replay memory reaches the number set in advance, the

algorithm starts training. It updates the local network parameters. Finally, the algorithm updates the temperature parameters. After running preset training iterations, the agents will upload their local trend network parameters to the federated center node. After receiving the local trend network parameters of agents, the federated center node integrates these parameters and updates the global trend network parameters. Afterward, the federated center node distributes the global trend network parameters to each agent, and agents update their backup with the latest global trend networks. Through this integrated and distributed processing method, the local network in the agents can reflect the real-time situation of the global environment.

---

#### Algorithm 1 DFSAC algorithm

---

- 1: Initialize  $\chi_{\theta_1}^n : S \rightarrow \mathbb{R}^{|A|}, \chi_{\theta_2}^n : S \rightarrow \mathbb{R}^{|A|}, \pi_{\phi}^n : S \rightarrow [0, 1]^{|A|}$  for  $n \in \{1, 2, \dots, N\}$   $\triangleright$  Initialize local network parameters
  - 2: Initialize  $\bar{\chi}_{\theta_1} : S \rightarrow \mathbb{R}^{|A|}, \bar{\chi}_{\theta_2} : S \rightarrow \mathbb{R}^{|A|}$   $\triangleright$  Initialize global trend networks parameters at the federated center node
  - 3:  $\theta_1^n \leftarrow \bar{\theta}_1, \theta_2^n \leftarrow \bar{\theta}_2$  for  $n \in \{1, 2, \dots, N\}$   $\triangleright$  Equalize global trend networks and local trend network parameters
  - 4:  $D^n \leftarrow \emptyset$  for  $n \in \{1, 2, \dots, N\}$   $\triangleright$  Initialize an empty replay memory
  - 5: **while** running **do**
  - 6:     **for** each agent  $n$  **do**
  - 7:         Get state  $s_t$  from the environment  $E^n$
  - 8:          $a_t \sim \pi(a_t | s_t)$   $\triangleright$  Sample action from the agent  $n$
  - 9:          $s_{t+1} \sim p(s_{t+1} | s_t, a_t)$   $\triangleright$  Sample transition from the environment  $E^n$
  - 10:          $D^n \leftarrow D^n \cup \{(s_t, a_t, r(s_t, a_t), s_{t+1})\}$   $\triangleright$  Store the transition in replay memory
  - 11:          $\theta_i^n \leftarrow \theta_i^n - \lambda_{\chi} \hat{\nabla}_{\theta_i^n} J_{\chi}(\theta_i^n)$  for  $i \in \{1, 2\}$  and  $n \in \{1, 2, \dots, N\}$   $\triangleright$  Update local trend networks use Eq (16)
  - 12:          $\phi^n \leftarrow \phi^n - \lambda_{\pi} \hat{\nabla}_{\phi^n} J_{\pi}(\phi^n)$  for  $n \in \{1, 2, \dots, N\}$   $\triangleright$  Update policy networks use Eq (18)
  - 13:          $\alpha^n \leftarrow \alpha^n - \lambda \hat{\nabla}_{\alpha^n} J(\alpha^n)$   $\triangleright$  Update temperature use Eq (20)
  - 14:         When running  $k$  iterations upload  $\theta_1^n, \theta_2^n$  to federated center node
  - 15:     **end for**
  - 16:     **if** in federated center node **then**
  - 17:          $\bar{\chi}_i \leftarrow \varepsilon \chi_i^n + (1 - \varepsilon) \bar{\chi}_i$  for  $i \in \{1, 2\}$  and  $n \in \{1, 2, \dots, N\}$   $\triangleright$  Aggregated global trend network use Eq. (17)
  - 18:     **end if**
  - 19: **end while**
- 

## V. EMPIRICAL STUDY

### A. Settings

In this section, we describe the simulation environment and evaluate our DFSAC-based traffic offloading method. The simulation environment is a dynamic four-layer SAGIN environment, which divides each device into three types: the source node device type, the relay node device type, and the destination node device type.



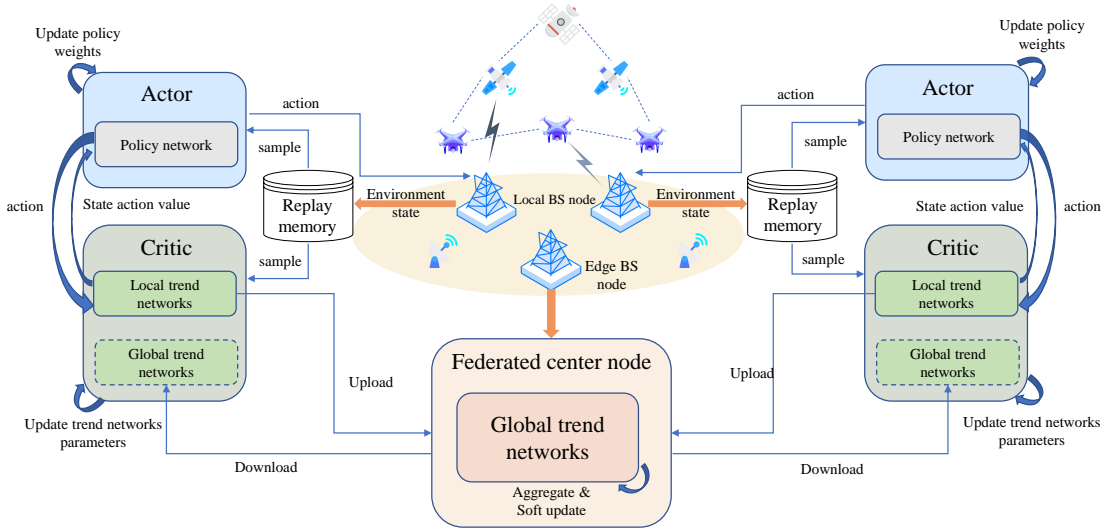


Fig. 5. The frame structure of DFSAC-based traffic offloading method in SAGIN.

The source node device is used to generate data packets, which are transmitted to the destination node device for reception. The source node device and the destination node device are modeled as user equipment (UE). The source node device generates data using a normally distributed data generation rate. The relay node device transfers data packets based on the Open Shortest Path First (OSPF) traffic routing protocol or the proposed traffic offloading method. The relay nodes are modeled as BS, UAV, LEO, and GEO. Except for BSs are static, other devices have their mobile models. The source node and destination node devices use a random-waypoint mobility model [45], which can simply randomize the device's location.

In our simulation, the UEs use the arbitrary movement model to move within a  $100km^2$  two-dimensional plane area. For the mobility model of UAV, we set the UAV to fly in a certain direction at a fixed speed in a predefined period. After a while, the UAV's moving direction will shift arbitrarily and fly at the same moving speed. We consider that LEO periodically covers the considered area as it moves around the earth. For GEO, due to its high coverage, we consider it to cover the target area all the time and have a stable connection (i.e., the air condition change is not considered). For this area,  $(\varphi, \omega_0, \eta, \gamma, k_0)$  is set to  $(3.04, -3.61, -23.29, 4.14, 20.7)$  [46]. For the rain attenuation of the satellite channel model, we set the  $F_{rain}$  to 6dB [46]. The remaining simulation parameters are shown in Table I.

The frame structure of the DFSAC-based traffic offloading method is shown in Fig. 5. The local base station nodes are training nodes for federated learning to collect the network environment state in replay memory for training. The neural network in each training node can be divided into Actor-network and Critic-network. The Actor-network is responsible for outputting the action of traffic offloading based on the policy network. The Critic network is responsible for judging the actions and the environmental trends, which are used to optimize and improve the performance of the policy network. Meanwhile, the edge base station node is used as a federated

TABLE I  
SIMULATION PARAMETERS

Parameters	Values
Number of BS nodes	8
Number of UAV nodes	6
Number of LEO nodes	2
Number of GEO nodes	1
Packet generation rate	$N(1e6, \sigma^2)$ Mbps
BS - UAV bandwidth	20MHz
BS - LEO bandwidth	37.5MHz
BS - GEO bandwidth	25MHz
UAV - LEO bandwidth	10MHz
UAV - GEO bandwidth	5MHz
Gamma	0.99
$\epsilon$	1e-2
Learning rate	5e-4
Learning rate of $\alpha$	1e-3
Target entropy	-4.0

center node to receive local trend networks uploaded by the local base station nodes, aggregate the parameters, and then distribute them to the local base station nodes to update their network parameters to the latest.

### B. Performance Evaluation

To evaluate our proposed traffic offloading method, we conducted experiments comparing it with existing approaches. Specifically, we run the traditional traffic offloading method used in [47]–[50] in the same environment. The traditional traffic offloading method is a greedy offloading method based on the shortest path. In addition, we simulated the distributed DDQN-based traffic offloading method in [8] and the traditional FRL method for experimental comparison.

Fig. 6 and Fig. 7 present a comprehensive analysis of throughput and packet loss rate variations as the number of source nodes incrementally increases from 200 to 380, involving the evaluation of four distinct methods. Fig. 6 clearly demonstrates the substantial advantages of our proposed method. As the number of source nodes increases,

our method exhibits a consistent improvement in throughput, contrasting with the unaltered performance of the other methodologies. This compelling observation underscores our method's remarkable ability to enhance network throughput and maximize its capacity compared to its counterparts. Fig. 7 shows the dynamics of packet loss rates as the source node count escalates. It becomes evident that, with the increasing number of source nodes, the packet loss rates of all methods tend to rise, a common trend attributed to network congestion. However, our proposed method emerges as the standout performer, consistently maintaining a lower packet loss rate in comparison to the other three methods. This noteworthy outcome underscores the efficacy of our approach in mitigating packet loss and ensuring reliable data transmission in dynamic and challenging network environments.

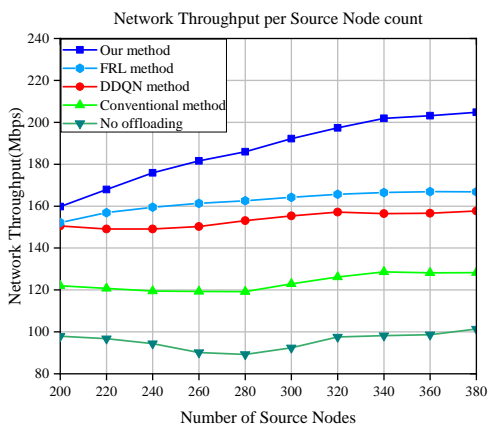


Fig. 6. Throughput per increasing source node count.

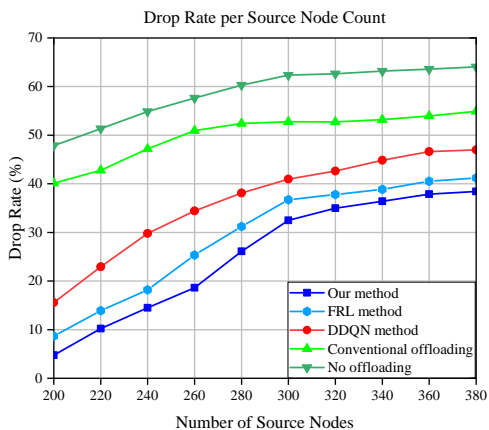


Fig. 7. Drop rate per increasing source node count.

Next, we focus on evaluating packet delay performance, as depicted in Fig. 8, when these methods are deployed while varying the number of source nodes. A clear trend emerges from the figure: when the offloading method is not utilized, there is a noticeable and consistent increase in average packet delay as the number of source nodes grows. The strategy proposed in this study significantly outperforms the methods based on DDQN and FRL in scenarios with source node counts ranging from 200 to 380. This empirical evidence highlights the remarkable capability of our method to maintain

the lowest delay, even amid increasing network loads and a significant rise in the number of nodes. Notably, such performance demonstrates the robust stability of our approach in maintaining low latency within dynamic heterogeneous network environments.

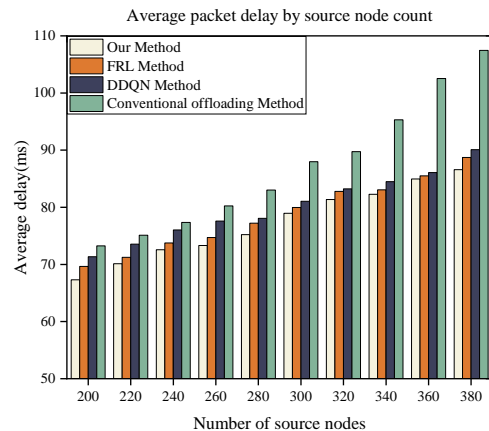


Fig. 8. Packet delay per increasing source node count.

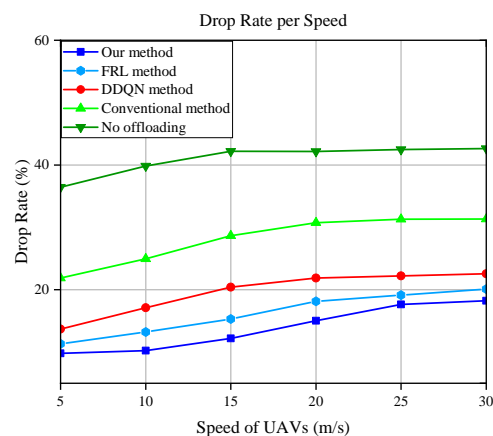


Fig. 9. Packet drop rate with increasing moving speed of UAVs.

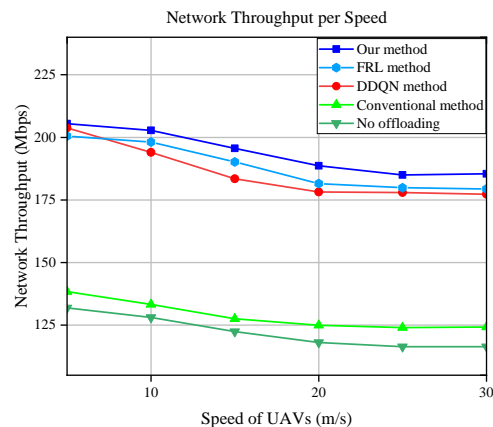


Fig. 10. Throughput with increasing moving speed of UAVs.

Further, to assess the robustness of our algorithm in the high-mobility setting of SAGIN, we conducted simulations

to evaluate packet drop rates and throughput. The movement characteristics of the satellites are predetermined within the environment. Therefore, we systematically varied the speeds of the UAVs, ranging from 5m/s to 30m/s, to assess the performance of each method. Fig. 9 and Fig. 10 show the packet drop rate increases, and throughput decreases with UAV moving speed increase. Nevertheless, our proposed method consistently outperforms the alternatives even as the UAV speed increases. This trend aligns with the earlier experimental results, affirming the overall superiority of our proposed traffic offloading method in the context of SAGIN, particularly in high-mobility scenarios.

### C. Validity of DFRL Framework

The distributed DDQN-based traffic offloading approach [8] is at the forefront of achieving optimal performance in highly dynamic SAGIN environments. Building upon this foundation, we undertake a comprehensive evaluation to validate the effectiveness of our proposed DFRL framework thoroughly. We meticulously compare and contrast the performance of three distinct approaches: the distributed DDQN-based traffic offloading approach [8], the FL-based DDQN approach, and the DFRL-based DDQN approach.

Fig. 11, Fig. 12 and Fig. 13 show throughput, packet loss rate, and packet delay fluctuation as the number of source nodes increases gradually from 200 to 380. The graphical representation clearly and decisively showcases the superior performance of the FL-based DDQN approach compared to the approach detailed in [8]. This advantage can be attributed to the heightened interaction among intelligent agents, which empowers them to make more effective decisions by harnessing global state information, resulting in consistently superior outcomes. Furthermore, as the number of source nodes continuously increases, leading to a corresponding rise in network load, the DFRL-based DDQN approach consistently outperforms the other two algorithms, exhibiting exceptional performance across three key performance indicators. This is attributed to our proposed DFRL framework, where agents develop an understanding of the uniqueness of each region based on information collected from multiple areas, thereby formulating more suitable local strategies for their respective regions. This serves as a demonstration of the robustness and effectiveness of our proposed DFRL framework, providing a viable solution to address the intricate challenges posed by the dynamic and heterogeneous SAGIN environment.

### D. Application Cases

We use CartPole from OpenAI Gym [51] as the experimental environment. CartPole is a cart-pole game with a cart and a pole erected. The algorithm must control the cart to move left or right to keep the pole upright while satisfying the constraints. Whenever an operation is performed without tilting the rod beyond the limit, the environment rewards the agent with a value of 1, it otherwise 0. Our goal is to verify the algorithm's performance in several differentiated environments. Therefore, the transfer probability of the environment is varied by changing the length of the pole to obtain multiple

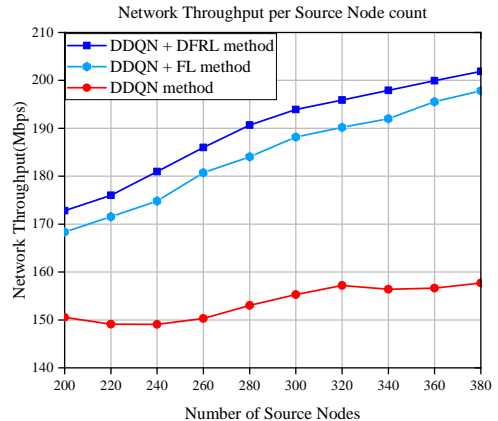


Fig. 11. Throughput per increasing source node count.

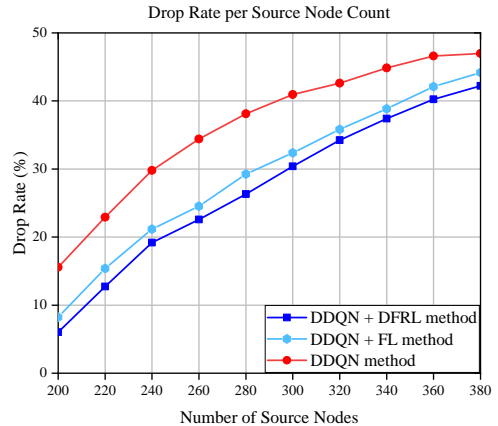


Fig. 12. Drop rate per increasing source node count.

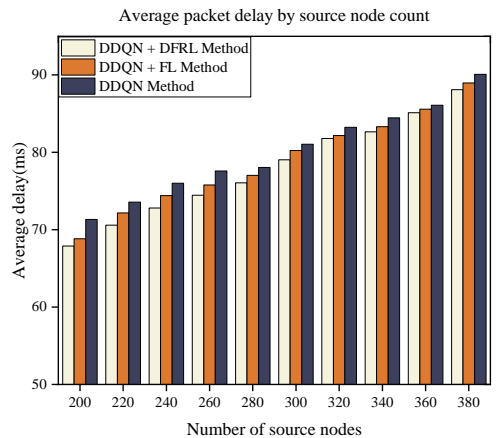


Fig. 13. Packet delay per increasing source node count.

differentiated environments and setting up agents in each of the differentiated environments.

For the DFSAC algorithm, we use the output of the policy network to control the agents. An additional node is set as a federated center node for model aggregation and distribution. As described in the previous section, agents in different environments share knowledge by sharing global trend networks. Considering the communication overhead between agents, each agent stores a backup copy of global trend networks locally to guide the local data training in the actual algorithm implementation. After a certain number of training sessions, i.e., soft update, the local trend networks are uploaded to the federal center node for aggregation, and the latest global trend networks are downloaded to update the local backup. This approach can reduce the communication overhead incurred during training.

The characteristic between the DFSAC algorithm and the traditional FRL algorithm is that the differentiation between different environments is considered in the learning process. Therefore, we use the SAC algorithm trained by traditional federated learning as our baseline, i.e. FedAVG [12], the agents in each environment share the model, including both value parameters and policy model parameters in the SAC algorithm, and the models are aggregated and distributed periodically using averaging model parameters. In addition, we implement a centralized RL algorithm to verify the performance of a non-federated learning algorithm in this scenario. The centralized RL algorithm uses a SAC agent to collect information directly in each environment and then train.

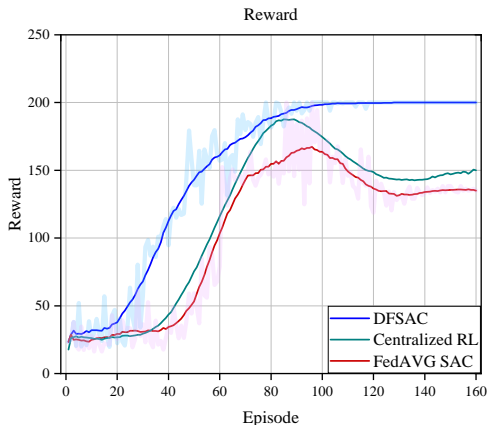


Fig. 14. The rewards achieved by different algorithms.

The reward obtained by the different algorithms in each episode is shown in Fig. 14, which results from averaging the reward from several differentiated environments. As can be seen, the DFSAC algorithm steadily increases and eventually converges as the training progresses. In contrast, the reward of other algorithms fluctuates significantly during the training process and is lower than that of the DFSAC algorithm. Fig. 15 shows the change in the average loss value of each agent's policy network during the training process. It can be seen from the figure that the DFSAC algorithm has higher learning efficiency than others. This indicates that the DFSAC algorithm can obtain policy models more suitable for agents in

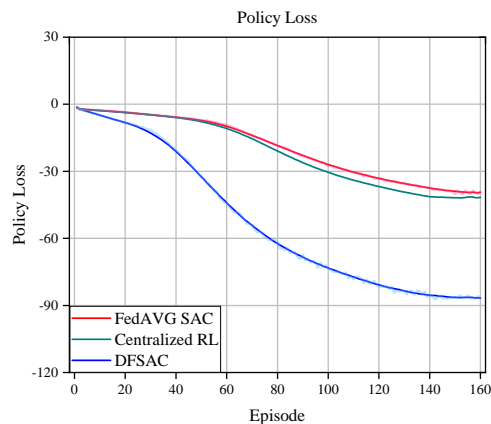


Fig. 15. The policy network loss evolution.

differentiated environments than the traditional FRL algorithm. The global policy model obtained by directly aggregating the traditional FRL algorithm is difficult to adapt to this differentiated environment. Further, the DFSAC algorithm can avoid the additional communication overhead incurred by the centralized RL algorithm when transmitting environmental information.

## VI. CONCLUSION

The dynamic and heterogeneous nature of the environment gives rise to distinct regions, creating a potential challenge for RL due to the risk of getting trapped in local strategies. This study introduces a novel concept called DFRL tailored for networks operating in dynamic and heterogeneous environments. DFRL isolates the local policy model from global integration and employs a trend model to discern regional variations. In contrast to conventional FRL, our proposal strongly emphasizes adjusting biases in differentiated regions and ensuring the independence of local policy models. Our approach demonstrates superior performance, particularly in scenarios characterized by heterogeneous network environments.

## REFERENCES

- [1] X. Xiong, K. Zheng, L. Lei, and L. Hou, "Resource allocation based on deep reinforcement learning in iot edge computing," *IEEE Journal on Selected Areas in Communications*, vol. 38, no. 6, pp. 1133–1146, 2020.
- [2] J. Wang, L. Zhao, J. Liu, and N. Kato, "Smart resource allocation for mobile edge computing: A deep reinforcement learning approach," *IEEE Transactions on Emerging Topics in Computing*, vol. 9, no. 3, pp. 1529–1541, 2021.
- [3] T. Dong, Z. Zhuang, Q. Qi, J. Wang, H. Sun, F. R. Yu, T. Sun, C. Zhou, and J. Liao, "Intelligent joint network slicing and routing via gcn-powered multi-task deep reinforcement learning," *IEEE Transactions on Cognitive Communications and Networking*, 2021.
- [4] F. Tang, Y. Kawamoto, N. Kato, K. Yano, and Y. Suzuki, "Probe delay based adaptive port scanning for iot devices with private ip address behind nat," *IEEE Network*, vol. 34, no. 2, pp. 195–201, 2019.
- [5] H. Guo, J. Liu, J. Zhang, W. Sun, and N. Kato, "Mobile-edge computation offloading for ultradense iot networks," *IEEE Internet of Things Journal*, vol. 5, no. 6, pp. 4977–4988, 2018.
- [6] F. Tang, B. Mao, N. Kato, and G. Gui, "Comprehensive survey on machine learning in vehicular network: technology, applications and challenges," *IEEE Communications Surveys & Tutorials*, 2021.

- [7] J. Zhang, H. Guo, J. Liu, and Y. Zhang, "Task offloading in vehicular edge computing networks: A load-balancing solution," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 2, pp. 2092–2104, 2020.
- [8] F. Tang, H. Hofner, N. Kato, K. Kaneko, Y. Yamashita, and M. Hangai, "A deep reinforcement learning-based dynamic traffic offloading in space-air-ground integrated networks (sagin)," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 1, pp. 276–289, 2021.
- [9] F. Tang, B. Mao, Y. Kawamoto, and N. Kato, "Survey on machine learning for intelligent end-to-end communication toward 6g: From network access, routing to traffic control and streaming adaption," *IEEE Communications Surveys & Tutorials*, vol. 23, no. 3, pp. 1578–1598, 2021.
- [10] J. Liu, Y. Shi, Z. M. Fadlullah, and N. Kato, "Space-air-ground integrated network: A survey," *IEEE Communications Surveys Tutorials*, vol. 20, no. 4, pp. 2714–2741, 2018.
- [11] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon, "Federated learning: Strategies for improving communication efficiency," *arXiv preprint arXiv:1610.05492*, 2016.
- [12] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial intelligence and statistics*. PMLR, 2017, pp. 1273–1282.
- [13] Z. Yuan, Z. Guo, Y. Xu, Y. Ying, and T. Yang, "Federated deep auc maximization for heterogeneous data with a constant communication complexity," in *International Conference on Machine Learning*. PMLR, 2021, pp. 12 219–12 229.
- [14] C. T. Dinh, N. Tran, and J. Nguyen, "Personalized federated learning with moreau envelopes," *Advances in Neural Information Processing Systems*, vol. 33, pp. 21 394–21 405, 2020.
- [15] Y. Ruan, X. Zhang, S.-C. Liang, and C. Joe-Wong, "Towards flexible device participation in federated learning," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2021, pp. 3403–3411.
- [16] Z. Qu, K. Lin, J. Kalagnanam, Z. Li, J. Zhou, and Z. Zhou, "Federated learning's blessing: Fedavg has linear speedup," *arXiv preprint arXiv:2007.05690*, 2020.
- [17] L. P. Kaelbling, M. L. Littman, and A. W. Moore, "Reinforcement learning: A survey," *Journal of artificial intelligence research*, vol. 4, pp. 237–285, 1996.
- [18] Y. Wang, M. Chen, T. Luo, W. Saad, D. Niyato, H. V. Poor, and S. Cui, "Performance optimization for semantic communications: An attention-based reinforcement learning approach," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 9, pp. 2598–2613, 2022.
- [19] Z. Luan, L. Lu, Q. Li, and Y. Jiang, "Epc-te: Explicit path control in traffic engineering with deep reinforcement learning," in *2021 IEEE Global Communications Conference (GLOBECOM)*. IEEE, 2021, pp. 1–6.
- [20] X. Gao, Y. Sun, H. Chen, X. Xu, and S. Cui, "Joint computing, pushing, and caching optimization for mobile edge computing networks via soft actor-critic learning," *IEEE Internet of Things Journal*, pp. 1–1, 2023.
- [21] M. Grounds and D. Kudenko, "Parallel reinforcement learning with linear function approximation," in *Adaptive Agents and Multi-Agent Systems III. Adaptation and Multi-Agent Learning*. Springer, 2005, pp. 60–74.
- [22] J. Qi, Q. Zhou, L. Lei, and K. Zheng, "Federated reinforcement learning: techniques, applications, and open challenges," *arXiv preprint arXiv:2108.11887*, 2021.
- [23] F. Hanzely and P. Richtárik, "Federated learning of a mixture of global and local models," *arXiv preprint arXiv:2002.05516*, 2020.
- [24] J. Shen, N. Cheng, X. Wang, F. Lyu, W. Xu, Z. Liu, K. Aldubaikhy, and X. Shen, "Ringsfl: An adaptive split federated learning towards taming client heterogeneity," *IEEE Transactions on Mobile Computing*, pp. 1–16, 2023.
- [25] X. Zhang, Y. Li, W. Li, K. Guo, and Y. Shao, "Personalized federated learning via variational bayesian inference," in *International Conference on Machine Learning*. PMLR, 2022, pp. 26 293–26 310.
- [26] A. Shamsian, A. Navon, E. Fetaya, and G. Chechik, "Personalized federated learning using hypernetworks," in *International Conference on Machine Learning*. PMLR, 2021, pp. 9489–9502.
- [27] Y. Deng, M. M. Kamani, and M. Mahdavi, "Adaptive personalized federated learning," *arXiv preprint arXiv:2003.13461*, 2020.
- [28] Y.-J. Liu, G. Feng, Y. Sun, S. Qin, and Y.-C. Liang, "Device association for ran slicing based on hybrid federated deep reinforcement learning," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 12, pp. 15 731–15 745, 2020.
- [29] D. Kwon, J. Jeon, S. Park, J. Kim, and S. Cho, "Multiagent ddpq-based deep learning for smart ocean federated learning iot networks," *IEEE Internet of Things Journal*, vol. 7, no. 10, pp. 9895–9903, 2020.
- [30] Y. Zhu, J. Xu, Y. Xie, J. Jiang, X. Yang, and Z. Li, "Dynamic task offloading in power grid internet of things: A fast-convergent federated learning approach," in *2021 IEEE 6th International Conference on Computer and Communication Systems (ICCCS)*, 2021, pp. 933–937.
- [31] S. Wang, S. Xu, J. Wang, P. Yu, Y. Wang, and F. Zhou, "Frl-assisted edge service offloading mechanism for iot applications in fiwi hetnets," in *NOMS 2023-2023 IEEE/IFIP Network Operations and Management Symposium*, 2023, pp. 1–5.
- [32] H. Guo and J. Liu, "Uav-enhanced intelligent offloading for internet of things at the edge," *IEEE Transactions on Industrial Informatics*, vol. 16, no. 4, pp. 2737–2746, 2020.
- [33] J. Li, K. Xue, D. S. L. Wei, J. Liu, and Y. Zhang, "Energy efficiency and traffic offloading optimization in integrated satellite/terrestrial radio access networks," *IEEE Transactions on Wireless Communications*, vol. 19, no. 4, pp. 2367–2381, 2020.
- [34] N. Cheng, F. Lyu, W. Quan, C. Zhou, H. He, W. Shi, and X. Shen, "Space/aerial-assisted computing offloading for iot applications: A learning-based approach," *IEEE Journal on Selected Areas in Communications*, vol. 37, no. 5, pp. 1117–1129, 2019.
- [35] S. Zhang, A. Liu, C. Han, X. Liang, X. Xu, and G. Wang, "Multiagent reinforcement learning-based orbital edge offloading in sagin supporting internet of remote things," *IEEE Internet of Things Journal*, vol. 10, no. 23, pp. 20 472–20 483, 2023.
- [36] W. Shi, J. Li, N. Cheng, F. Lyu, S. Zhang, H. Zhou, and X. Shen, "Multi-drone 3-d trajectory planning and scheduling in drone-assisted radio access networks," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 8, pp. 8145–8158, 2019.
- [37] X. Lyu, W. Ni, H. Tian, R. P. Liu, X. Wang, G. B. Giannakis, and A. Paulraj, "Optimal schedule of mobile edge computing for internet of things using partial information," *IEEE Journal on Selected Areas in Communications*, vol. 35, no. 11, pp. 2606–2615, 2017.
- [38] S. A. Kanellopoulos, C. I. Kourgiorgas, A. D. Panagopoulos, S. N. Livieratos, and G. E. Chatzarakis, "Channel model for satellite communication links above 10ghz based on weibull distribution," *IEEE communications letters*, vol. 18, no. 4, pp. 568–571, 2014.
- [39] D. S. Bernstein, R. Givan, N. Immerman, and S. Zilberstein, "The complexity of decentralized control of markov decision processes," *Mathematics of operations research*, vol. 27, no. 4, pp. 819–840, 2002.
- [40] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing atari with deep reinforcement learning," *arXiv preprint arXiv:1312.5602*, 2013.
- [41] V. Konda and J. Tsitsiklis, "Actor-critic algorithms," *Advances in neural information processing systems*, vol. 12, 1999.
- [42] P. Christodoulou, "Soft actor-critic for discrete action settings," *arXiv preprint arXiv:1910.07207*, 2019.
- [43] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," in *International conference on machine learning*. PMLR, 2018, pp. 1861–1870.
- [44] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel *et al.*, "Soft actor-critic algorithms and applications," *arXiv preprint arXiv:1812.05905*, 2018.
- [45] M. D. Soltani, A. A. Purwita, Z. Zeng, H. Haas, and M. Safari, "Modeling the random orientation of mobile devices: Measurement, analysis and lifi use case," *IEEE Transactions on Communications*, vol. 67, no. 3, pp. 2157–2172, 2018.
- [46] C. Zhou, W. Wu, H. He, P. Yang, F. Lyu, N. Cheng, and X. Shen, "Delay-aware iot task scheduling in space-air-ground integrated network," in *2019 IEEE Global Communications Conference (GLOBECOM)*. IEEE, 2019, pp. 1–6.
- [47] X. Kang, Y.-K. Chia, S. Sun, and H. F. Chong, "Mobile data offloading through a third-party wifi access point: An operator's perspective," *IEEE Transactions on Wireless Communications*, vol. 13, no. 10, pp. 5340–5351, 2014.
- [48] H. Yu, M. H. Cheung, G. Iosifidis, L. Gao, L. Tassiulas, and J. Huang, "Mobile data offloading for green wireless networks," *IEEE Wireless Communications*, vol. 24, no. 4, pp. 31–37, 2017.
- [49] A. Y. Ding, Y. Liu, S. Tarkoma, H. Flinck, H. Schulzrinne, and J. Crowcroft, "Vision: Augmenting wifi offloading with an open-source collaborative platform," in *Proceedings of the 6th International Workshop on Mobile Cloud Computing and Services*, 2015, pp. 44–48.
- [50] J. Korhonen, T. Savolainen, A. Y. Ding, and M. Kojo, "Toward network controlled ip traffic offloading," *IEEE Communications Magazine*, vol. 51, no. 3, pp. 96–102, 2013.
- [51] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, "Openai gym," *arXiv preprint arXiv:1606.01540*, 2016.



**Yeguang Qin** received the M.E. degree from the Department of Software Engineering, Xinjiang University, in 2023. He is currently the Ph.D. Student with the Department of Computer Science and Technology, Central South University, advised by Prof. Ming Zhao. His research work focuses on wireless networks, network traffic control, and machine learning algorithm.



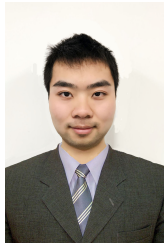
**Ming Zhao** received the M.Sc. and Ph.D. degrees in computer science from Central South University, Changsha, China, in 2003 and 2007, respectively. He is currently a Professor with the School of Computer Science and Engineering, Central South University. His main research focuses on wireless networks. He is also a Member of the China Computer Federation.



**Yilin Yang** received the M.E. degree from the Department of Computer Science and Technology, Northeast Forestry University, in 2020 and the M.S. degree in software engineering from the School of Computer Science and Engineering, Central South University, Changsha, China, in 2023. His research work focuses on network traffic control and machine learning algorithm.



**Nei Kato** is a Full Professor and the Dean with the Graduate School of Information Sciences, Tohoku University. He has researched on computer networking, wireless mobile communications, satellite communications, ad hoc and sensor and mesh networks, UAV networks, smart grid, AI, IoT, Big Data, and pattern recognition. He has published more than 500 papers in prestigious peerreviewed journals and conferences. He served as the Vice-President (Member and Global Activities) of IEEE Communications Society from 2018 to 2021, and the Editor-in-Chief of IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY from 2017 to 2021. He is the Editor-in-Chief of IEEE INTERNET OF THINGS JOURNAL. He is a Fellow of the Engineering Academy of Japan and IEICE.



**Fengxiao Tang** received the B.E. degree in measurement and control technology and instrument from the Wuhan University of Technology, Wuhan, China, in 2012 and the M.S. degree in software engineering from the Central South University, Changsha, China, in 2015. He received the Ph.D. degrees from the Graduate School of Information Science, Tohoku University, Japan. Currently, He is an full professor in the School of Computer Science and Engineering of Central South University. He has been an Assistant Professor from 2019 to 2020

and an Associate Professor from 2020 to 2021 at the Graduate School of Information Sciences (GSIS) of Tohoku University. His research interests are unmanned aerial vehicles system, IoT security, game theory optimization, network traffic control and machine learning algorithm. He was a recipient of the prestigious Dean's and President's Awards from Tohoku University in 2019, and several best paper awards at conferences including IC-NIDC 2018, GLOBECOM 2017/2018. He was also a recipient of the prestigious Funai Research Award in 2020, IEEE ComSoc Asia-Pacific (AP) Outstanding Paper Award in 2020 and IEEE ComSoc AP Outstanding Young Researcher Award in 2021.



**Xin Yao** received the B.S. degree in computer science from Xidian University, Xi'an, China, in 2011, the M.S. degree in software engineering and the Ph.D. degree in computer science and technology from Hunan University, Changsha, China, in 2013 and 2018, respectively. From 2015 to 2017, he worked as a Visiting Scholar with Arizona State University. He is currently an Assistant Professor with Central South University. His research interests include security and privacy issues in social network, Internet of Things, cloud computing, and big data