

Efficient Communications in Multi-Agent Reinforcement Learning for Mobile Applications

Zefang Lv, *Student Member, IEEE*, Liang Xiao, *Senior Member, IEEE*, Yousong Du, Yunjun Zhu,
Shuai Han, *Senior Member, IEEE*, and Yong-Jin Liu, *Senior Member, IEEE*

Abstract—The environment observations and learning experiences shared by the cooperative learning agents accelerate multi-agent reinforcement learning (MARL) with partial observations for mobile applications but the performance degrades due to the redundant and outdated observations under severe channel fading in wireless networks. In this paper, we propose an efficient communication scheme in MARL for mobile applications that enables each learning agent to optimize the cooperative agents and the learning parameters to integrate the shared information. The cooperative agents are chosen according to the learning environment observations, the channel states, and the task similarity with neighboring agents. The learning parameters are chosen based on the attention mechanism that exploits the correlation with the local observation to enhance the agent receptive field for efficient policy exploration. Neural networks with weights updated based on the learning factors determined by the task similarity are designed to further improve the learning efficiency. The performance bounds including the information gain from the learning agent cooperation, the communication cost and the utility are provided based on the Nash equilibrium of the cooperative MARL communication game. The proposed scheme is implemented in the anti-jamming video transmission of the unmanned aerial vehicle swarms to optimize the transmit channel and power and experimental results verify the performance gain over the benchmark.

Index Terms—Multi-agent reinforcement learning, communications, unmanned aerial vehicle, mobile applications, wireless networks.

I. INTRODUCTION

THE latest development of multi-agent reinforcement learning (RL) has enabled mobile devices as the learning agents to optimize their decisions such as the radio resource allocation and trajectory design to support a variety of mobile applications and services in wireless networks [1], [2]. However, most learning agents only have partial observations on the environment such as the radio channel states and suffer from non-stationarity caused by the changing policies of other agents, which in turn degrades the RL performance and

This work was supported in part by the Natural Science Foundation of China under Grant U21A20444 and the National Key Research and Development Program of China under Grant 2023YFB3107603, and in part by the Natural Science Foundation of China under Grants 62332019 and U2336214. (*Corresponding author: Liang Xiao*)

Zefang Lv, Liang Xiao, Yousong Du and Yunjun Zhu are with the Department of Informatics and Communication Engineering & Key Laboratory of Multimedia Trusted Perception and Efficient Computing, Xiamen University, Xiamen 361005, China. Email: lxiao@xmu.edu.cn

Shuai Han is with the School of Electronics and Information Engineering, Harbin Institute of Technology, China. Email: hanshuai@hit.edu.cn

Yong-Jin Liu is with the Department of Computer Science and Technology, Tsinghua University. Email: liuyongjin@tsinghua.edu.cn

the quality of service (QoS) for partially observable Markov decision process (POMDP) in dynamic networks [3].

As the neighboring learning agents share a common environment and can exchange their environment observations and learning experiences via the wireless networks, the multi-agent RL (MARL) algorithms can exploit the learning information shared from the cooperative agents in the state formulation and learning parameter update to enhance the agent receptive field and improve the RL policy optimization efficiency [4], [5]. For example, the multi-agent deep Q-network (DQN)-based interference mitigation scheme in [5] exploits the shared observations such as the signal-to-interference-plus-noise ratio (SINR) and the weights of associated user equipment to optimize the downlink transmit power and user scheduling and thus improves the data rate in wireless networks against interference from neighboring access points.

Both the multi-agent communication efficiency and RL performance depend on the cooperative agent selection, which in turn relies on the accurate network states that are rarely known by most mobile devices. As a seminal work on multi-agent communications, I2C proposed in [6] applies a prior network consisting of a feed-forward neural network to evaluate the belief of the neighboring agents based on local observations, such as the positions of the neighboring agents in the navigation task, to decide which agent to communicate with. However, the performance of I2C is hindered by long communication delay and observation redundancy due to the limited bandwidth and time-varying channel states, leading to a degradation in RL performance and application QoS in large-scale wireless networks.

In this paper, we propose an efficient communication scheme in MARL for mobile applications that optimizes the cooperative agents to share environment observations and learning parameters to integrate the shared learning information to accelerate task learning with restricted communication overhead. The cooperative agents are chosen based on a communication state that consists of the local environment observation of the learning agent, the radio channel states, the task similarity with neighboring agents and previous learning performance to maximize the expected long-term discounted communication reward as the weighted sum of RL task reward and communication cost.

A scaled dot-product attention-based observation encoder compares the shared and local observations to dynamically extract observation correlations and enhance the receptive field for learning agents, while reducing the redundancy in the RL task state formulation to accelerate the task policy selection.

The local task learning parameters are updated based on both the local task reward and the shared learning parameters with learning factors determined by the task similarity to further improve the task policy exploration efficiency.

A communication scheme for deep MARL that exploits both the observations and learning parameters shared from the selected cooperative agent is proposed to accelerate learning of neural networks in the task policy optimization. A sequential architecture with fully-connected layer based neural networks in the communication learning layer estimates the expected long-term communication reward to choose the cooperative agents for network weights and observation sharing. The neural network for task policy optimization with weights updated based on the shared learning experiences and their corresponding importance adaptively determined by the task similarity are designed to further improve the learning efficiency.

We formulate a cooperative MARL communication game among learning agents based on the informational model presented in [7] to provide the communication learning performance bounds. In this game, each learning agent chooses whether to communicate with neighboring agents with the goal of maximizing the RL task performance with less communication cost. Based on the Nash equilibrium of the game, the performance bounds including the information gain, the communication cost and the utility are provided and verified via simulations, in which the information gain obtained from the learning agent cooperation increases with the task similarity with neighboring agents and shared information of the selected cooperative agents.

As a case study, we implement the proposed communication schemes in MARL-based anti-jamming video transmission of unmanned aerial vehicle (UAV) swarm to optimize the transmit channel and power. Simulations are performed in a 10-UAV swarm based on the collected channel states of the swarm communication and jammer-UAV links to show the performance gain of the proposed scheme over the benchmark I2C in [6]. In addition, the proposed scheme is also implemented in a 5-UAV swarm equipped with Raspberry Pi 4 against a jammer, i.e., a universal software radio peripheral (USRP) N210 controlled by a laptop. Both the simulation and experimental results show that the proposed communication scheme improves the MARL-based UAV swarm anti-jamming video transmission performance with reduced packet loss rate (PLR), transmission delay, UAV transmit power and convergence time over the benchmark I2C in [6].

The main contributions of this paper are summarized as follows:

- We propose an efficient communication scheme to improve both the learning and communication efficiency for MARL-based mobile applications, which enables each learning agent to optimize the cooperative agents for environment observation and learning experience sharing and the parameters to integrate the shared information in the RL task policy learning.
- We formulate the interactions among learning agents as a cooperative MARL communication game based on the informational model and provide the performance

bounds including the information gain from the learning information sharing, the communication cost and utility based on the Nash equilibrium of the game.

- We implement the MARL with the proposed scheme in the UAV swarm anti-jamming video transmission system to optimize the transmit channel and power. Both simulations and experiments are performed to show the performance gain in terms of PLR and transmission delay over the benchmark.

The rest of this paper is organized as follows. We review the related work in Section II, and present the system model in Section III. An efficient communication scheme for MARL and a version for deep MARL are presented in Sections IV and V. The performance analysis and a case study is given in Sections VI and VII, respectively. Simulation and experimental results are reported in Sections VIII and IX, followed by the conclusion in Section X.

II. RELATED WORK

MARL such as multi-agent deep deterministic policy gradient (DDPG) enables multi-agent systems in wireless networks to optimize transmission policies such as resource allocation and trajectory design [1], [2], [8]–[11]. For example, the multi-UAV assisted mobile edge computing system in [2] applies the multi-agent twin delayed DDPG to optimize the trajectory design, task allocation, and power management and thus decreases the total system cost. The full-duplex multi-UAV network system in [8] uses a clip and count-based proximal policy optimization algorithm to optimize the de-coupled uplink-downlink association and trajectory design in a distributed manner. However, fully decentralized learning or centralized training and decentralized executing have learning performance degradation under partial observation in large-scale wireless networks.

Learning information exchange including environment observations and learning experiences among learning agents such as mobile devices enhances the RL performance under partial observation [5], [7], [12]–[18]. For example, the MARL-based resource management scheme in [5] enables the radio devices to exchange observations with their neighboring agents to improve the learning performance of power control and user association. The multi-agent wireless system in [7] shares information with other agents via neural network parameters and analyzes the effect of the parameter sharing frequency on the convergence speed. A multi-agent DQN based cellular offloading scheme in [16] optimizes trajectory design and power allocation with experiences shared via a neural network and decreases the convergence time. However, the full communication protocols above that share observations or learning parameters among all the learning agents may lead to higher communication overhead and introduce information redundancy in the task learning process, thus degrading the learning performance and speed.

Efficient communication is promising for MARL to enhance learning performance with attention experience sharing, and as examples, RIAL and DIAL in [3] and CommNet in [19] learn to communicate between agents in cooperative environments

with partial observations. For example, the communication protocol in [3] uses deep Q-learning with recurrent networks to encode past observations, actions and local observations into binary messages and thus reduces communication overhead but has information loss in large-scale wireless networks. The communication information mainly consists of original observations, encoded observations [4], [20] and intended actions [21]. In particular, an intention sharing scheme in [21] applies the attention mechanism to generate an imagined trajectory based on the shared messages from other agents and enhance the coordination of multi-agent systems.

Learning cooperative agent selection and integration of shared information further enhance the learning efficiency [6], [22]. For example, an attentional communication scheme in [22] uses a bidirectional long short-term memory unit to learn when to communicate and how to integrate the received information for cooperative decision making. An individually inferred communication scheme in [6] applies causal inference to pretrain a prior network for agent-agent communications, which maps the local observation to a belief value to determine which agent to communicate with. However, the limited radio bandwidth and time-varying channel states resulting in long communication delay may degrade the RL performance and application QoS in large-scale wireless networks.

Existing communication schemes in MARL also consider realistic constraints such as communication cost and dynamic environments [23]–[27]. For example, an information bottleneck based communication scheme in [24] jointly optimizes which agent is communicating what message and to whom under limited bandwidth and estimates the posterior distribution of messages based on the variational information bottleneck. An efficient communication scheme in [25] investigates communications among multiple agents over a noisy channel and jointly optimizes the channel coding, source coding and joint source-channel coding problems to enhance coordination and cooperation among agents. However, these communication schemes based on fully connected structure suffer from higher communication overhead and QoS degradation for mobile applications in large-scale wireless networks.

Compared with MARL communication schemes mentioned above, we optimize the cooperative agents for learning information sharing based on task similarity and radio channel states with neighboring agents and the information integration method to enhance both the learning and communication efficiency for mobile applications in large-scale dynamic wireless networks. In our previous work in [28], an efficient communication scheme for MARL is proposed, which enables each learning agent to optimize the cooperative agent selection and the task state formulation to improve learning performance and QoS of RL-based applications in wireless networks. In this paper, we further propose a sequential architecture to choose cooperative agents with lower computational complexity and design neural networks with weight updated with learning factor determined by the task similarity to improve the learning efficiency. The proposed schemes are implemented in the UAV swarm anti-jamming video transmission to optimize transmit channel and power and experimental results verify the performance gain over the benchmark.

TABLE I: Summary of Symbols and Notations.

| Symbol | Description |
|----------------------------|--|
| N | Number of neighboring agents |
| $\mathbf{o}_i^{(k)}$ | Local observation of the learning agent at time slot k |
| $\hat{\mathbf{o}}_i^{(k)}$ | Received observation from agent i |
| $\bar{\mathbf{o}}_i^{(k)}$ | Encoded observation of agent i |
| $\bar{\mathbf{o}}^{(k)}$ | Previous received observations |
| $\eta_i^{(k)}$ | RL task similarity with agent i |
| $h_i^{(k)}$ | Channel gain with agent i |
| $\tilde{s}^{(k)}$ | Communication learning state |
| $\mathbf{x}^{(k)}$ | Selected cooperative agents |
| $s^{(k)}$ | Task learning state |
| $\mathbf{a}^{(k)}$ | Task action |
| $\rho^{(k)}$ | Task learning performance |
| $r^{(k)}$ | Task reward |
| $\xi^{(k)}$ | Number of selected cooperative agents |
| $\tilde{r}^{(k)}$ | Communication reward |
| $p^{(k)}$ | Transmit power of the UAV |
| $f^{(k)}$ | Transmit channel of the UAV |
| $b^{(k)}$ | UAV message size |
| $\zeta^{(k)}$ | Jamming power received at the control center |
| $\rho_1^{(k)}$ | SINR of the UAV signals received by the control center |
| $\rho_2^{(k)}$ | Transmission delay |
| $\rho_3^{(k)}$ | Packet loss rate |
| $\varrho^{(k)}$ | Frequency hopping cost |

III. SYSTEM MODEL

A. MARL Model for Mobile Applications

In a POMDP, N learning agents communicate via wireless channels to cooperatively perform learning tasks of mobile applications such as target tracking and collaborative sensing. At time slot k , learning agent i obtains local observation $\mathbf{o}_i^{(k)}$, selects an action $\mathbf{a}_i^{(k)} \in \mathcal{A}_i$ that follows the policy distribution $\pi_i \in [0, 1]^{\lvert \mathcal{A}_i \rvert}$, receives a reward $r_i^{(k)}$ from the environment and transitions to a new state $s_i^{(k+1)}$. The goal is to learn the policy π_i that maximizes the expected long-term discounted reward $R_i^{(k)} = \sum_{n=k}^{\infty} \gamma^{n-k} r_i^{(n)}$, where γ is the discount factor that represents the importance of the immediate reward and the future reward. The POMDP is modeled by a tuple $\mathcal{G} = \{\mathcal{S}, \mathcal{O}_i, \mathcal{A}_i, \mathcal{P}, r_i, \pi_i\}_{1 \leq i \leq N}$, where

- State space \mathcal{S} represents all possible state configurations of the environment and learning agents.
- Observation set \mathcal{O}_i consists of all local observations of learning agent i .
- Action set \mathcal{A}_i consists of all feasible task policies of learning agent i .
- $\mathcal{P}: \mathcal{S} \times \mathcal{A}_1 \times \cdots \times \mathcal{A}_N \times \mathcal{S} \rightarrow [0, 1]$ describes the distribution of possible next state given the current state and actions of all learning agents.
- Reward $r_i: \mathcal{S} \times \mathcal{A}_1 \times \cdots \times \mathcal{A}_N \rightarrow \mathbb{R}$ is the immediate reward as the weighted sum of task performance such as

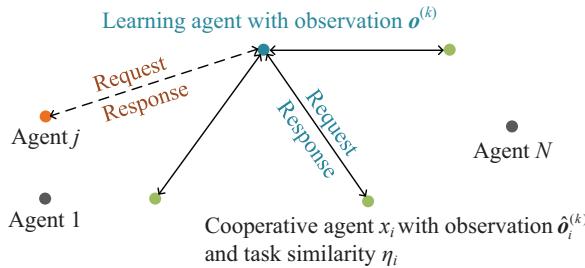


Fig. 1. A wireless network with $N + 1$ cooperative agents performing RL-based tasks such as target tracking, in which the learning agent selects cooperative agents $\mathbf{x}^{(k)}$ to request learning information such as environment observations and learning experiences from neighboring agents to enhance the performance of mobile applications.

video transmission quality and target tracking accuracy, which depends on the global state and actions of all agents.

- Policy distribution $\pi_i: \mathcal{S} \times \mathcal{A}_i \rightarrow [0, 1]^{|\mathcal{A}_i|}$ represents the probability of learning agent i to choose the action $a_i \in \mathcal{A}_i$ in state s_i .

B. Network and Communication Model for Learning Agents

Without loss of generality, each learning agent is assumed to have N neighboring agents $\mathcal{N} = \{1, 2, \dots, N\}$ and exchanges learning messages including observations and learning parameters with neighboring agents within the communication range based on the carrier sense multiple access with collision avoidance access via the wireless channels. As shown in Fig. 1, at time slot k , the learning agent observes local observation $\mathbf{o}^{(k)} \in \mathbb{R}^O$, estimates the channel states $\mathbf{h}^{(k)} = [h_{1 \leq i \leq N}^{(k)}]$ and evaluates the task similarity with neighboring agents $\boldsymbol{\eta}^{(k)} = [\eta_{1 \leq i \leq N}^{(k)}] \in [0, 1]^N$ based on the Jensen-Shannon distance between local and shared policy distribution or Bisimulation relation of previous shared observations as [29], [30]. For example, an encoder with convolutional layers learns the task-relevant information from local observations $\mathbf{o}^{(k)}$ and previous shared observations $\bar{\mathbf{o}}_i$ to evaluate the similarity as a Bisimulation metric $d(\mathbf{o}^{(k)}, \bar{\mathbf{o}}_i)$ as presented in [31].

The learning agent chooses at most M cooperative agents $\mathbf{x}^{(k)} \subset \mathcal{N}$ to request the learning information and formulate the task learning state for more efficient RL task policy optimization. Specifically, the learning agent collects information about neighboring agents when participating in a new RL task, such as the device identities, learning state configuration and task policy from the previous learning response messages of other agents. After determining the cooperative agents $\mathbf{x}^{(k)}$ from N neighboring agents $\mathcal{N} = \{1, 2, \dots, N\}$, the selected agents ID and required information such as the observation and current learning parameters, e.g., the task Q-values or policy distribution, are broadcasted to the selected agents $\mathbf{x}^{(k)}$ in the learning request message.

Upon receiving the learning request message from the learning agent, the cooperative agent $i \in \mathbf{x}^{(k)}$ obtains its local observation $\hat{o}_i^{(k)}$ and current learning parameters and sends the learning information in the learning response message back to

the learning agent. If neighboring agent i receives multiple cooperative requests from other agents, the current learning information will be sent to these request learning agents in the learning response message sequentially.

The cooperative states $\{\hat{o}_i^{(k)} | i \in \mathbf{x}^{(k)}\}$ are exploited to formulate the task learning state $\mathbf{s}^{(k)}$, which is used as the basis for the learning agent to choose the task action $\mathbf{a}^{(k)} \in \mathcal{A}$. Based on the task performance $\rho^{(k)}$ such as the PLR of video transmission and message bit error rate (BER), the learning agent evaluates the task reward $r^{(k)}$ and communication cost, i.e., the number of selected cooperative agents denoted by $\xi^{(k)}$, to guide the cooperative agent selection. Time slot k and agent ID i are omitted if no confusion occurs in the subsequent sections because our proposed algorithms are applicable to each learning agent and only depend on the latest learning information such as observations, learning performance and parameters instead of all the learning information of previous time slots.

IV. EFFICIENT COMMUNICATIONS FOR MARL

We propose a learning based efficient multi-agent communication scheme with sequential architecture named ECOM for MARL-based mobile applications in wireless networks that enables learning agents to choose their cooperative agents $\mathbf{x}^{(k)}$ and task action $\mathbf{a}^{(k)}$ for faster task learning under partial observation. The cooperative agents $\mathbf{x}^{(k)}$ are chosen in each time slot k based on a communication state $\tilde{\mathbf{s}}^{(k)}$ including local observation \mathbf{o} , channel states \mathbf{h} and task similarity $\boldsymbol{\eta}$ with neighboring agents, and previous communication cost ξ . According to the attention mechanism similar to [32], the correlation between the local and shared observations denoted by \tilde{o}_i , $i \in \mathbf{x}^{(k)}$, is calculated based on the attention function, i.e., $\tilde{o}_i = m(\mathbf{o}, \mathbf{o}_i)$, to reduce the information redundancy in the formulation of the task learning state $\mathbf{s}^{(k)}$ for task policy selection. In addition, the shared learning parameters such as task Q-values are exploited to update the local task Q-values based on the learning factors determined by the task similarity for efficient task policy exploration.

As shown in Fig. 2, the MARL communication scheme ECOM jointly chooses the cooperative agents $\mathbf{x}^{(k)}$ and task action $\mathbf{a}^{(k)}$ to maximize the discounted long-term reward of the learning agent. The optimal cooperative agent selection that determines both the multi-agent communication efficiency and RL performance depends on the accurate network states and perfect knowledge on the learning information of neighboring agents, which are rarely known by most mobile devices accurately. Therefore, the communication learning layer adaptively chooses the cooperative agents $\mathbf{x}^{(k)}$ based on the communication state $\tilde{\mathbf{s}}^{(k)}$ and communication policy distribution π^C for mobile devices under time-varying channel states via ϵ -greedy based learning method. Specifically, the communication state $\tilde{\mathbf{s}}^{(k)}$ consists of local observations \mathbf{o} , channel states \mathbf{h} , task similarity $\boldsymbol{\eta}$, previous observations of the W neighboring agents $\bar{\mathbf{o}}$ and previous communication cost ξ , i.e.,

$$\tilde{\mathbf{s}}^{(k)} = [\mathbf{o}, \mathbf{h}, \boldsymbol{\eta}, \bar{\mathbf{o}}, \xi] \in \hat{\mathbf{S}} \quad (1)$$

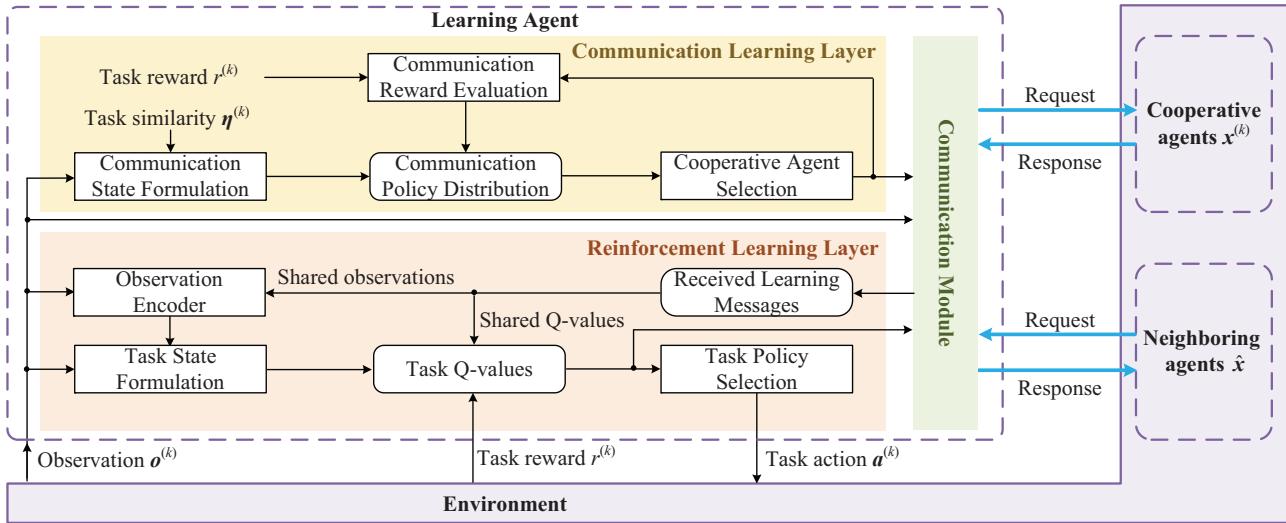


Fig. 2. Illustration of the MARL communication scheme, in which the learning agent chooses the cooperative agents $x^{(k)}$ in the communication learning layer, request and send learning information to neighboring agents via the communication module and selects the task action $a^{(k)}$ in the reinforcement learning layer based on the aggregated task learning state $s^{(k)}$.

Algorithm 1 Efficient communications for MARL

```

1: Initialize  $\bar{o}$ ,  $\xi$ ,  $\lambda$ ,  $\gamma$ ,  $\varpi$  and  $\omega$ ,  $Q^R$  and  $Q_{1 \leq j \leq M}^C$ 
2: for  $k = 1, \dots, K$  do
3:   Obtain local observation  $\mathbf{o}^{(k)}$ 
4:   Estimate channel gain  $h^{(k)}$  and task similarity  $\eta^{(k)}$  with neighboring agents
5:   Formulate communication state  $\tilde{s}^{(k)}$  via (1)
6:   for  $j = 1, \dots, M$  do
7:     Select  $x_j^{(k)}$  via (2)
8:   end for
9:   Send learning request to agents  $\mathbf{x}^{(k)}$  and receive  $\{\hat{o}_i, \hat{Q}_i^R\}_{i \in \mathbf{x}^{(k)}}$ 
10:  Calculate  $\tilde{o}_i = m(\mathbf{o}, \hat{o}_i)$ ,  $\forall i \in \mathbf{x}^{(k)}$ 
11:  Formulate task state  $s^{(k)} = [\mathbf{o}^{(k)}, [\tilde{o}_i]_{i \in \mathbf{x}^{(k)}}]$ 
12:  Select  $a^{(k)}$  via (4)
13:  Receive task reward  $r^{(k)}$ 
14:  Update  $Q^R(s^{(k)}, a^{(k)})$  via (5)
15:  Update  $Q^R$  via (6)
16:  Calculate  $\tilde{r}^{(k)}$  via (7)
17:  for  $j = 1, \dots, M$  do
18:    Update  $Q_j^C(\tilde{s}^{(k)}, x_j^{(k)})$  similar to (5)
19:  end for
20: end for

```

A sequential cooperative agent selection architecture is designed for large-scale networks to accelerate the communication policy learning for choosing a large number of cooperative agents. This architecture decomposes the communication policy $\mathbf{x}^{(k)}$ into M levels to sequentially choose at most M cooperative agents from neighboring agents $\mathcal{N} = \{1, 2, \dots, N\}$, which significantly decreases the computational complexity of the communication policy selection and learning. Specifically, level j consists of an $|\mathcal{S}| \times (N + 1)$ Q-table to output the

Q-values of the j -th cooperative agent under current state $\tilde{s}^{(k)}$ denoted by $Q_j^C(\tilde{s}^{(k)}, \cdot)$ to choose $x_j \in \{0, 1, 2, \dots, N\}$, where $x_j = n \in \mathcal{N}$ indicates that neighboring agent n is selected in level j and $x_j = 0$ indicates no selected agent in this level. Based on the ϵ -greedy method, the communication policy distribution in level j is designed to make a trade-off between exploration and exploitation in the communication policy learning, which is formulated as

$$\Pr(x_j = \hat{x}) = \begin{cases} 1 - \epsilon, & \hat{x} = \underset{x \in \{0, 1, \dots, N\}}{\operatorname{argmax}} Q_j^C(\tilde{s}^{(k)}, x) \\ \frac{1}{N}, & \hat{x} \neq \underset{x \in \{0, 1, \dots, N\}}{\operatorname{argmax}} Q_j^C(\tilde{s}^{(k)}, x) \end{cases} \quad (2)$$

Upon receiving the learning response message from the selected cooperative agents $\mathbf{x}^{(k)}$, the shared observation \hat{o}_i and task Q-values \hat{Q}_i^R , $i \in \mathbf{x}^{(k)}$, are extracted to assist task learning. The observation encoder $m(\mathbf{o}, \hat{o}_i)$ applies the scaled dot-product attention mechanism in [32] to compare the local observation \mathbf{o} with the shared observation \hat{o}_i , $i \in \mathbf{x}^{(k)}$, and extract the importance of the shared observation, which is used to encode observation \hat{o}_i in the formulation of the task learning state $s^{(k)}$. Specifically, local observation \mathbf{o} and shared observations \hat{o}_i from cooperative agent i both serve as the input of the observation encoder. The local observation \mathbf{o} is transformed to *query* with $\mathbf{W}^Q \in \mathbb{R}^{O \times V}$ and the shared observations \hat{o}_i is transformed to *key* and *value* with \mathbf{W}^K and $\mathbf{W}^V \in \mathbb{R}^{O \times V}$, respectively. The encoded observation denoted by \tilde{o}_i depends on the weighted sum of *value* and the dot product and softmax function is applied to calculate the weight for cooperative agent $i \in \mathbf{x}^{(k)}$ denoted by $\alpha_i = [\alpha_{i,v}]_{v=1}^V$, i.e.,

$$\alpha_{i,v} = \frac{\exp(\mathbf{o} \mathbf{W}_v^Q \hat{o}_i \mathbf{W}_v^K)}{\sum_{v=1}^V \exp(\mathbf{o} \mathbf{W}_v^Q \hat{o}_i \mathbf{W}_v^K)} \quad (3)$$

Based on the shared observation \hat{o}_i and its transforming matrix \mathbf{W}^V and weights α_i , the encoded observation $\tilde{o}_i = \hat{o}_i \mathbf{W}^V \alpha_i^T$.

The reinforcement learning layer optimizes the specific RL task action $\mathbf{a}^{(k)}$ based on the task learning state $\mathbf{s}^{(k)}$, task Q-values $Q^R(\mathbf{s}^{(k)}, \cdot)$ and task policy distribution π^R . Specifically, the task state $\mathbf{s}^{(k)}$ consists of local observation \mathbf{o} and encoded shared observations of the cooperative agents, i.e., $\mathbf{s}^{(k)} = [\mathbf{o}, [\tilde{\mathbf{o}}_i]_{i \in \mathbf{x}^{(k)}}]$. The task policy $\mathbf{a} \in \mathcal{A}$ is chosen based on the ϵ -greedy method according to the Q-values of task state-action pairs $Q^R(\mathbf{s}^{(k)}, \cdot)$ or the task policy distribution π^R with

$$\Pr(\mathbf{a}^{(k)} = \hat{\mathbf{a}}) = \pi^R(\hat{\mathbf{a}} | \mathbf{s}^{(k)}) \quad (4)$$

Upon receiving the task feedback from the environment, task reward $r^{(k)}$ is evaluated based on the task performance ρ such as the video transmission quality. The Q-value in the task layer $Q^R(\mathbf{s}^{(k)}, \mathbf{a}^{(k)})$ is updated with reward $r^{(k)}$ by Bellman equation

$$Q^R(\mathbf{s}^{(k)}, \mathbf{a}^{(k)}) \leftarrow \lambda \left(r^{(k)} + \gamma \max_{\hat{\mathbf{a}} \in \mathcal{A}} Q^R(\mathbf{s}^{(k+1)}, \hat{\mathbf{a}}) \right) + (1 - \lambda) Q^R(\mathbf{s}^{(k)}, \mathbf{a}^{(k)}), \quad (5)$$

where λ and γ represents the importance of the received reward and the future reward in the learning process, respectively. In addition, the shared learning information such as Q-values \hat{Q}_i^R are also exploited to update local Q-values, i.e.,

$$Q^R \leftarrow \frac{Q^R + \sum_{i \in \mathbf{x}^{(k)}} \omega_i \hat{Q}_i^R}{1 + \sum_{i \in \mathbf{x}^{(k)}} \omega_i} \quad (6)$$

The learning factor ω_i for cooperative agent i is adaptively determined based on the task similarity or previous communication performance to further improve the task learning efficiency. For example, the learning factor ranging from 0 to 1 increases linearly with the task similarity between the learning agent and cooperative agent i .

The communication cost, i.e., the number of the selected cooperative agents $\xi^{(k)}$, is evaluated to calculate the communication reward $\tilde{r}^{(k)}$, which increases with task reward $r^{(k)}$ and decreases with the communication cost $\xi^{(k)}$ and the corresponding importance factor ϖ^C , i.e.,

$$\tilde{r}^{(k)} = r^{(k)} - \varpi^C \xi^{(k)} \quad (7)$$

The Q-values in the M levels of the communication learning layer $Q_j^C(\tilde{\mathbf{s}}^{(k)}, \mathbf{x}_j^{(k)})$, $1 \leq j \leq M$ are updated with communication reward $\tilde{r}^{(k)}$ similar to (5).

V. EFFICIENT COMMUNICATIONS FOR DEEP MARL

We further propose an efficient communication scheme with sequential architecture for the deep MARL that exploits the shared learning parameters to further accelerate the communication and task policy learning. The communication Q-values of each level in the cooperative agent selection are estimated via the fully connected neural networks to compress the high-dimension state for the large-scale wireless networks. The shared learning parameters such as the neural network weights are exploited to update the task network for fast learning.

The communication learning layer chooses the cooperative agent $\mathbf{x}^{(k)}$ via the sequential cooperative agent selection architecture consisting of M levels. Each level j consists

Algorithm 2 Efficient communications for deep MARL

```

1: Initialize  $\bar{\mathbf{o}}, \xi, \lambda, \gamma, \varpi, \theta, \omega$  and  $\alpha_{1 \leq j \leq M}$ 
2: for  $k = 1, \dots, K$  do
3:   Same as lines 3-5 in Algorithm 1
4:   for  $j = 1, \dots, M$  do
5:     Input  $\tilde{\mathbf{s}}^{(k)}$  to C-Network  $j$  and obtain  $Q_j^C(\tilde{\mathbf{s}}^{(k)}, \cdot; \alpha_j)$ 
6:     Select  $x_j^{(k)}$  via (2)
7:   end for
8:   Send learning request to agents  $\mathbf{x}^{(k)}$  and receive
    $\{\{\hat{\mathbf{o}}_i, \hat{\theta}_i\} | i \in \mathbf{x}^{(k)}\}$ 
9:   Formulate  $\mathbf{s}^{(k)}$  similar to Algorithm 1
10:  Input  $\mathbf{s}^{(k)}$  to task network and obtain  $Q^R(\mathbf{s}^{(k)}, \cdot)$ 
11:  Select task action  $\mathbf{a}^{(k)}$  via (4)
12:  Receive task reward  $r^{(k)}$ 
13:   $\mathcal{D}^R \leftarrow \mathcal{D}^R \cup z^R$ 
14:  if  $|\mathcal{D}^R| \geq G$  then
15:    Randomly sample  $G$  experiences from  $\mathcal{D}^R$  to update
         $\theta$  via (9)
16:  end if
17:  Update  $\theta$  via (10)
18:  Evaluate communication reward  $\tilde{r}^{(k)}$  via (7)
19:   $\mathcal{D}^C \leftarrow \mathcal{D}^C \cup z^C$ 
20:  if  $|\mathcal{D}^C| \geq G$  then
21:    Randomly sample  $G$  experiences from  $\mathcal{D}^C$ 
22:    for  $j = 1, \dots, M$  do
23:      Update  $\alpha_j$  similar to (9)
24:    end for
25:  end if
26: end for

```

of a neural network named C-Network j with weight α_j including four fully-connected layers, i.e., an input layer with $(W + 1)O + 2N + 1$ nodes, two hidden layers with $g_{j,1}$ and $g_{j,2}$ nodes and an output layer with $N + 1$ nodes. With the communication state $\tilde{\mathbf{s}}^{(k)}$ similar to (1) as input, the C-Network j estimates the Q-values $Q_j^C(\tilde{\mathbf{s}}^{(k)}, \cdot; \alpha_j)$ to calculate the communication policy distribution via (2).

The task learning state $\mathbf{s}^{(k)}$ is formulated as in Algorithm 1 and input to the task network such as a fully-connected neural network to obtain the Q-values of state-action pairs $Q^R(\mathbf{s}^{(k)}, \cdot)$. The task action $\mathbf{a} \in \mathcal{A}$ is chosen according to the task policy distribution $\pi^R(\cdot | \mathbf{s}^{(k)}; \theta)$ formulated based on $Q^R(\mathbf{s}^{(k)}, \mathbf{a})$ with ϵ -greedy method. Upon receiving the task feedback, task reward $r^{(k)}$ is evaluated to update the network weights θ with the experience replay technique. The task learning experience including the task state $\mathbf{s}^{(k)}$, task action $\mathbf{a}^{(k)}$, task reward $r^{(k)}$ and next task state $\mathbf{s}^{(k+1)}$, i.e.,

$$z^R = \{\mathbf{s}^{(k)}, \mathbf{a}^{(k)}, r^{(k)}, \mathbf{s}^{(k+1)}\} \quad (8)$$

The task memory pool \mathcal{D}^R stores the experience z^R and G experiences are randomly and uniformly sampled from \mathcal{D}^R to formulate the minibatch \mathcal{B}^R . The Q-values $Q^R(s, a; \theta)$ in the reinforcement learning layer estimated by the task network are used to update the network weights θ via stochastic gradient descent algorithm such as Adam optimizer with learning rate

β by minimizing the loss function given by

$$\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{\{\mathbf{s}, \mathbf{a}, r, \bar{s}\} \in \mathcal{B}^R} \left[\left(r + \gamma \max_{\mathbf{a} \in \mathcal{A}} Q^R(\bar{s}, \mathbf{a}; \boldsymbol{\theta}) \right. \right. \\ \left. \left. - Q^R(s, \mathbf{a}; \boldsymbol{\theta}) \right)^2 \right] \quad (9)$$

The shared network weights $\{\hat{\theta}_i | i \in \mathbf{x}^{(k)}\}$ are also applied to update the local task network weights, i.e.,

$$\boldsymbol{\theta} \leftarrow \frac{\boldsymbol{\theta} + \sum_{i \in \mathbf{x}^{(k)}} \omega_i \hat{\theta}_i}{1 + \sum_{i \in \mathbf{x}^{(k)}} \omega_i} \quad (10)$$

Upon receiving the task reward $r^{(k)}$, the communication reward $\tilde{r}^{(k)}$ is evaluated based on the task reward $r^{(k)}$ and the number of the selected cooperative agents $\xi^{(k)}$. The communication experience is formulated based on the communication state sequence $\tilde{\mathbf{s}}^{(k)}$, communication policy $\mathbf{x}^{(k)}$, communication reward $\tilde{r}^{(k)}$ and next communication state $\tilde{\mathbf{s}}^{(k+1)}$, i.e.,

$$\mathbf{z}^C = \left\{ \tilde{\mathbf{s}}^{(k)}, \mathbf{x}^{(k)}, \tilde{r}^{(k)}, \tilde{\mathbf{s}}^{(k+1)} \right\}, \quad (11)$$

which is stored in the memory pool \mathcal{D}^C to update network weights $\alpha_{1 \leq j \leq M}$ in the communication layer. By sampling G experiences from \mathcal{D}^C to formulate the minibatch \mathcal{B}^C , the neural network weight α_j of level j is updated by minimizing the loss function similar to (9).

VI. PERFORMANCE ANALYSIS

We formulate the interactions among the learning agents as a cooperative MARL communication game and analyze the performance bounds including the information gain from the learning agent cooperation, the communication cost and utility based on the Nash equilibrium of the game.

A. Game Model

In the cooperative MARL communication game, N learning agents $\mathcal{N} = \{1, 2, \dots, N\}$ as players communicate with each other to share learning information. Specifically, learning agent i chooses whether to communicate with neighboring agents \mathcal{N}_{-i} , i.e., $\mathbf{x}_i \in \mathcal{X}_i = \{0, 1\}^{N-1}$, to request learning information denoted by $\{\mathcal{M}_j | j \in \mathcal{N}_{-i}\}$ including the observations and learning parameters and chooses the RL action $\mathbf{a}_i \in \mathcal{A}_i$ according to state $s_i \in \mathcal{S}$. The goal of learning agent $i \in \mathcal{N}$ is to maximize the communication utility u_i as a weighted sum of the RL performance and the communication cost ξ_i , i.e., the number of the selected cooperative agents. Thus, the cooperative communication game model is formulated as

$$\mathcal{G} = \left\{ \mathcal{N}, \mathcal{S}, \{\mathcal{X}_i\}_{i \in \mathcal{N}}, \{\mathcal{A}_i\}_{i \in \mathcal{N}}, \{u_i\}_{i \in \mathcal{N}} \right\} \quad (12)$$

According to [33], a profile of strategies denoted by the $N \times (N - 1)$ matrix $\mathbf{X}^* = [\mathbf{x}_i^*]_{1 \leq i \leq N}$ is a Nash equilibrium of the cooperative N -player game \mathcal{G} , where $\forall i \in \mathcal{N}$, \mathbf{x}_i^* is the best response to the strategies of the other $N - 1$ agents denoted by \mathbf{X}_{-i}^* . That is, $\forall \mathbf{x}_i \in \{0, 1\}^{N-1}$, we have

$$u_i(\mathbf{x}_i, \mathbf{X}_{-i}^*) \leq u_i(\mathbf{x}_i^*, \mathbf{X}_{-i}^*), \forall i \in \mathcal{N} \quad (13)$$

indicating that agents have no incentive to deviate from the equilibrium \mathbf{X}^* .

B. Performance Bound

The performance bound of the MARL communication scheme is derived based on the Nash equilibrium of the cooperative communication game in (12). Without loss of generality, the communication among N learning agents is assumed to follow the informational model as proposed in [7], which analyzes the effect of the learning experience exchange on the RL performance. Specifically, the RL performance of learning agent i at each time slot increases with the total information gain obtained from the cooperative communication process denoted by \mathcal{I}_i , which consists of the local information gain denoted by $\mathcal{H}_{i,i}$ in local state s_i and the shared information gain from the neighboring agent $j \in \mathcal{N}_{-i}$ denoted by $\mathcal{H}_{i,j}$, and the information loss denoted by $\mathcal{L}_{i,j}$ in the experience sharing caused by the changing policies of neighboring agents. For the sake of convenience, the communication policy \mathbf{x}_i is transformed to N -dimension vector that includes $x_{i,i} = 1$ in the i -th variable, which has no influence on the performance analysis. Thus, the total information gain \mathcal{I}_i of learning agent i is given by

$$\mathcal{I}_i = \sum_{j=1}^N x_{i,j} \mathcal{H}_{i,j} + (x_{i,i} - 1) \mathcal{L}_{i,j} \quad (14)$$

For the learning agent i , the information gain $\mathcal{H}_{i,j}$ from the learning message \mathcal{M}_j shared by neighboring agent j depends on the previous information gain from agent j denoted by $\hat{\mathcal{H}}_{i,j}$, current information fraction of agent j in the total environment information of the learning agent denoted by $C_{i,j}$, the learning function $\Lambda(\cdot)$ on the learning factor $\omega_{i,j}$ and the state s_j , and the information transfer function $\Phi(\cdot)$ with property $\Phi(y) \leq y$ due to the information loss. Thus, the information gain from agent j is given by

$$\mathcal{H}_{i,j} = \Lambda(\omega_{i,j}, s_j) \Phi(C_{i,j} - \hat{\mathcal{H}}_{i,j}) \quad (15)$$

Similarly, the information gain from the local state s_i is $\mathcal{H}_{i,i} = \Lambda(\omega_{i,i}, s_i) \Phi(C_{i,i} - \hat{\mathcal{H}}_{i,i})$.

The RL performance of learning agent i also suffers from an information loss at the agents that are chosen not to communicate with in current time slot, which results in outdated knowledge from previous shared learning messages. Specifically, the information loss $\mathcal{L}_{i,j}$ of the learning agent i at neighboring agent j with $x_{i,j} = 0$ depends on the information gain of agent j from the environment at current time slot denoted by \mathcal{T}_j , previous total information gain at agent j denoted by $\hat{\mathcal{T}}_j$, and previous shared information gain from agent j denoted by $\hat{\mathcal{H}}_{i,j}$, i.e.,

$$\mathcal{L}_{i,j} = \frac{\hat{\mathcal{H}}_{i,j} \mathcal{T}_j}{\hat{\mathcal{T}}_j + \mathcal{T}_j} \quad (16)$$

By (14), (15) and (16), the communication utility u_i that increases with the total information gain \mathcal{I}_i and decreases with

the communication cost ξ_i is given by

$$u_i = \sum_{j \neq i}^N x_{i,j} \Lambda(\omega_{i,j}, s_j) \Phi(C_{i,j} - \hat{\mathcal{H}}_{i,j}) - \frac{(1 - x_{i,j}) \hat{\mathcal{H}}_{i,j} \mathcal{T}_j}{\hat{\mathcal{I}}_j + \mathcal{T}_j} - \varpi^C x_{i,j} + \Lambda(\omega_{i,i}, s_i) \Phi(C_{i,i} - \hat{\mathcal{H}}_{i,i}) \quad (17)$$

For simplicity, the learning agent is assumed to have no information loss in the learning process, i.e., $\Phi(C_{i,j} - \hat{\mathcal{H}}_{i,j}) = C_{i,j} - \hat{\mathcal{H}}_{i,j}, \forall j \in \mathcal{N}$. The shared information from neighboring agents is assumed to have positive effect on the learning performance of learning agent i , which increases with the task similarity $\eta_{i,j}$, i.e., $\Lambda(\omega_{i,j}, s_j) = \omega_{i,j} \eta_{i,j}, \forall j \in \mathcal{N}_{-i}$. Thus, by (17), the communication utility u_i is given by

$$u_i = \sum_{j \neq i}^N x_{i,j} \omega_{i,j} \eta_{i,j} (C_{i,j} - \hat{\mathcal{H}}_{i,j}) - \frac{(1 - x_{i,j}) \hat{\mathcal{H}}_{i,j} \mathcal{T}_j}{\hat{\mathcal{I}}_j + \mathcal{T}_j} - \varpi^C x_{i,j} (\mathcal{C}_{i,i} - \hat{\mathcal{H}}_{i,i}) \quad (18)$$

Theorem 1: The performance bound of MARL communication schemes in terms of the information gain \mathcal{I}_i , communication cost ξ_i and utility $u_i, \forall i \in \mathcal{N}$, is given by

$$\bar{\mathcal{I}}_i = \bar{\omega} \bar{\eta} \sum_{j=1, j \neq i}^{M_i} (C_{i,j} - \hat{\mathcal{H}}_{i,j}) - \sum_{j=M_i+1, j \neq i}^N \frac{\hat{\mathcal{H}}_{i,j} \mathcal{T}_j}{\hat{\mathcal{I}}_{i,j} + \mathcal{T}_j} + \bar{\omega} (\mathcal{C}_{i,i} - \hat{\mathcal{H}}_{i,i}) \quad (19a)$$

$$\bar{\xi}_i = M_i \quad (19b)$$

$$\bar{u}_i = \bar{\omega} \bar{\eta} \sum_{j=1, j \neq i}^{M_i} (C_{i,j} - \hat{\mathcal{H}}_{i,j}) - \sum_{j=M_i+1, j \neq i}^N \frac{\hat{\mathcal{H}}_{i,j} \mathcal{T}_j}{\hat{\mathcal{I}}_{i,j} + \mathcal{T}_j} + \bar{\omega} (\mathcal{C}_{i,i} - \hat{\mathcal{H}}_{i,i}) - \varpi^C M_i, \quad (19c)$$

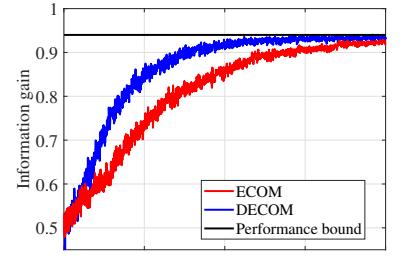
if the $N \geq M_i$ learning agents have identical learning factor $\bar{\omega}$ and task similarity $\bar{\eta}$, and $\forall j \in \mathcal{N}_{-i}$, have

$$\bar{\eta} > \frac{\varpi^C - \hat{\mathcal{H}}_{i,j}}{\bar{\omega} (C_{i,j} - \hat{\mathcal{H}}_{i,j})} \quad (20)$$

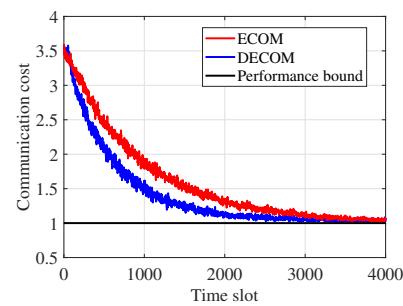
$$\begin{aligned} \frac{\max_{M_i < m \leq N} \mathcal{T}_m}{\min_{1 \leq m \leq N} \hat{\mathcal{I}}_m} &\leq \frac{\varpi^C - \bar{\omega} \bar{\eta} (C_{i,j} - \hat{\mathcal{H}}_{i,j})}{\hat{\mathcal{H}}_{i,j} + \bar{\omega} \bar{\eta} (C_{i,j} - \hat{\mathcal{H}}_{i,j}) - \varpi^C} \\ &\leq \frac{\min_{1 \leq m \leq M_i} \mathcal{T}_m}{\max_{1 \leq m \leq N} \hat{\mathcal{I}}_m} \end{aligned} \quad (21)$$

Proof. See Appendix A. \square

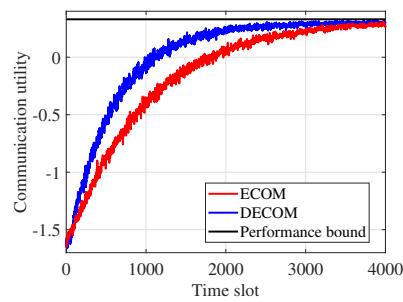
Remark 1: If the task similarity among learning agents is identical and higher than a bound that decreases with the learning factor and information fraction given by (20), learning agent i chooses to communicate with M_i neighboring agents that have more newly learned information from the environment at current time slot as given by (21). In this case, the bound of the information gain obtained from the environment and experience sharing with the neighboring agents given by (19) increases with the task similarity and the information fraction of the selected cooperative agents and decreases with the newly learnt information of the unchosen



(a) Information gain



(b) Communication cost



(c) Communication utility

Fig. 3. Performance of the MARL communication scheme with 8 learning agents, in which each agent chooses whether to communicate with neighboring agents for learning information sharing.

agents $\mathcal{T}_i, M_i + 1 \leq i \leq N$, as given by (21). Our proposed ECOM and DECOM converges to the performance bound provided in (19) after 3000 and 4000 time slots, respectively as shown in Fig. 3. For example, DECOM increases 84.0% information gain after convergence, which is about 1.1% smaller than the bound given by (19a).

The communication energy consumption increases with the number of the selected cooperative agents and their transmit power to share the learning information. The communication cost of our proposed scheme decreases over time due to the optimization of the cooperative agents, as shown in Fig. 3. In addition, with the information gain obtained from the learning information exchange increasing, the communication cost will decrease to make a trade-off between learning and communication performance. Therefore, the total communication cost in the learning process increases with the network size N and the learning time K before convergence.

VII. CASE STUDY: MARL-BASED UAV SWARM VIDEO TRANSMISSION AGAINST JAMMING

As a case study, the proposed communication scheme can be implemented in the MARL-based anti-jamming UAV swarm transmission to improve the task performance such as PLR and energy efficiency with QoS guarantee. This work is also applicable to other MARL-based mobile applications. More specifically, at time slot k , each UAV in the swarm as the learning agent is assumed to have N neighboring UAVs and chooses the transmit channel $f^{(k)} \in \{1, 2, \dots, F\}$ and transmit power $p^{(k)} \in [P_{\min}, P_{\max}]$ to send $b^{(k)}$ -byte images or videos to the control center and support mobile applications such as traffic surveillance or target tracking. A reactive jammer applies the signal detection techniques to sense the ongoing swarm transmission and chooses the jamming channel and power to degrade the swarm transmission performance.

As an example, the RL algorithms including Q-learning and DQN are applied to choose the task action, i.e., anti-jamming transmission policy $\mathbf{a}^{(k)} = [p^{(k)}, f^{(k)}]$ from action set $\mathcal{A} = \{1, 2, \dots, F\} \times \{jP_{\max}/J | 1 \leq j \leq J\}$ to improve the video transmission quality with less energy consumption. The proposed MARL communication scheme enables the UAV to choose the cooperative UAVs $\mathbf{x}^{(k)}$ for cooperative information sharing of policy learning and formulate the learning state $s^{(k)}$ of anti-jamming policy selection to improve the transmission performance under partial observation. The shared learning parameters of cooperative UAVs are also exploited to update Q-values and neural network weights to further improve the policy optimization efficiency. This work is also applicable to other RL algorithms.

A. ECOM-enhanced UAV Swarm Transmission

At time slot k , UAV as the learning agent evaluates the message size and obtains the signal SINR received at control center ρ_1 , the previous transmission delay ρ_2 , the packet PLR ρ_3 and the received jamming power ς from the feedback message sent from the control center. The local observation $\mathbf{o}^{(k)}$ consists of the size of transmitted message b , the signal SINR received at control center ρ_1 , the previous transmission delay ρ_2 , the PLR ρ_3 and the received jamming power at control center ς , i.e.,

$$\mathbf{o}^{(k)} = [b, \varsigma, \rho] \quad (22)$$

Based on the local observation $\mathbf{o}^{(k)}$, the channel state \mathbf{h} and task similarity η with neighboring agents, previous W observations of neighboring UAVs $\bar{\mathbf{o}}$ and the number of selected agents in the last time slot, the communication learning state $\tilde{\mathbf{s}}^{(k)}$ is formulated as

$$\tilde{\mathbf{s}}^{(k)} = [b, \varsigma, \rho, \mathbf{h}, \eta, \bar{\mathbf{o}}], \quad (23)$$

where each element \tilde{s}_i is quantized into \tilde{q}_i levels, $1 \leq i \leq 5 + 5W + 2N$. With the quantized $\tilde{\mathbf{s}}^{(k)}$ as the input, the M Q-tables in the sequential architecture choose cooperative agents $\mathbf{x}^{(k)}$ based on $Q_{1 \leq j \leq M}^C(\tilde{\mathbf{s}}^{(k)}, \cdot)$ via (2) to request the learning information including the observation $\hat{\mathbf{o}}_i$ and learning parameters \hat{Q}_i^R of each cooperative agent $i \in \mathbf{x}^{(k)}$.

The state $s^{(k)}$ of task learning, i.e., the anti-jamming transmission policy selection, consists of the local observation $\mathbf{o}^{(k)}$ in (22) and the encoded observations $\tilde{\mathbf{o}}_i$ of each cooperative agent $i \in \mathbf{x}^{(k)}$ based on the observation encoder, i.e.,

$$\mathbf{s}^{(k)} = [b, \varsigma, [\rho_{1 \leq i \leq 3}], [\tilde{\mathbf{o}}_{i \in \mathbf{x}^{(k)}}]] \quad (24)$$

Each element s_i in the anti-jamming policy learning state $\mathbf{s}^{(k)}$ is quantized into q_i levels, $1 \leq i \leq 5 + M$. The anti-jamming transmission policy $\mathbf{a}^{(k)}$ consists of the transmit channel $f^{(k)} \in \{1, 2, \dots, F\}$ and transmit power $p^{(k)} \in \{jP_{\max}/J | 1 \leq j \leq J\}$ and is chosen based on the task Q-values $Q^R(s^{(k)}, \cdot)$ according to the ϵ -greedy method.

Upon receiving the feedback from the control center, the task reward $r^{(k)}$ is evaluated based on signal SINR ρ_1 , transmission delay ρ_2 , packet PLR ρ_3 , UAV transmit power $p^{(k)}$ and the frequency hopping cost ϱ , given by

$$r^{(k)} = \rho_1 - \varpi_1^T \rho_2 - \varpi_2^T \rho_3 - \varpi_3^T p^{(k)} - \varrho \mathbb{I}(f^{(k)} \neq f^{(k-1)}) \quad (25)$$

Based on the task reward $r^{(k)}$ and the number of selected cooperative agents ξ , the communication reward $\tilde{r}^{(k)}$ is evaluated as (7). The task Q-values are updated with $r^{(k)}$ and shared Q-values via (5) and (6). The communication Q-values are also updated with reward $\tilde{r}^{(k)}$ similar to (5).

B. DECOM-enhanced UAV Swarm Transmission

The proposed communication scheme for deep MARL is implemented in the DQN-based anti-jamming UAV swarm transmission to further enhance the performance for the high-dimension anti-jamming state space and improve the policy exploration and learning speed with shared learning observations and neural network weights from the cooperative agents. With the communication state $\tilde{\mathbf{s}}^{(k)}$ in (23) as the input, each level j with same network architecture, i.e., four fully-connected layers including a $(5+5W+2N)$ -node input layer, two hidden layers with g_1 and g_2 nodes and an $(N+1)$ -node output layer, estimates the communication Q-values $Q_j^C(\tilde{\mathbf{s}}^{(k)}, \cdot; \alpha_j)$ to choose the cooperative UAV x_j , $j \in \{1, 2, \dots, M\}$ in the sequential architecture.

After receiving the learning response including the observations $\{\hat{\mathbf{o}}_i | i \in \mathbf{x}^{(k)}\}$ and neural network weights of DQN $\{\hat{\theta}_i | i \in \mathbf{x}^{(k)}\}$ from cooperative agents $\mathbf{x}^{(k)}$, the anti-jamming policy transmission selection state $s^{(k)}$ is formulated as (24) and input to the Q-network with four fully-connected layers, each having $5 + M$, ν_1 , ν_2 and FJ nodes, respectively, to obtain the Q-values $Q^R(s^{(k)}, \cdot; \theta)$ and choose the transmit channel and power $\mathbf{a}^{(k)} = [p^{(k)}, f^{(k)}]$.

Upon receiving the feedback from the control center, the anti-jamming transmission reward $r^{(k)}$ and communication reward $\tilde{r}^{(k)}$ is evaluated via (25) and (7), respectively. The anti-jamming transmission experience \mathbf{z}^R and communication experience \mathbf{z}^C are stored in the memory pool $\mathcal{D} \leftarrow \mathcal{D} \cup \{\mathbf{z}^R, \mathbf{z}^C\}$. The experience replay technique is used to update network weights with the stochastic gradient descent algorithm by sampling G experiences $\tilde{\mathbf{z}} = \{\mathbf{z}^R, \mathbf{z}^C\}$ to formulate the minibatch \mathcal{B} . The network weights are updated by minimizing

the loss function in (9) and also updated based on the shared network weights from the cooperative agents via (10).

C. Computational Complexity

The computational complexities of the proposed ECOM and DECOM increase with the number of multiplications in the cooperative agent and task policy selection and the update of the reinforcement learning layer and the communication learning layer. According to [34], the computational complexity of ECOM in the cooperative agent selection and task policy selection is $O(MN)$ and $O(FJ)$, respectively. The proposed ECOM updates the communication Q-values with computational complexity $O(MN)$ and updates task Q-values based on task reward and shared Q-values with computational complexity $O(FJ \prod_{i=1}^{5+M} q_i)$, where q_i is the quantized level of i -th element in task state. Therefore, the computational complexity of ECOM with a large anti-jamming policy space is $O(FJ \prod_{i=1}^{5+M} q_i)$.

In DECOM, M C-Networks in the sequential architecture of the communication learning layer and the task network in the reinforcement learning layer perform forward propagation and back propagation to choose the cooperative agents and anti-jamming policy and update network weights $\alpha_{1 \leq j \leq M}$ and θ . The M C-Networks in the communication learning layer perform $\phi_{C,F}$ multiplications in the forward propagation of the cooperative agent selection that linearly increases with the number of network weights in four fully-connected layers, i.e.,

$$\begin{aligned} \phi_{C,F} = & M(5 + 5W + 2N)g_1 + M(g_1 + 1)g_2 \\ & + M(g_2 + 1)(N + 1) \end{aligned} \quad (26)$$

The observation encoder performs $\phi_{E,F} = 11V$ multiplications in the forward propagation to extract correlation between local and shared observations. The task network performs forward propagation to choose the anti-jamming policy and has $\phi_{T,F}$ multiplications given by

$$\phi_{T,F} = (M + 5)(\nu_1 + 1) + \nu_1(\nu_2 + 1) + (\nu_2 + 1)FJ \quad (27)$$

The C-Networks and task network perform forward propagation to estimate the target value in the loss function to update the network weights $\alpha_{1 \leq j \leq M}$ and θ with $2\phi_{C,F}$ and $2\phi_{T,F}$ multiplications, respectively. To calculate the gradient the loss function and update neural network weights, M C-Networks perform $\phi_{C,B}$ multiplications in the back propagation given by

$$\begin{aligned} \phi_{C,B} = & 2M(5 + 5W + 2N)g_1 + 2M(g_1 + 1)g_2 \\ & + 3M(g_2 + 1)(N + 1) \end{aligned} \quad (28)$$

Similarly, the number of the multiplications in the back propagation of task network is given by

$$\begin{aligned} \phi_{T,B} = & 2(M + 5)(\nu_1 + 1) + 2\nu_1(\nu_2 + 1) \\ & + 3(\nu_2 + 1)FJ \end{aligned} \quad (29)$$

In addition, the task network also takes $M\phi_{T,F}$ multiplications to update local neural network weights based on the shared information. Thus, the total number of multiplications of

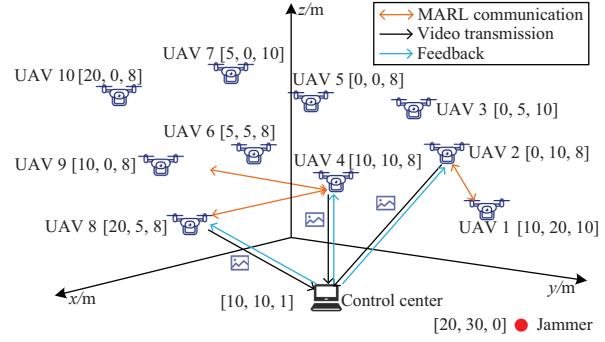


Fig. 4. Simulation setting of the swarm anti-jamming video transmission, in which each UAV sends 2-Mb videos or images to the control center with up to 100 mW transmit power and chooses the cooperative agents from the neighboring UAVs for observation and learning information sharing to accelerate anti-jamming transmission policy learning.

DECOM with G -experience minibatch per time slot is

$$\begin{aligned} \psi = & (2G + 1)\phi_{C,F} + (2G + M + 1)\phi_{T,F} + \phi_{E,F} \\ & + G(\phi_{C,B} + \phi_{T,B}) \end{aligned} \quad (30)$$

Theorem 2: The computational complexity of the ECOM and DECOM-based anti-jamming UAV swarm transmission scheme is given by $O(FJ \prod_{i=1}^{5+M} q_i)$ and $\mathcal{O}(GKFJ)$, respectively.

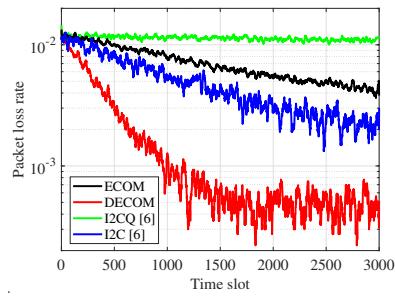
Proof. See Appendix B. \square

Remark 2: The computational complexity of ECOM increases with quantized level of observation $q_{1 \leq i \leq M+5}$ and transmit power level J as well as the number of selected cooperative agents M that determines the task state size. The computational complexity of DECOM depends on the number of transmit channels F and the number of quantized power level J that determine the task action space size as well as the batch size G and the number of learning samples K .

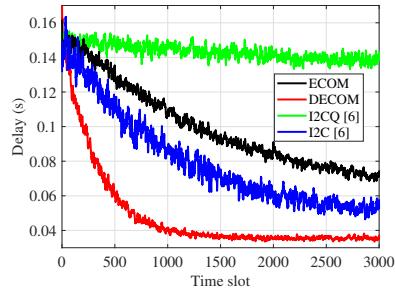
VIII. SIMULATION RESULTS

The performance is evaluated in the DQN-based UAV swarm transmission against jamming attacks as shown in Fig. 4 compared with benchmark I2C in [6]. Each UAV sends 2-Mb video data to the control center modulated with quadrature phase-shift keying and chooses the transmit power with up to 100 mW on 13 channels at 2.4 GHz each with the bandwidth of 10 MHz according to 802.11n [35], [36]. The jammer sends Gaussian signals with jamming power chosen from $\{20j | 0 \leq j \leq 5\}$ mW at 2.4 GHz with jamming bandwidth of 30 MHz to block 3 UAV channels of swarm transmission according to the jamming model in [37].

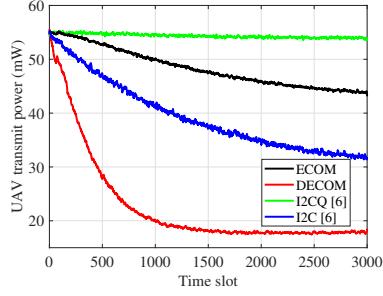
The C-Network in the communication learning layer and the task network in the task learning layer are instantiated as the four-layer neural networks including the input and output layer and two hidden layers each with 128 neural nodes. The UAV transmit power level $J = 10$, the learning rate of the Q-value update $\lambda = 0.4$, the discount factor of reward $\gamma = 0.3$, and the greedy parameter ϵ decays from 1 to 0.01 after 1000 time slots



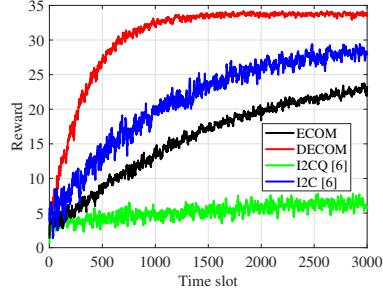
(a) Packet loss rate



(b) Delay



(c) UAV transmit power

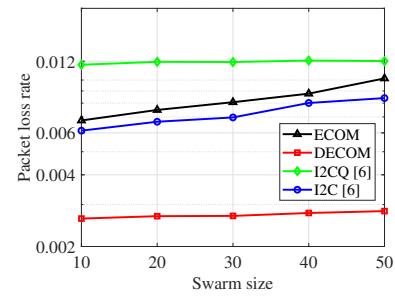


(d) Reward

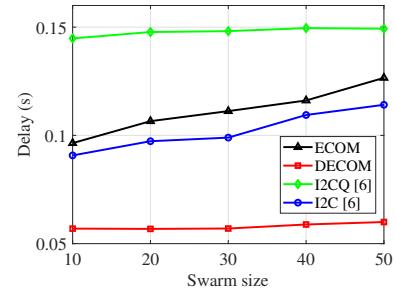
Fig. 5. Performance of the 10-UAV swarm video transmission against a jammer as shown in Fig. 4.

to make a trade-off between exploration and exploitation. In addition, the learning rate of Adma optimizer for the update of neural network weights $\beta = 0.01$ and the number of sampled experiences in the minibatch $G = 64$.

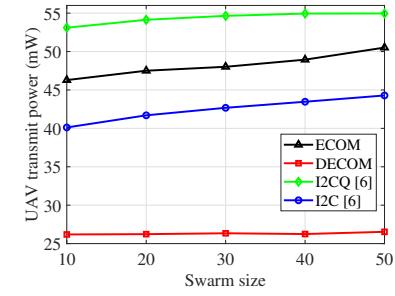
Compared with the benchmark MARL communication scheme I2C in [6], our proposed communication schemes accelerate the learning speed and improve the RL-based anti-jamming transmission performance as shown in Fig. 5. For example, our proposed ECOM reduces 62.5% PLR to



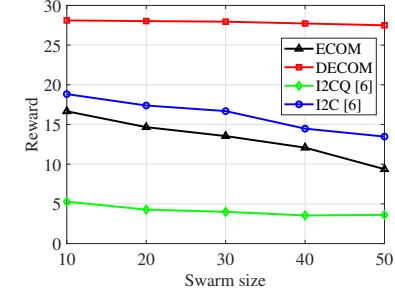
(a) Packet loss rate



(b) Delay



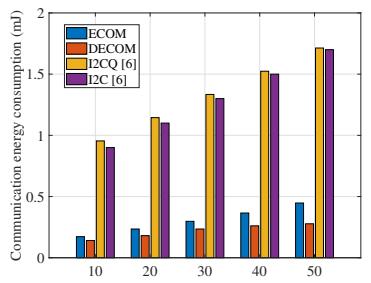
(c) UAV transmit power



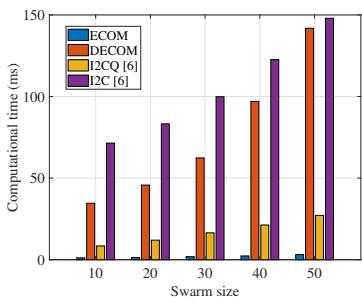
(d) Reward

Fig. 6. Average performance of the swarm video transmission against jamming attack with different swarm sizes.

4.2×10^{-3} with 18.9% less transmit power compared with I2C-based Q-learning algorithm named I2CQ after 3000 time slots. The reason is that learning based communication mechanism enables learning agent to choose the most relevant cooperative agent for observation sharing and thus reduces the information redundancy in the task policy selection. Our proposed communication scheme for deep MARL further improves the performance with 80.1% reduced PLR to 4.3×10^{-4} and 43.8% less UAV transmit power after 2000 time slots and



(a) Communication energy consumption



(b) Computational time

Fig. 7. MARL communication energy consumption and computational time for the UAV swarm ranging from 10 to 50 averaged over 2000 time slots.

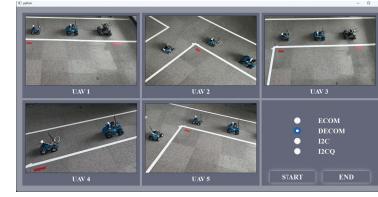
saves 44.0% learning time compared with I2C-based DQN algorithm, because that the neural network in the communication layer accelerates the cooperative UAV selection and the shared network weights further improves the anti-jamming policy optimization.

The average performance as shown in Fig. 6 verifies the performance gain and the robustness of our proposed schemes for the swarm ranging between 10 and 50 UAVs. For example, compared with benchmark I2CQ, our proposed ECOM decreases 32.2% PLR and 25.0% transmission delay for the 30-UAV swarm, due to the adaptive cooperative agent selection that results in more efficient anti-jamming policy learning. DECOM further reduces 61.4% PLR, 42.4% transmission delay and saves 38.4% UAV transmit power and the performance gain increases with the swarm size to 66.6%, 47.4% and 40.2%, respectively, because the neural network based sequential architecture accelerates the cooperative agent selection under the large swarm size.

As shown in Fig. 7, our proposed schemes also improve the multi-agent communication efficiency and decrease the computational time of the communication and task learning under different swarm size compared with benchmark I2C. For example, our proposed ECOM reduces larger than 70.4% MARL communication energy consumption, which is less than 0.5 mJ for the UAV swarms ranging from 10 to 50, and 86.6% computational time due to the adaptive cooperative agent selection that reduces the information redundancy in the RL state formulation. In addition, our proposed DECOM reduces larger than 83.7% communication energy communication compared with I2C due to the neural networks that

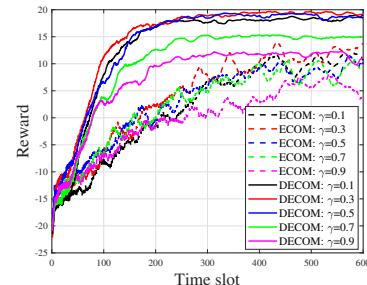


(a) Experiment setting

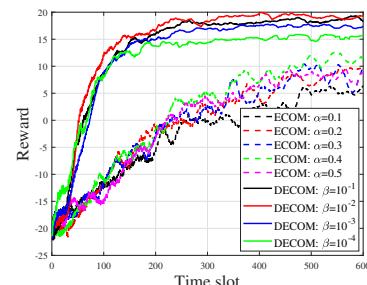


(b) Graphical user interface

Fig. 8. Experimental setting and graphical user interface of the UAV swarm transmission against jamming, in which each of 5 UAVs as the learning agent sends captured videos of the sensing area to the control center with transmit power up to 100 mW.



(a) Reward vs. discount factor



(b) Reward vs. learning rate

Fig. 9. Convergence performance of the anti-jamming video transmission of the proposed schemes for 5 UAVs as shown in Fig. 8 versus discount factor and learning rate averaged over 100 runs.

accelerate the cooperative agent learning.

IX. EXPERIMENTAL RESULTS

Experiments were performed to evaluate the performance gain of the proposed communication scheme in the 5-UAV swarm anti-jamming video transmission with topology as shown in Fig. 8. Each UAV equipped with Raspberry Pi sends 1-Mb captured videos of the sensing area to the control center

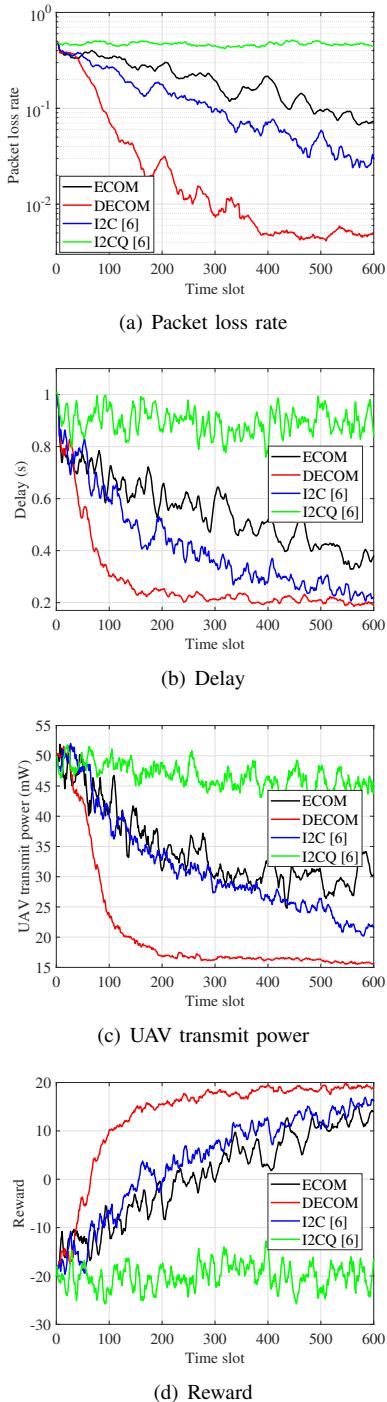


Fig. 10. Performance of the swarm video transmission to the control center against jamming attack with topology as shown in Fig. 8, in which each UAV chooses a cooperative agent from the other 4 UAVs to share observations for faster learning.

per second with up to 100 mW transmit power from one of the 6 channels at 2.4 GHz each with bandwidth 20 MHz according to IEEE 802.11n [35]. A laptop with an Intel i7-11800H 8 cores at 2.3 GHz CPU as the control center can switch the MARL-based anti-jamming transmission policy learning algorithms and display the real-time transmission performance in the graphical user interface. The jammer, i.e., a USRP N210

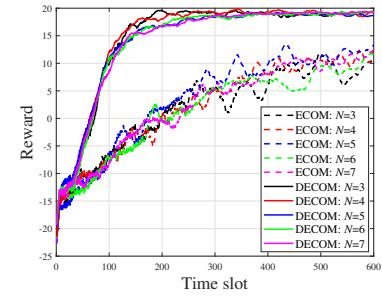


Fig. 11. Convergence performance of the anti-jamming video transmission for the swarm with $N \in [3, 7]$ UAVs averaged over 100 runs in the experimental setting as shown in Fig. 8.

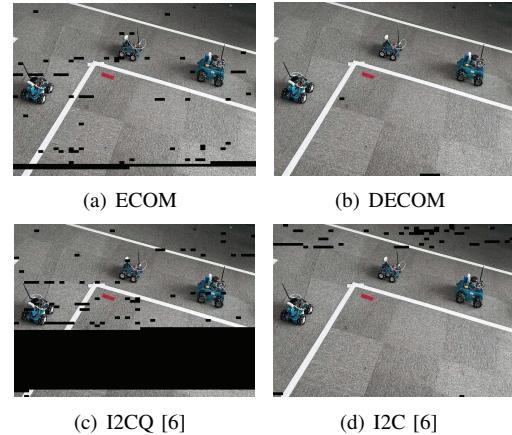


Fig. 12. Anti-jamming transmission performance of UAV 5 to the control center after 500 time slots as shown in Fig. 8.

controlled by a laptop, sends Gaussian signals with up to 100 mW jamming power at center frequency 2.412 GHz with bandwidth 20 MHz.

As shown in Fig. 9, experimental results of the UAV swarm anti-jamming video transmission with topology in Fig. 8 are provided to show the effect of the discount factor and learning rate on the convergence of our proposed ECOM and DECOM. The convergence performance averaged over 100 runs under different discount factor $0.1 \leq \gamma \leq 0.9$ shows that the discount factor $\gamma = 0.3$ makes a tradeoff between the exploration and exploitation in the learning process. In addition, the learning rate $\lambda = 0.4$ of the Q-value update and the learning rate $\beta = 0.01$ in the Adam optimizer of the neural network weight update have been shown to improve the anti-jamming video transmission performance.

The performance gain over the benchmark scheme I2C in [6] is also verified via experimental results as shown in Fig. 10. For example, our proposed ECOM decreases 84.5% PLR, 55.4% transmission delay and 33.5% UAV transmit power after 500 time slots compared with I2CQ. The reason is that the adaptive cooperative agent selection chooses the agent with similar learning task in the dynamic swarm compared with the pre-trained prior network. In addition, the neural network based sequential architecture in DECOM further accelerates the communication policy learning and improves the video

transmission performance with 81.4% PLR, 18.7% transmission delay and 26.8% transmit power reduced compared with the benchmark I2C. Experimental results as shown in Fig. 11 also show that our proposed schemes are robust against the swarm size ranging from 3 to 7. In addition, both our proposed ECOM and DECOM improve the anti-jamming transmission quality with lower packet PLR compared with I2C as shown in Fig. 12 due to the MARL with enhanced communication scheme that accelerates the UAV anti-jamming transmission policy learning.

X. CONCLUSION

In this paper, we have proposed an efficient communication scheme to choose the cooperative agents to share the environment observations and learning experiences for MARL-based mobile applications. The correlation with the observations of the neighboring learning agents extracted based on the attention mechanism is exploited in the task state formulation to reduce the redundancy of the shared observations and enhance the RL policy exploration speed. Both the shared Q-values of state-action pairs and neural network weights are used in the policy distribution and task network update to accelerate RL task policy learning. The performance bound regarding the utility of the learning agent increases with the task similarity and the shared information of the selected agents. The communication enhanced multi-agent RL has been implemented in the UAV swarm video transmission to choose the transmit channel and power against jamming. Both the simulation and experimental results show the performance gain over the benchmark I2C, e.g., with PLR, transmission delay, transmit power and learning time reduced by 81.4%, 18.7%, 26.8% and 50.0%, respectively, in the experimental results of the 5-UAV swarm.

APPENDIX A PROOF OF THEOREM 1

Proof: By (18), $\forall x_{i,j} \in \{0, 1\}$ and $j \in \{1, 2, \dots, M_i\}$, we have

$$\begin{aligned} & u_i([x_{i,j}, \mathbf{x}_{i,-j}^*], \mathbf{X}_{-i}^*) \\ &= y_{i,j} x_{i,j} + \omega_0 (C_{i,i} - \hat{\mathcal{H}}_{i,i}) - \frac{\hat{\mathcal{H}}_{i,j} \mathcal{T}_j}{\hat{\mathcal{I}}_j + \mathcal{T}_j} - \varpi^c (M_i - 1) \\ &+ \sum_{j \neq i}^{M_i} \bar{\omega} \bar{\eta} (C_{i,j} - \hat{\mathcal{H}}_{i,j}) - \sum_{j=M_i+1}^N \frac{\hat{\mathcal{H}}_{i,j} \mathcal{T}_j}{\hat{\mathcal{I}}_j + \mathcal{T}_j}, \end{aligned} \quad (31)$$

where

$$\mathbf{x}_i^* = [\underbrace{1, 1, \dots, 1}_{M_i}, \underbrace{0, 0, \dots, 0}_{N-M_i}] \quad (32)$$

$$y_{i,j} = \bar{\omega} \bar{\eta} (C_{i,j} - \hat{\mathcal{H}}_{i,j}) + \frac{\hat{\mathcal{H}}_{i,j} \mathcal{T}_j}{\hat{\mathcal{I}}_j + \mathcal{T}_j} - \varpi^c \quad (33)$$

By (33), (20) and (21), we have

$$y_{i,j} = \frac{(\bar{\omega} \bar{\eta} (C_{i,j} - \hat{\mathcal{H}}_{i,j}) - \varpi^c) (\hat{\mathcal{I}}_j + \mathcal{T}_j) + \hat{\mathcal{H}}_{i,j} \mathcal{T}_j}{\hat{\mathcal{I}}_j + \mathcal{T}_j} \geq 0 \quad (34)$$

Thus, by (31) and (34), for $x_{i,j} \in \{0, 1\}$ and $j \in \{1, 2, \dots, M_i\}$, we have

$$u_i([x_{i,j}, \mathbf{x}_{i,-j}^*], \mathbf{X}_{-i}^*) \leq u_i(\mathbf{x}_i^*, \mathbf{X}_{-i}^*) \quad (35)$$

Similarly, by (18) and (21), for $x_{i,j} \in \{0, 1\}$ and $\forall j \in \{M_i + 1, \dots, N\}$, we have $y_{i,j} \leq 0$ and

$$u_i([x_{i,j}, \mathbf{x}_{i,-j}^*], \mathbf{X}_{-i}^*) \leq u_i(\mathbf{x}_i^*, \mathbf{X}_{-i}^*) \quad (36)$$

By (35) and (36), $\forall \mathbf{x}_i \in \{0, 1\}^N$, we have

$$u_i(\mathbf{x}_i, \mathbf{X}_{-i}^*) \leq u_i(\mathbf{x}_i^*, \mathbf{X}_{-i}^*), \forall i \in \mathcal{N} \quad (37)$$

Thus, by (37), $\mathbf{X}^* = [\mathbf{x}_i^*]_{1 \leq i \leq N}$ with \mathbf{x}_i^* given by (32) is a Nash equilibrium of the game \mathcal{G} and the performance bound of MARL communication scheme is given by (19) according to [33].

APPENDIX B PROOF OF THEOREM 2

Proof: According to [38], the numbers of the nodes in two hidden layers depend on the number of nodes in the output layer and the number of learning samples K . Specifically, the number of nodes in the second hidden layer $g_2 = \sqrt{K(N+1)}$. The number of nodes in the first hidden layer g_1 increases with the g_2 and the number of learning samples K , i.e.,

$$g_1 = \sqrt{K(N+1)} + 2\sqrt{\frac{K}{N+1}} \quad (38)$$

Similarly, the numbers of nodes in the hidden layers of the task network are $\nu_2 = \sqrt{KFJ}$ and $\nu_1 = \nu_2 + 2\sqrt{K/(FJ)}$.

By (26)-(30) and (38), the computational complexity of DECOM is given by

$$\begin{aligned} \Upsilon &= \mathcal{O}((2G+1)\phi_{C,F} + (2G+M+1)\phi_{T,F} \\ &\quad + G(\phi_{C,B} + \phi_{T,B}) + \phi_{E,F}) \end{aligned} \quad (39)$$

$$= \mathcal{O}(G\nu_1\nu_2 + GFJ\nu_2 + Gg_1g_2 + GNg_1 + GNg_2) \quad (40)$$

$$\begin{aligned} &= \mathcal{O}\left(G\sqrt{KFJ}\left(\sqrt{\frac{K}{FJ}} + \sqrt{KFJ} + FJ\right) + GN\sqrt{KN}\right. \\ &\quad \left.+ \left(G\sqrt{\frac{K}{N}} + GN\right)\left(\sqrt{KN} + 2\sqrt{\frac{K}{N}}\right)\right) \end{aligned} \quad (41)$$

$$= \mathcal{O}(GKFJ), \quad (42)$$

where (40) is obtained by (26)-(29), (41) is obtained by (38), and (42) is obtained as $K \gg FJ > N$.

REFERENCES

- [1] H. Peng and X. Shen, "Multi-agent reinforcement learning based resource management in MEC- and UAV-assisted vehicular networks," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 1, pp. 131–141, Jan. 2021.
- [2] N. Zhao, Z. Ye, Y. Pei *et al.*, "Multi-agent deep reinforcement learning for task offloading in UAV-assisted mobile edge computing," *IEEE Trans. Wireless Commun.*, vol. 21, no. 9, pp. 6949–6960, Sep. 2022.
- [3] J. Foerster, I. A. Assael, N. De Freitas *et al.*, "Learning to communicate with deep multi-agent reinforcement learning," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, Barcelona, Spain, Dec. 2016, pp. 2145–2153.
- [4] A. Das, T. Gervet, J. Romoff *et al.*, "TarMAC: Targeted multi-agent communication," in *Proc. Int. Conf. Mach. Learn. (ICML)*, Long Beach, CA, Jun. 2019, pp. 1538–1546.

- [5] N. Naderizadeh, J. J. Sydir, M. Simsek *et al.*, "Resource management in wireless networks via multi-agent deep reinforcement learning," *IEEE Trans. Wireless Commun.*, vol. 20, no. 6, pp. 3507–3523, Jun. 2021.
- [6] Z. Ding, T. Huang, and Z. Lu, "Learning individually inferred communication for multi-agent cooperation," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, Dec. 2020, pp. 22069–22079.
- [7] F. Hu, Y. Deng, and A. Hamid Aghvami, "Scalable multi-agent reinforcement learning for dynamic coordinated multipoint clustering," *IEEE Trans. Commun.*, vol. 71, no. 1, pp. 101–114, Jan. 2023.
- [8] Y. Yu, S. C. Liew, and T. Wang, "Multi-agent deep reinforcement learning multiple access for heterogeneous wireless networks with imperfect channels," *IEEE Trans. Mobile Comput.*, vol. 21, no. 10, pp. 3718–3730, Oct. 2022.
- [9] C. Dai, K. Zhu, and E. Hossain, "Multi-agent deep reinforcement learning for joint decoupled user association and trajectory design in full-duplex multi-UAV networks," *IEEE Trans. Mobile Comput.*, vol. 22, no. 10, pp. 6056–6070, Oct. 2023.
- [10] H. Yang, J. Zhao, K.-Y. Lam *et al.*, "Distributed deep reinforcement learning-based spectrum and power allocation for heterogeneous networks," *IEEE Trans. Wireless Commun.*, vol. 21, no. 9, pp. 6935–6948, Sep. 2022.
- [11] J. Hu, H. Zhang, L. Song *et al.*, "Cooperative Internet of UAVs: Distributed trajectory design by multi-agent deep reinforcement learning," *IEEE Trans. Commun.*, vol. 68, no. 11, pp. 6807–6821, Nov. 2020.
- [12] N. Zhao, Y.-C. Liang, D. Niyato *et al.*, "Deep reinforcement learning for user association and resource allocation in heterogeneous cellular networks," *IEEE Trans. Wireless Commun.*, vol. 18, no. 11, pp. 5141–5152, Nov. 2019.
- [13] A. A. Khan and R. S. Adve, "Centralized and distributed deep reinforcement learning methods for downlink sum-rate optimization," *IEEE Trans. Wireless Commun.*, vol. 19, no. 12, pp. 8410–8426, Dec. 2020.
- [14] Y. Xiao, L. Xiao, K. Wan *et al.*, "Reinforcement learning based energy-efficient collaborative inference for mobile edge computing," *IEEE Trans. Commun.*, vol. 71, no. 2, pp. 864–876, Feb. 2023.
- [15] X. Zhang, H. Zhao, J. Wei *et al.*, "Cooperative trajectory design of multiple UAV base stations with heterogeneous graph neural networks," *IEEE Trans. Wireless Commun.*, vol. 22, no. 3, pp. 1495–1509, Mar. 2023.
- [16] R. Zhong, X. Liu, Y. Liu *et al.*, "Multi-agent reinforcement learning in NOMA-aided UAV networks for cellular offloading," *IEEE Trans. Wireless Commun.*, vol. 21, no. 3, pp. 1498–1512, Mar. 2022.
- [17] X. Liu, G. Chuai, X. Wang *et al.*, "QoE-driven antenna tuning in cellular networks with cooperative multi-agent reinforcement learning," *IEEE Trans. Mobile Comput.*, pp. 1–15, Dec. 2022.
- [18] Z. Lv, L. Xiao, Y. Du *et al.*, "Multi-agent reinforcement learning based UAV swarm communications against jamming," *IEEE Trans. Wireless Commun.*, vol. 22, no. 12, pp. 9063–9075, Apr. 2023.
- [19] S. Sukhbaatar and R. Fergus, "Learning multiagent communication with backpropagation," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, Barcelona, Spain, Dec. 2019, pp. 2252–2260.
- [20] T. Lin, J. Huh, C. Stauffer *et al.*, "Learning to ground multi-agent communication with autoencoders," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, Dec. 2021, pp. 15230–15242.
- [21] W. Kim, J. Park, and Y. Sung, "Communication in multi-agent reinforcement learning: Intention sharing," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, Long Beach, CA, Apr. 2020, pp. 1538–1546.
- [22] J. Jiang and Z. Lu, "Learning attentional communication for multi-agent cooperation," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, Montréal, Canada, Dec. 2018, pp. 2252–2260.
- [23] S. Q. Zhang, Q. Zhang, and J. Lin, "Efficient communication in multi-agent reinforcement learning via variance based control," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, Vancouver, Canada, Dec. 2019, pp. 3235–3244.
- [24] R. Wang, X. He, R. Yu *et al.*, "Learning efficient multi-agent communication: An information bottleneck approach," in *Proc. Int. Conf. Mach. Learn. (ICML)*, Jul. 2020, pp. 9908–9918.
- [25] T.-Y. Tung, S. Kobus, J. P. Roig *et al.*, "Effective communications: A joint learning and communication framework for multi-agent reinforcement learning over noisy channels," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 8, pp. 2590–2603, Aug. 2021.
- [26] H. Mao, Z. Zhang, Z. Xiao *et al.*, "Learning agent communication under limited bandwidth by message pruning," in *Proc. the AAAI Conference on Artificial Intelligence*, New York, NY, Apr. 2020, pp. 5142–5149.
- [27] Y. Wang, J. Xu, Y. Wang *et al.*, "ToM2C: Target-oriented multi-agent communication and cooperation with theory of mind," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, Apr. 2022.
- [28] Z. Lv, Y. Du, Y. Chen *et al.*, "Efficient communications for multi-agent reinforcement learning in wireless networks," in *Proc. IEEE Global Commun. Conf. (Globecom)*, Kuala Lumpur, Malaysia, Dec. 2023, pp. 583–588.
- [29] A. Narayan and T. Y. Leong, "Effects of task similarity on policy transfer with selective exploration in reinforcement learning," in *Proc. Int. Conf. Autonomous Agents and Multiagent Systems*, Montreal, Canada, May 2019, pp. 2132–2134.
- [30] Q. Zhou, Y. Niu, P. Xiang *et al.*, "Intra-domain knowledge reuse assisted reinforcement learning for fast anti-jamming communication," *IEEE Trans. Inf. Forensics Security*, vol. 18, pp. 4707–4720, June 2023.
- [31] A. Zhang, R. T. McAllister, R. Calandra *et al.*, "Learning invariant representations for reinforcement learning without reconstruction," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, Vienna, Austria, May. 2021, pp. 1–12.
- [32] A. Vaswani, N. Shazeer, N. Parmar *et al.*, "Attention is all you need," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, Long Beach, CA, Dec. 2017, pp. 6000–6010.
- [33] K. Zhang, Z. Yang, and T. Başar, "Multi-agent reinforcement learning: A selective overview of theories and algorithms," *Handbook of reinforcement learning and control*, pp. 321–384, Jun. 2021.
- [34] C. Jin, Z. Allen-Zhu, S. Bubeck *et al.*, "Is Q-learning provably efficient?" in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, Montreal, QC, Canada, Dec. 2019, pp. 1–11.
- [35] LATITUDE MATRICE 200 V2 SERIES technical parameters, Accessed: Feb. 15, 2024. [Online]. Available: <https://www.dji.com/cn/matrice-200-series-v2/info/#specs>
- [36] X. Lu, L. Xiao, G. Niu *et al.*, "Safe exploration in wireless security: A safe reinforcement learning algorithm with hierarchical structure," *IEEE Trans. Inf. Forensics Security*, vol. 17, pp. 732–743, 2022.
- [37] M. K. Hanawal, M. J. Abdel-Rahman, and M. Krunz, "Joint adaptation of frequency hopping and transmission rate for anti-jamming wireless systems," *IEEE Trans. Mobile Comput.*, vol. 15, no. 9, pp. 2247–2259, Sep. 2016.
- [38] G.-B. Huang, "Learning capability and storage capacity of two-hidden-layer feedforward networks," *IEEE Trans. Neural Netw.*, vol. 14, no. 2, pp. 274–281, Mar. 2003.



Zefang Lv (Student Member, IEEE) received her B.S. degree in statistics from Shandong University in 2016 and her M.S. degree in applied statistics from North China Electric Power University in 2020. She is currently pursuing a Ph.D. degree with the Department of Informatics and Communication Engineering, Xiamen University. Her research interests include network security, wireless communications and reinforcement learning.



Liang Xiao (Senior Member, IEEE) received the B.S. degree in communication engineering from the Nanjing University of Posts and Telecommunications, China, in 2000, the M.S. degree in electrical engineering from Tsinghua University, China, in 2003, and the Ph.D. degree in electrical engineering from Rutgers University, NJ, USA, in 2009. She was a Visiting Professor with Princeton University, Virginia Tech, and the University of Maryland, College Park. She is currently a Professor with the Department of Informatics and Communication Engineering, Xiamen University, Xiamen, China. She was a recipient of the Best Paper Award for 2016 INFOCOM Big Security WS and 2017 ICC. She has served as an Associate Editor for *IEEE TRANSACTIONS ON INFORMATION FORENSICS AND SECURITY* and a Guest Editor for *IEEE JOURNAL OF SELECTED TOPICS IN SIGNAL PROCESSING*.



Yousong Du received the B.S. degree in Electronic Information Engineering from Beijing University of Chemical Technology, China, in 2021. He is currently pursuing the M.S. degree in the Department of Informatics and Communication Engineering, Xiamen University. His research interests include network security and wireless communications.



Yunjun Zhu received the B.S. degree in Electrical Engineering and Automation from South China University of Technology, China, in 2021. He is currently pursuing a Master's degree in the Department of Informatics and Communication Engineering at Xiamen University, China. His research interests include network security and wireless communications.



Shuai Han (S'11–M'12–SM'17) is currently a full Professor at the Department of Electronics and Communication Engineering, Harbin Institute of Technology. Shuai Han's research interests include wireless communications, satellite IoT and integrated satellite-terrestrial communication networks. Over his academic career, his students and he have contributed in various fields in wireless networks.

He has authored or co-authored over 100 technical papers in major journals and conferences. As PI, he has more than twenty grants on wireless networks

and positioning. He is an associate editor of *IEEE China Communications*, *IEEE ACCESS*, *Journal of Communications and Information Networks (JCIN)*, *Journal of Telemetry, Tracking and Command*, *Journal of Signal Processing*. And has served as guest editor for many IEEE magazines and journals. He has served as a co-chair for technical symposia of international conference, *IEEE GC 2023*, *ICC 2023*, *IEEE GC 2021*, *IEEE GC 2019*, *IEEE ICC 2018*, *IEEE VTC FALL 2016*. He has also served as the TPC Chair for some international conferences, including the *AICON2019* and *MLICOM2018*. He is a member of 2020-2021 R10 Awards & Recognition Committee. Also, he is a senior member of *IEEE Communication Society*, Chair of IEEE BTS Chapter, Vice Chair of *IEEE Harbin ComSoc Chapter*, Vice Chair of *IEEE Harbin VTS Chapter*, and Vice Chair of *IEEE IoT-AHSN TC*.

Shuai Han began his university studies in Communication Engineering in 2000 at the Harbin Institute of Technology. He received his ME and Ph.D. degrees in Information and Communication Engineering from the Harbin Institute of Technology in 2007 and 2011, respectively. And he completed his post-doctoral work in 2012 in Electrical and Computer Engineering at Memorial University of Newfoundland in Canada.



Yong-Jin Liu (Senior Member, IEEE) is a full professor with the Department of Computer Science and Technology, Tsinghua University, China. He received the BEng degree from Tianjin University, China, in 1998, and the PhD degree from the Hong Kong University of Science and Technology, Hong Kong, China, in 2004. His research interests include machine learning, cognitive computation, computer graphics and computer-aided design. For more information, visit <https://cg.cs.tsinghua.edu.cn/people/Yongjin/Yongjin.htm>