# Attention-based Reinforcement Learning for Real-Time UAV Semantic Communication

Won Joon Yun, Byungju Lim, Soyi Jung, Young-Chai Ko, †Jihong Park, Joongheon Kim, and ‡Mehdi Bennis

*Abstract*—In this article, we study the problem of air-to-ground ultra-reliable and low-latency communication (URLLC) for a moving ground user. This is done by controlling multiple unmanned aerial vehicles (UAVs) in real time while avoiding inter-UAV collisions. To this end, we propose a novel multi-agent deep reinforcement learning (MADRL) framework, coined a graph attention exchange network (GAXNet). In GAXNet, each UAV constructs an attention graph locally measuring the level of attention to its neighboring UAVs, while exchanging the attention weights with other UAVs so as to reduce the attention mismatch between them. Simulation results corroborates that GAXNet achieves up to 4.5x higher rewards during training. At execution, without incurring inter-UAV collisions, GAXNet achieves 6.5x lower latency with the target 0.0000001 error rate, compared to a state-of-the-art baseline framework.

## I. INTRODUCTION

Unmanned aerial vehicles (UAVs) can provision agile and mobile network infrastructure [1], [2], and enable ultra-reliable and low-latency communication (URLLC) even under time-varying environments and tasks, such as moving target sites in disaster scenes and transformable assembly lines in smart factories, to mention a few. As opposed to the stationary and fixed nature of 5G URLLC, such non-terrestrial URLLC systems in beyond 5G are time-varying and physically reconfigured, mandating to co-design control and communication in real time [3].

To this end, machine intelligence is envisaged to play a crucial role. Without any central coordination, intelligent agents can promptly react to local environmental changes, thereby reducing latency while saving radio resources [4]. To imbue the intelligence into multiple interactive agents, multi-agent deep reinforcement learning (MADRL) is a promising solution [5], [6], wherein each agent runs deep learning so as to maximize its expected long-term reward by carrying out actions for its given state observations, i.e., decision-makings on state-action policies. Depending on how to train and execute the deep learning models, MADRL is broadly categorized into three types as elaborated next.

First, *centralized MADRL* is the case wherein all agents send their observations to a central entity to build a global policy determining the actions of the entire agents. MADDPG [5], CommNet [7], and G2ANet [6] fall in this category. These

W. Yun, B. Lim, S. Jung, Y.-C. Ko, and J. Kim are with the School of Electrical Engineering, Korea University, Seoul 02841, Korea (email: {ywjoon95 ,limbj93 ,jungsoyi ,koyc ,joongheon}@korea.ac.kr).

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

‡M. Bennis is with the Centre for Wireless Communications, University of Oulu, Oulu 90014, Finland (email: mehdi.bennis@oulu.fi).

J. Park and J. Kim are the corresponding authors of this paper.

algorithms presumably achieve the highest reward, albeit at the cost of frequently exchanging a non-negligible amount of information on the states, actions, and rewards, which are thus far from meeting stringent latency constraints in URLLC. Second, *fully decentralized MADRL* is the type in which every agent trains and executes its policy without exchanging any information with other entities. These include IQL [8] and I-DQN [9], which however may not guarantee reliability due to the lack of understanding the interactions among agents. Lastly, *centralized training and decentralized execution (CTDE)* MARL lies between the aforementioned two extremes. Following the actor-critic architecture [10], an actor model stored at each agent determines its policy for both training and execution, while a centralized critic model evaluates the reward of all agents only during training.

In this work we aspire to build a novel CTDE MADRL framework for UAV aided beyond-5G URLLC, as visualized in Fig. 1. To this end, we first delve into the opportunities and limitations of the existing CTDE MARL frameworks. In Differentiable Inter-Agent Learning (DIAL) [11], while taking actions, the agents exchange clean-slate messages passed through their actor models. During training, these messages are progressively turned into meaningful representations for better inter-agent collaboration, hereafter referred to as *semantic representations (SRs)*, which is an analog of children's developing language-based cues as they grow. These emergent SRs are effective in achieving higher rewards at execution, yet its starting from clean-slate messages is too inefficient to outperform other state-of-the-art CTDE frameworks. Alternatively, during training, the agents in QMIX [12] and CSGA [13] exchange global states and local graph attention, respectively, thereby achieving competitive performance. However, they do not share any local information during execution, so cannot reach the full potential of CTDE MADRL.

By integrating DIAL's emergent SR learning into graph attentive CSGA, in this work we propose a novel CTDE MADRL framework, termed a *graph attention exchange netowrk (GAXNet)*. In CSGA, each agent $n$ at every time $t$ locally constructs a star-topological graph as shown in Fig. 2a. The edge weight $w_{n,m}^t$ increases with the level of attention to the agent $m$ when the agent $n$ takes its action. For mutually symmetric agent interactions (e.g., collision, repulsion, etc.), the rationale being what I attend to you should ideally be the same as what you attend to me, i.e., $w_{n,m}^t = w_{m,n}^t$. This condition is often violated in CSGA due to its local attention graph construction. Inspired by DIAL, we overcome such inter-agent attention mismatch via developing and exchanging emergent SRs.

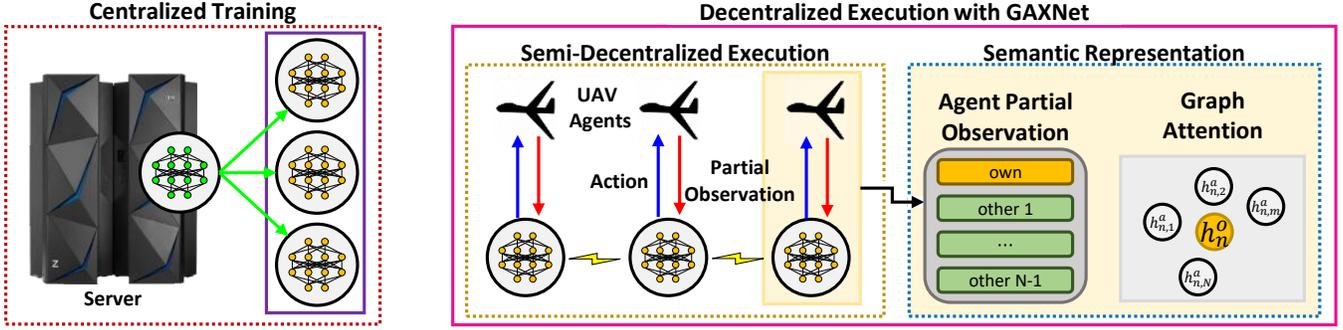To be precise, the agent $n$ runs an SR encoder whose

Fig. 1. A schematic illustration of the proposed graph attention exchange network (GAXNet) and its application to a multi-UAV emergency network.

input is the agent $n$'s current attention to the agent $m$ and its counterpart attention to the agent $n$ from the agent $m$ in the previous time slot, resulting in the output SR $\bar{w}_{n,m}^t$, i.e., $\bar{w}_{n,m}^t \overset{\text{SR}}{\leftarrow} \{w_{n,m}^t, \bar{w}_{n,m}^{t-1}\}$ as illustrated in Fig. 2b. This SR $\bar{w}_{n,m}^t$ is exchanged across agents for constructing SR $\bar{w}_{n,m}^{t+1}$ in the next time slot. In contrast to DIAL that initially exchanges meaningless messages, GAXNet exchanges semantically meaningful attention weights from the beginning, promoting better SR emergence. Consequently, the agent $n$ determines its action at time $t$ based on the emergent SR $\bar{w}_{n,m}^t$, as opposed to $w_{n,m}^t$ in CSGA, thereby reflecting the semantic importance of the attention at both agent sides, aiming at ensuring $\hat{w}_{n,m}^t \approx \hat{w}_{m,n}^t$ for all interacting agents.

To show the effectiveness of GAXNet in UAV aided beyond-5G URLLC, we study a scenario where a moving ground URLLC user is supported at least by one of the multiple fixed-wing UAVs. These UAVs aim to hover within a range guaranteeing URLLC requirements, referred to as *URLLC range*, while avoiding inter-UAV collisions. Simulations validate that compared to QMIX [12], the proposed GAXNet significantly reduces inter-UAV collision occurrences, and achieves by up to $4.5$x higher reward for $5,000$ epochs during training, where the reward is increased when the agent satisfies the latency and reliability requirements, and is penalized when an inter-UAV collision occurs. At execution, GAXNet achieves up to $6.5$x lower latency with the target $10^{-7}$ error rate, compared to QMIX that fails to guarantee the target error rate.

## II. PRELIMINARIES: MADRL AND SELF ATTENTION

GAXNet relies on two key principles, CTDE MADRL and the self-attention mechanism, which are described in the following subsections.

### A. Self-Attention

Let's denote the input as $\mathcal{X} = \{\boldsymbol{x}_1, \cdots, \boldsymbol{x}_k, \cdots, \boldsymbol{x}_K\}$. The input goes through the encoding process that converts into a trainable form. The process appears in the form of a linear combination, and the formula is summarized:

$$\mathcal{H}^{enc} = \boldsymbol{W}^{enc} \cdot \mathcal{X} + \boldsymbol{b}^{enc}, \tag{1}$$

where $\boldsymbol{W}^{enc}$ and $\boldsymbol{b}^{enc}$ stand for the weight and bias of the encoding layer, respectively. Self-attention is the process of determining the coefficient between inputs using scale-dot
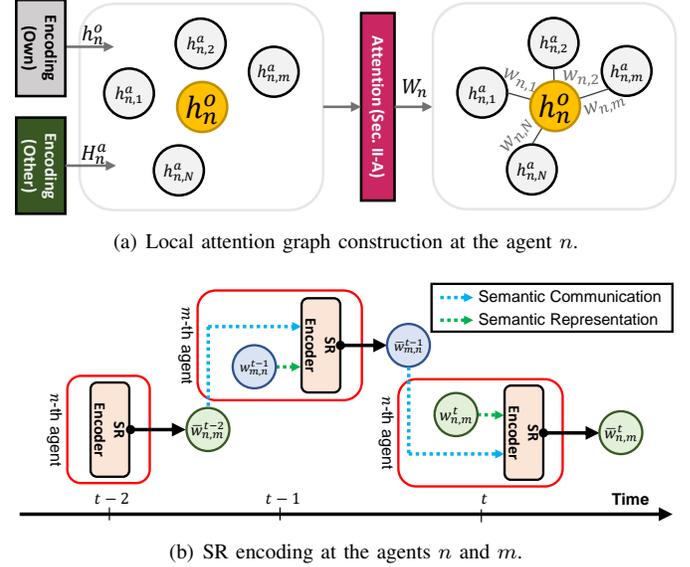


(a) Local attention graph construction at the agent $n$.



(b) SR encoding at the agents $n$ and $m$.

Fig. 2. The operations of GAXNet: (a) local attention graph construction at the agent $n$ and (b) semantic representation (SR) encoding at the agents $n$ and $m$ for 3 consecutive time slots.

production [14]. It consists of query layer, key layer and, value layer. Note that query layer, key layer and, value layer are denoted as $l^q$, $l^k$, $l^v$. Self-attention converts the received hidden variable into query, key, and value. This process can be expressed as follows:

$$\mathbf{Q} = l^q(\mathcal{H}^{enc}) = \{\boldsymbol{q}_1, \cdots, \boldsymbol{q}_k, \cdots, \boldsymbol{q}_K\}, \tag{2}$$

$$\mathbf{K} = l^k(\mathcal{H}^{enc}) = \{\boldsymbol{k}_1, \cdots, \boldsymbol{k}_k, \cdots, \boldsymbol{k}_K\}, \tag{3}$$

$$\mathbf{V} = l^v(\mathcal{H}^{enc}) = \{\boldsymbol{v}_1, \cdots, \boldsymbol{v}_k, \cdots, \boldsymbol{v}_K\}, \tag{4}$$

where $\boldsymbol{q}_k, \boldsymbol{k}_k$, and $\boldsymbol{v}_k$ stand for transformer of query, key, and value of the input $\boldsymbol{x}_k$. By taking dot product of query $\boldsymbol{q}_k$ and other input keys $\boldsymbol{K}_{-k} \triangleq \cup_{k' \neq k}^K \{\boldsymbol{k}_{k'}\}$, the key component of other input information is obtained. After that, we multiply by the $k$-th 's value, and the weighted message matrix is obtained. Formally, $W$ is given as:

$$W = \mathsf{Attention}(\mathbf{Q}, \mathbf{K}^{-1}, \mathbf{V}) = \frac{\mathbf{Q} \cdot \mathbf{K}^{-1} \cdot \mathbf{V}}{\sqrt{K}}, \tag{5}$$

where $W$ is $K$-by-$(K-1)$ matrix which represents the weights of adjacency matrix, $\sqrt{K}$ is a scaling factor that prevents the attention value diverge and note that $\mathbf{K}^- \triangleq$

$\{\boldsymbol{K}_{-1}, \cdots, \boldsymbol{K}_{-n}, \cdots, \boldsymbol{K}_{-N}\}.$

### B. Centralized Training and Decentralized Execution MADRL

In this subsection, we introduces CTDE MADRL in briefly. CTDE MADRL is mathematically modeled with Dec-POMDP. Dec-POMDP is a stochastic decision making model that is to consider the uncertainty between agents [15]. Dec-POMDP is defined as a tuple $G = \langle S, A, P, r, Z, O, n, \gamma \rangle$. In environment, the finite set of agents is denoted as $U = \{u^1, u^2, \cdots, u^N\}$, where $N \in \mathbb{N}^1$, respectively. The global state of the environment is defined as $s \in S$. The action of $n$-th agent is defined as $a^n \in A$. Since all agents can observe local information $z \in Z$, the observation function is defined as $O(s, u) : S \times U \to Z$.

The policy of the agent $\pi^n$ takes the observation-action history $\tau^n \in T$ and decides action as $\pi^n(a^n|\tau^n) : T \times A$. In addition, the joint action of all agents is $\boldsymbol{a} \in \boldsymbol{A}$. The state transition function is denoted as $P(s', s, \boldsymbol{a}) : S \times A \times S \to [0, 1]$. All agents obtain same reward $r(s, \boldsymbol{a}) : S \times \boldsymbol{A} \to \mathbb{R}$. The joint action value function at time $t$ is defined as follows:

$$Q^\pi(s_t, \boldsymbol{a}_t) = \mathbb{E}\left[\sum_{i=t+1}^{\infty} \gamma^{i-t-1} r_i \mid s_t, \boldsymbol{a}_t\right], \quad (6)$$

where $\pi$ and $\gamma \in [0, 1)$ stand for joint policy and discount factor, respectively.

The objective of MADRL system is to maximize the cumulative reward. In order to maximize cumulative reward under the consideration of centralized training and decentralized execution and according to [12], the joint action value function is decomposed by agent's individual utility function. The optimal joint action value function is as follows:

$$\arg\max_{\boldsymbol{a}} Q^\pi(\boldsymbol{\tau}, \boldsymbol{a}, s; \theta) = [\cdots, \arg\max_{a^n} Q^n(\tau^n, a^n; \theta^n), \cdots], \quad (7)$$

where $\theta$ and $\theta^n$ are the parameters of $\pi$ and $\pi^n$, respectively. To summarize CTDE, the joint action value function leads all agents' optimal utility function, all agents determine their actions via their policies.

### III. UAV AUTOMATION FRAMEWORK FOR URLLC

### A. Network and Channel Model

As illustrated in Fig. 1, the network under study consists of a set $\mathcal{U} \triangleq \{u_1, \cdots, u_n, \cdots, u_N\}$ of $N$ fixed-wing UAV agents and a moving target location $c$. For simplicity, we hereafter focus only on UAV-to-ground channels, while assuming the inter-UAV communications are always successful with negligible delays. This is not a strong assumption in our experimental settings with $N = 4$ (see Sec. IV) wherein the payload size to be received during a unit time slot is only 576 bits.

According to [16], the air-to-ground mean path loss for $xy$-direction distance $d$ and $z$-direction distance $h$, is given as:

$$\text{PL}(d, h) = \frac{\eta_{LoS} - \eta_{NLoS}}{1 + \alpha \cdot \exp\left(-\beta\left(\frac{180}{\pi}\tan^{-1}\left(\frac{h}{d}\right) - \alpha\right)\right)} \quad (8)$$
$$+ 10\log_{10}(h^2 + d^2) + 20 * \log_{10}\left(\frac{4\pi f_c}{c}\right) + \eta_{NLoS},$$

where $\eta_{LoS}$ and $\eta_{NLoS}$ stand for additional pathloss of LoS and NLoS. Also, $\alpha$ and $\beta$ represent environmental constant,

respectively. Then, Signal-to-noise ratio (SNR) can be calculated as follows:

$$\text{SNR}(d, h) = \frac{P_t}{P_n} * 10^{-\text{PL}/10}, \quad (9)$$

where $P_t$, $P_n$, and $W$ stand for transmit power, noise power, and bandwidth of A2G, respectively. According to [17], the error rate is obtained as follows:

$$\varepsilon = \frac{1}{\sqrt{2\pi}} \int_{f(\text{SNR})}^{\infty} e^{-\frac{t^2}{2}} \mathrm{d}t = Q_f(f(\text{SNR})), \quad (10)$$

where $Q_f(\cdot)$ is the Q-function and $f(\text{SNR})$ is given by $f(\text{SNR}) = \sqrt{\frac{W \cdot T_{\max}}{1 - (1 + \text{SNR})^{-2}}}\left(\ln(1 + \text{SNR}) - \frac{R_s \ln 2}{W}\right)$, and $T_{\max} = L_B/W$ and $R_s$ represent the maximum transmission time and achievable throughput, respectively. The achievable throughput $R_s$ can be calculated as follows:

$$R_s \simeq \ln(1 + \text{SNR}) - \frac{1 - (1 + \text{SNR})^{-2}}{L_B}Q_f(\varepsilon_0), \quad (11)$$

where $\varepsilon_0$ is the required reliability.

### B. MADRL Model

Next, we introduce the multi-agent reinforcement learning formulation (*e.g.*, state, action, and reward).

*1) State space:* The state space of our proposed UAV automation framework of UAVs/target location information, and relative position information with other UAVs/targets. The position of $u_n$ at time $t$ is defined as $l_t^n$, where $\forall l_n = (x_t^n, y_t^n), n \in [1, N]$. In addition, $l_t^{n,c} = (x_t^{n,c}, y_t^{n,c})$ and $d_t^{n,c}$ denote the relative positions and distance for the $u_n$ with $c$, respectively. The relative position and distance of $u_n$ and $u_m$, $m \in [1, N] \backslash n$ is denoted as $l_t^{n,m}$ and $d_t^{n,m}$. In addition, when $u_n$ is observable of $u_m$, $e_t^{n,m} = 1$, and $e_t^{n,m} = 0$ when unobservable. The observation of $u_n$ at $t$ consists of position information, relative position information and distance information about the target and other agents as follows:

$$o_t^n = \underbrace{\bigcup_{i=0}^{1}\{l_{t-i}^n, l_{t-i}^{n,c}, d_{t-i}^{n,c}\}}_{\text{own}} \cup \underbrace{\bigcup_{m \neq n}^{N}\{l_t^{n,m}, d_t^{n,m}, e_t^{n,m}\}}_{\text{other agents}},$$

Note that when $u_n$ is not observable (*i.e.*, $e_t^{n,m} = 0$), $l_t^{n,m}, d_t^{n,m}, e_t^{n,m}$ are zero. The true state information $s_t$ in the environment represents the observation information of all agents and the absolute position information of the target in time $t$ and $t - 1$.

$$s_t = \bigcup_{i=0}^{1}\left\{\bigcup_{n=1}^{N}\{o_{t-i}^n\} \cup l_{t-i}^c\right\}, \quad (12)$$

where $l_t^c$ stands for the absolute position of the target.

*2) Action space:* The action space of the UAV automation framework consists of 8 discrete actions, which is defined as $A \triangleq \{(x+\nu, y), (x-\nu, y), (x, y+\nu), (x, y-\nu), (x+\frac{1}{\sqrt{2}}\nu, y+\frac{1}{\sqrt{2}}\nu), (x-\nu, y+\frac{1}{\sqrt{2}}\nu), (x+\frac{1}{\sqrt{2}}\nu, y-\frac{1}{\sqrt{2}}\nu), (x-\frac{1}{\sqrt{2}}\nu, y-\frac{1}{\sqrt{2}}\nu)\}$, where $\nu$ is the speed of UAV. At time $t$, the action of $u_n$ is denoted as $a_t^n \in A$, where the joint action is $\boldsymbol{a_t} = \{a_t^1, \cdots, a_t^n, \cdots, a_t^N\}$.

*3) Reward space:* The reward space is designed to make UAV agent quickly reach the target area and prevent collision. For $u^n$ to guarantee the URLLC requirements, $u^n$ should reach to the target quickly, the URLLC reliability reward is

$$r_t^{n,c} = \begin{cases} 1, & \text{if. } d_t^{n,c} < d^c, \\ 0.05, & \text{if. } d^c \le d_t^{n,c} < d_{t-1}^{n,c} \\ -0.01, & \text{otherwise.} \end{cases} \quad (13)$$

where $d^c$ represents the URLLC range for URLLC reliability, respectively. If $d_t^{n,c}$ is less than $d^c$, the error rate of $u^n$ is less than the required error rate for URLLC and vice versa. In addition, the collision reward for minimizing collision between agents is presented as follows:

$$r_t^{i,n,m} = \begin{cases} -0.5, & \text{if. } d_t^{n,m} < d^i, \\ 0, & \text{otherwise.} \end{cases} \quad (14)$$

When the distance between $u_n$ and $u_m$ is less than $d^i$, collision occurs and vice versa.

According to (13) and (14), the final reward at time $t$ is as follows:

$$r_t = \underbrace{\sum_{n=1}^{N} r_t^{n,c}}_{\text{UAV-to-ground URLLC}} + \underbrace{\sum_{n=1}^{N} \sum_{m \ne n}^{N} r_t^{i,n,m}}_{\text{Inter-UAV collision}} \quad (15)$$

Note that all agents have a common reward.

### C. Policy Updates Using Graph Attention and SR Encoding

In this section, we introduce how to configure policies to generate and transfer semantic messages. Since the agent's observation consists of its own unique information and partial information of other agents in Dec-POMDP, the agent information can be operated by categorization of the agent's observation into the agent's unique information and the other agent's information.

$o_t^n$ is encoded by taking input in the form taken by the flatten on the out characteristic encoding layer, while partial information about other agents is stacked to encode in the other agents encoding layer as follows:

$$h_n^o = f_{own}^{enc,n}(\text{Flatten}(o_t)), \quad (16)$$
$$H_n^a = f_{oth}^{enc,n}(\text{Stack}(o^{oth,n,m})), \forall m \in [1, N] \backslash n, \quad (17)$$
$$\text{where,} \quad o_t^{oth,n,m} \triangleq \{l_t^n, l_t^{n,d}, l_t^{n,m}, d_t^{n,c}, d_t^{n,m}, e_t^{n,m}\},$$

and $f_{own}^{enc,n}$, $f_{oth}^{enc,n}$ and $o_t^{oth,n,m}$ stand for $n$-th agent's own characteristics encoding layer, other agent's characteristics encoding layer, and observation about $u_m$ of $u_n$, respectively. In addition, $H_n^a$ is a set of vector, which represents to $H_n^a = \cup_{m \ne n}^{N} \{h_m^a, \}$. $h_n^o$, and $h_m^a$ have a vector size of $1 \times J_1$, and $H_n^a$ has a size of $N - 1 \times J_1$, where $J_1$ stands for the output dimension of encoding layer.

Fig. 2(a) illustrates the method that utilizes self-attention using $h_n^o$ and $H_n^a$ which are provided via (16)–(17). For $\mathcal{H}^{enc}$ in Sec. II-A, $\mathcal{H}^{enc}$ is redefined as $\mathcal{H}^{enc} = \{h_1^a, \cdots, h_{n-1}^a, h_n^o, h_{n+1}^a, \cdots, h_N^a\}$. In (2)–(4), $l^q$, $l^k$, $l^v$ are replaced by the query layer $l_n^q$, key layer $l_n^k$, and value layer $l_n^v$ of the $u_n$. Therefore the query, key, and value can be obtained through $q_n = l_n^q(h_n^o)$, $k_n = l_n^k(H_n^a)$, $v_n = l_n^v(H_n^a)$.

The size of each query, key, and value represents $1 \times J_2$, $N - 1 \times J_2$, and $N - 1 \times J_2$. $J_2$ represents the attention dimension, respectively. The weight of semantic representation for other agents can be obtained by taking (5), which is denoted as $W_n = \{w_{n,1}, \cdots w_{n,n-1}, w_{n,n+1}, \cdots, w_{n,N}\}$.

Let's assume the exchange of the semantic weights $w_{n,m}$ and $w_{m,n}$ of two agents $u_n$ and $u_m$ for semantic communication. Fig. 2(b) shows the overall process of semantic communication between $u_n$ and $u_m$.

The weight information $w_{n,m}$ which is generated by $u_n$, is the semantic weight created from the complete information of the $u_n$ and the partial information of the $u_m$ and vice versa. At time $t$, $u_n$ receives the $w_{m,n}^{t-1}$ which is created by the $u_m$ in the previous step $t-1$, through the semantic channel. The synthesized weight information is created by aggregating $w_{n,m}$, and $w_{m,n}$ through the SR encoder. The synthesized weight is denoted as $\bar{w}_{n,m}^t$. The process of semantic communication is as follows:

$$\bar{w}_{n,m}^t = \begin{cases} \text{RNN}^n(w_{n,m}^t, \bar{w}_{m,n}^{t-1}), & \text{if. } e_t^{n,m} = 1 \\ \text{RNN}^n(w_{n,m}^t, \vec{0}), & \text{otherwise.} \end{cases} \quad (18)$$

where $\text{RNN}^n$ and $\vec{0}$ stand for SR encoder, and $N - 1 \times 1$ size zero vector, respectively. `GRUCell` is used for SR encoder [18]. Note that $u_n$ obtains the semantic weight $w_{m,n}$ from all observable agents excluding itself (*i.e.*, $e_t^{n,m} = 1$). Finally, $\pi^n$ determines $Q^n(\tau_t^n, a_t^n; \theta^n)$ which indicates the utility function of $u_n$, utilizing $h_t^{o,n}$ and $\bar{W}_t^n$.

### D. Centralized Training and Decentralized Execution

In this subsection, we introduce a CTDE method for UAV aided URLLC systems, consisting of $N$ decentralized actors and one centralized critic denoted as $\Phi(\theta)$. QMIX architecture is adopted as the centralized critic [12]. The utility function of all agents, and environment true state are the input of $\Phi(\theta)$. In addition, the observation and action of all agents, current state, reward, the next observation of all agents, and the next state are stored as tuple $b = \langle O, S, \boldsymbol{A}, R, O', S' \rangle$ to the replay buffer. In training phase, $B = \{b_1, \cdots, b_i, \cdots, b_I\}$ are sampled from the replay buffer. $B$ is training data to train the parameters of actors and critic. Temporal difference error is used as loss function [19], i.e.,

$$L(\theta) = \sum_{b \in B} \left[ r + \gamma \max_{\boldsymbol{a}} Q^{\Psi}(\boldsymbol{\tau}', \boldsymbol{a}', s'; \theta^-) - Q^{\Phi}(\boldsymbol{\tau}, \boldsymbol{a}, s; \theta) \right]^2, \quad (19)$$

where $\Phi(\theta^-)$ stands for target network. The agent is trained with the direction to maximize the joint-action value function, and if the agent can communicate with others, it can reinforce cooperation through semantic communications.

## IV. PERFORMANCE EVALUATION

In order to verify the potential of the proposed approach, the environment is configured as shown in Fig. 3. The environment has a 2D $3,750$m $\times 3,750$m grid. 4 UAV agents are located in the grid (i.e., $N = 4$), a server is also located in the center of the grid, and users are also configured for moving around the server. The experiment is devised by comparing the GAXNet

TABLE I
THE PARAMETERS OF ENVIRONMENT AND TRAINING.

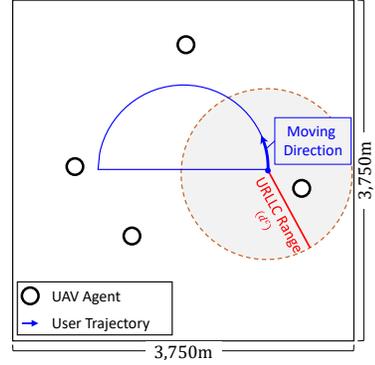| Simulation Parameter | Value |
|---|---|
| UAV speed ($\nu$) | 45 km/h |
| Target speed | 36 km/h |
| URLLC range ($d^c$) | 938 m |
| Collision distance ($d^i$) | 563 m |
| Carrier frequency ($f_c^g$) | 2 GHz |
| Bandwidth ($W$) | 20 MHz |
| Payload size ($L_B$) | 576 bits |
| Transmit Power ($P_t$) | 46 dBm |
| Noise Power ($P_n$) | $-99$ dBm |
| Time step per episode | 20 |
| **Training Parameter** | **Value** |
| Default number of nodes ($J_1$) | 64 |
| Number of node in attention layers ($J_2$) | 32 |
| epsilon-greedy constant ($\delta$) | 0.3 |
| Annealing step | 1,000 |
| Training iterations | 5,000 |
| Batch size ($I$) | 64 |
| Learning rate of GAXNet | $8 \times 10^{-4}$ |
| Learning rate of QMIX | $1 \times 10^{-4}$ |



Fig. 3. The simulation environment

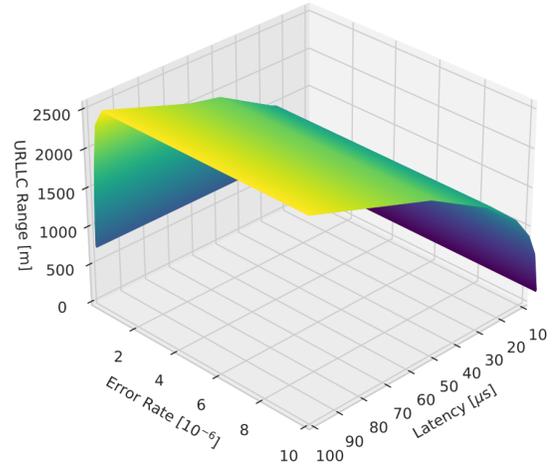

Fig. 4. A2G network specification for URLLC: The relationship between three components; reward distance, error rate, and latency.

with QMIX where the QMIX is the state of the art of CTDE. The experiment was conducted using Python 3.6 and Pytorch. In addition, Table I presented for specific experimental parameters. Moreover we set $\alpha = 9.61$, $\beta = 0.16$, $\eta_{LoS} = 1$ [dB], and $\eta_{NLoS} = 20$ [dB]. For fixed $L_B$, the transmission time is given as $T_{\max}$ and it can be assumed that $T_{\max}$ is less than the latency requirements. Actually, $T_{\max}$ can be calculated as 28.8 $\mu$s given $L_B = 576$ bits and $W = 20$MHz and it is within the 5G NR transmit latency [20]. To meet the reliability requirements, $\varepsilon \leq \varepsilon_0$ needs to be satisfied. By (10) and the property of monotonic increasing function $f(\text{SNR})$, the required SNR is

$$\text{SNR} \geq f^{-1}(Q_f^{-1}(\varepsilon_0)), \tag{20}$$

where $f^{-1}(\cdot)$ and $Q_f^{-1}(\cdot)$ is the inverse function of $f(\cdot)$ and $Q_f(\cdot)$, respectively. Due to (9), (20) can be rewritten as

$$\text{PL}(d, h) \leq -10 \log_{10}\left(\frac{P_n}{P_t} f^{-1}\left(Q_f^{-1}(\varepsilon_0)\right)\right). \tag{21}$$

If the equality holds in (21) when $d = d^*$, the reliability requirement is always satisfied with $d \leq d^*$ because the path loss in (8) monotonically increases with $d$. Therefore, we can obtain the URLLC range $d^c$ to satisfy the reliability requirements, $\varepsilon_0$, using (21). Fig. 4 shows the relationship between reward distance, required latency, and required error rate for URLLC via (20)–(21). For simulation, we adopt the required error rate as $10^{-7}$, the required latency as $39\mu$s, and the URLLC range as 938m for URLLC reliability.

In reinforcement learning, we observe reward convergences, collisions, and target arrivals for two algorithm, i.e., GAXNet and QMIX-based algorithms. Fig. 5 shows the sum of total reward. As shown in Fig. 5, the reward of GAXNet converges around to 26 at $3,100$ iterations. However, the reward of QMIX fluctuates between $-3$ and 11. Fig. 6 represents the visual trajectory dynamics of UAV agents. Fig. 6(a) illustrates the UAV trajectory path planning with QMIX-based algorithm whereas Fig. 6(b) visually presents the UAV trajectory dy-

namics with GAXNet. As shown in Fig. 6(a), all agents do not trace the target, and there are collisions between agents. However, all agents using our purposed algorithm chase the target with out any collision as shown in Fig. 6(b). In addition, collision between agents with GAXNet does not occur.

In terms of URLLC, Fig. 7 shows the reliability of both latency and error rate. The GAXNet always satisfy the latency reliability ($39\mu$s) and error rate reliability ($10^{-7}$). However, QMIX only satisfy at the time slots [3; 4]. During the rest of time slots, the QMIX-based UAV agents communicate under high latencies and error rates. This implies that QMIX-based UAV agents may not be able to guarantee URLLC reliability.

Therefore GAXNet outperforms the baseline regarding the objectives of the system.

We investigate the differences between untrained GAXNet without exchanging semantic weights, and trained GAXNet with exchanging. Fig. 8 represents the weighted adjacency matrix of 4 agents. In the case of GAXNet without exchanging semantic weights, the agents which have the maximum difference between $\bar{w}_{n,m}^t$ and $w_{m,n}^{t-1}$ are $u^2$ and $u^4$. The maximum difference is 0.18. Calculating the mean squared error (MSE) of all anti-diagonal elements (i.e., $\bar{w}_{n,m}^t$ and $w_{m,n}^{t-1}$), the result is 0.043. On the other hand, in the case of GAXNet with exchange, the maximum difference of two weights is 0.07, and MSE is 0.014. There is correlation of $o^{oth,n,m}$ with $o^{oth,m,n}$,
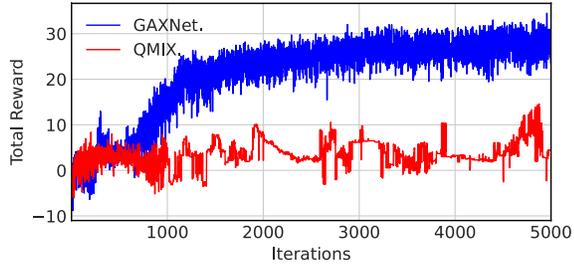
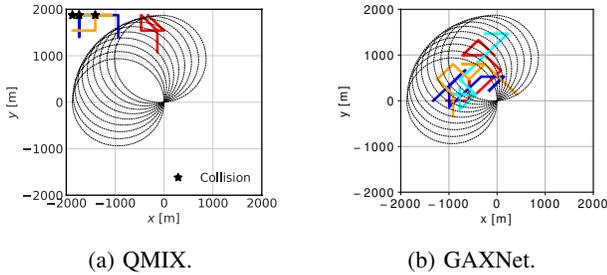Fig. 5. The learning curve of two algorithms.



(a) QMIX.

(b) GAXNet.

Fig. 6. The trajectory of UAV agents in 10 time slots. The *dotted circle*, *star marker*, *bold line* represent the target area, collision area, the trajectory of UAV agents, respectively.
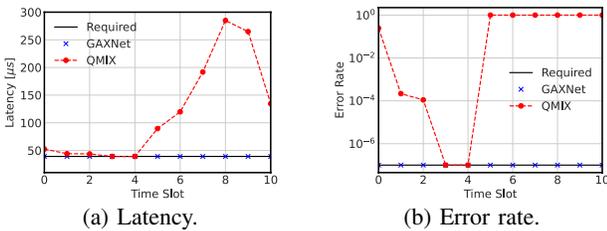


(a) Latency.

(b) Error rate.

Fig. 7. Comparison of two algorithms in respect to URLLC reliability in 10 time slots.



(a) GAXNet without exchange.
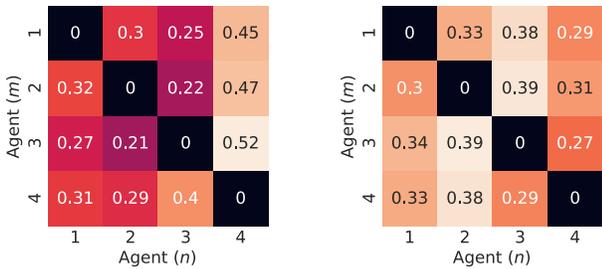
(b) GAXNet with exchange.

Fig. 8. Semantic attention weights.

because the relative position of $(u^n, u^m)$ has a relationship (*e.g.*, $l_t^{n,m} = -l_t^{m,n}$) and the distance and connectivity information between $u^n$ and $u^m$ are mutually same (*e.g.*, $d_t^{n,m} = d_t^{m,n}$ and $e_t^{n,m} = e_t^{m,n}$). For these reason, the attention weight between two agents (*e.g.*, $w_{n,m}^t$ and $w_{m,n}^t$ should be similar. In other words, the weighted adjacency matrix as shown in Fig. 8 should be diagonal. As shown in Fig. 8(a)/(b), GAXNet with exchange makes weighted adjacency matrix more symmetric than GAXNet without exchange. Because the SR encoder constructs symmetric matrix, we conclude that SR encoder is able to build semantic representation, successfully.

## V. CONCLUSION

In this paper, we developed a novel CTDE MADRL framework for UAV aided URLLC. The proposed solution, GAXNet, was shown to achieve lower latency with higher reliability compared to a state-of-the-art CTDE method, QMIX. To generalize and improve and extend GAXNet, incorporating realistic inter-UAV channels as well as considering a continuous UAV control action space could be interesting topics for future research.

## REFERENCES

[1] S. Jung, W. J. Yun, M. Shin, J. Kim, and J.-H. Kim, "Orchestrated scheduling and multi-agent deep reinforcement learning for cloud-assisted multi-UAV charging systems," *IEEE Transactions on Vehicular Technology*, pp. 1–1, 2021.

[2] H. Kim, J. Park, M. Bennis, and S.-L. Kim, "Massive UAV-to-ground communication and its stable movement control: A mean-field approach," in *Proc. of IEEE SPAWC*, Kalamata, Greece, Jun. 2018.

[3] J. Park, S. Samarakoon, H. Shiri, M. K. Abdel-Aziz, T. Nishio, A. Elgabli, and M. Bennis, "Extreme URLLC: Vision, challenges, and key enablers," *arXiv preprint arXiv:2001.09683*, 2020.

[4] J. Park, S. Samarakoon, M. Bennis, and M. Debbah, "Wireless network intelligence at the edge," *Proceedings of the IEEE*, vol. 107, no. 11, pp. 2204–2239, October 2019.

[5] R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," in *Proceedings of the Advances in Neural Information Processing Systems (NIPS)*, NY, USA, December 2017, pp. 6382–6393.

[6] Y. Liu, W. Wang, Y. Hu, J. Hao, X. Chen, and Y. Gao, "Multi-agent game abstraction via graph attention neural network," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 05, NY, USA, February 2020, pp. 7211–7218.

[7] S. Sukhbaatar, A. Szlam, and R. Fergus, "Learning multiagent communication with backpropagation," in *Proceedings of the Advances in Neural Information Processing Systems (NIPS)*, NY, USA, December 2016, pp. 2252–2260.

[8] M. Tan, "Multi-agent reinforcement learning: Independent vs. cooperative agents," in *Proceedings of the International Conference on Machine Learning (ICML)*, MA, USA, June 1993, pp. 330–337.

[9] A. Tampuu, T. Matiisen, D. Kodelja, I. Kuzovkin, K. Korjus, J. Aru, J. Aru, and R. Vicente, "Multiagent cooperation and competition with deep reinforcement learning," *PLOS ONE*, vol. 12, no. 4, pp. 1–15, 04 2017.

[10] V. R. Konda and J. N. Tsitsiklis, "Actor-critic algorithms," in *Proceedings of the Advances in Neural Information Processing Systems (NIPS)*, CO, USA, January 2000, pp. 1008–1014.

[11] J. Foerster, I. A. Assael, N. de Freitas, and S. Whiteson, "Learning to communicate with deep multi-agent reinforcement learning," in *Proceedings of the Advances in Neural Information Processing Systems (NIPS)*, vol. 29, NY, USA, December 2016, pp. 2145–2153.

[12] T. Rashid, M. Samvelyan, C. Schroeder, G. Farquhar, J. Foerster, and S. Whiteson, "QMIX: Monotonic value function factorisation for deep multi-agent reinforcement learning," in *Proceedings of the International Conference on Machine Learning (ICML)*, Stockholmsmässan, Sweden, July 2018, pp. 4295–4304.

[13] W. J. Yun, S. Yi, and J. Kim, "Multi-agent deep reinforcement learning using attentive graph neural architectures for real-time strategy games," *submitted to IEEE SMC 2021*.

[14] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proceedings of the Advances in Neural Information Processing Systems (NIPS)*, CA, USA, December 2017, pp. 5998–6008.

[15] F. A. Oliehoek and C. Amato, *A Concise Introduction to Decentralized POMDPs*. Springer Publishing Company, Incorporated, 2016.

[16] Z. Zhu, L. Li, and W. Zhou, "QoS-aware 3D deployment of UAV base stations," in *Proceedings of the IEEE International Conference on Wireless Communications and Signal Processing (WCSP)*, October 2018, pp. 1–6.

[17] C. She, C. Sun, Z. Gu, Y. Li, C. Yang, H. V. Poor, and B. Vucetic, "A tutorial on ultrareliable and low-latency communications in 6G: Integrating domain knowledge into deep learning," *Proceedings of the IEEE*, vol. 109, no. 3, pp. 204–246, Mar 2021.

[18] R. Dey and F. M. Salem, "Gate-variants of gated recurrent unit (GRU) neural networks," in *Proceedings of International Midwest Symposium on Circuits and Systems (MWSCAS)*. IEEE, August 2017, pp. 1597–1600.

[19] 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*, December 2013.

[20] H. Ji, S. Park, J. Yeo, Y. Kim, J. Lee, and B. Shim, "Ultra-reliable and low-latency communications in 5G downlink: Physical layer aspects," *IEEE Wireless Communications*, vol. 25, no. 3, pp. 124–130, June 2018.