# Federated Multi-Agent DRL for Radio Resource Management in Industrial 6G in-X subnetworks

Bjarke Madsen and Ramoni Adeogun

*Department of Electronic Systems, Aalborg University, Denmark*

E-mail:ra@es.aau.dk

*Abstract*—Recently, 6G in-X subnetworks have been proposed as low-power short-range radio cells to support localized extreme wireless connectivity inside entities such as industrial robots, vehicles, and the human body. Deployment of in-X subnetworks within these entities may result in rapid changes in interference levels and thus, varying link quality. This paper investigates distributed dynamic channel allocation to mitigate inter-subnetwork interference in dense in-factory deployments of 6G in-X subnetworks. This paper introduces two new techniques, Federated Multi-Agent Double Deep Q-Network (F-MADDQN) and Federated Multi-Agent Deep Proximal Policy Optimization (F-MADPPO), for channel allocation in 6G in-X subnetworks. These techniques are based on a client-to-server horizontal federated reinforcement learning framework. The methods require sharing only local model weights with a centralized gNB for federated aggregation thereby preserving local data privacy and security. Simulations were conducted using a practical indoor factory environment proposed by 5G-ACIA and 3GPP models for in-factory environments. The results showed that the proposed methods achieved slightly better performance than baseline schemes with significantly reduced signaling overhead compared to the baseline solutions. The schemes also showed better robustness and generalization ability to changes in deployment densities and propagation parameters.

## I. INTRODUCTION

Wireless communication has evolved significantly, with advancements in 5G technology enabling faster and more reliable communication. However, the need for more efficient and reliable networks has led to the development of 6G technology, which promises ultra-reliable communication, higher data rates, and lower latency. Wireless communication has limitless applications, including virtual and augmented reality in entertainment, education, healthcare, and manufacturing, personalized healthcare solutions, and real-time environmental monitoring. To achieve these benefits, new communication technologies must be developed, such as 6G short-range low-power in-X subnetworks [1], which can support demanding requirements inside entities like robots, production modules, vehicles, and human-body scenarios. Ensuring data privacy and security is crucial for advanced wireless communication networks, especially for sensitive personal or environmental data. Techniques for resource allocation must intelligently guarantee privacy and security while adapting to changing wireless environments.

**Related work:** Interference management solutions for wireless systems have been proposed in the literature, see e.g., [2]. There has also been a significant amount of work targetting 6G in-X subnetworks within the last few years. These solutions can be categorized into heuristics, machine learning (ML) methods, and reinforcement learning (RL) techniques. Heuristic methods generate near-optimal solutions based on simple rules or algorithms, providing good results with low computational complexity. In [3], three heuristic-based algorithms are presented: $\epsilon$-greedy channel selection, nearest neighbor conflict avoidance, and minimum Signal to Interference plus Noise Ratio (SINR) guarantee. Centralized Graph Coloring (CGC) is presented as a benchmark algorithm, but it is impractical due to its high signaling and computation overhead. A distributed interference-aware dynamic channel selection is presented in [4]. In [5], channel selection is approached based on centralized selective graph constructions. ML-based solutions may offer a higher quality of service compared to heuristics by leveraging historical data and adapting to evolving network conditions.

In [6], a Deep Neural Network (DNN) was successfully trained in offline simulations using CGC with mobile subnetworks. Subsequently, this DNN was deployed for real-time distributed channel selection. A novel solution for centralized power control is presented in [7], where decisions are based on positioning information using Graph Neural Network (GNN).

In [8], joint allocation of channel and transmit power based on a distributed multi-objective optimization problem is addressed, and an approach of Q-learning for multiple agents based on limited sensing information is proposed. The GA-Net framework presented in [9] is a distributed framework Multi-Agent Reinforcement Learning (MARL) resource management, based on GNN.

**Contributions:** This paper investigates the potential of federated MARL techniques for efficient, privacy-preserving channel allocation in 6G in-X subnetworks. The authors present their first study on federated MARL for radio resource management (RRM) in short-range low-power 6G in-X subnetworks, presenting the first contribution on FRL for RRM in 6G in-X subnetworks. They propose two novel federated MARL solutions - F-MADDQN and F-MADPPO - for dynamic channel allocation in 6G in-X subnetworks. These techniques enable privacy-preserving collaborative training of a Double Deep Q-Network (DDQN) policy and proximal Policy Optimization (PPO) agent for dynamic channel selec-

tion by multiple subnetworks in the presence of an umbrella network. The training phase requires only offline connectivity to the umbrella network. The paper also presents performance evaluation in industrial factory environment simulations based on a realistic factory floor plan from 5G-ACIA and 3GPP channel models for in-factory environments.

## II. SYSTEM MODEL AND PROBLEM FORMULATION

**System Model:** We consider downlink transmissions in a wireless network comprising $N$ independent and mobile subnetworks, each responsible for serving one or more devices, such as sensors and actuators. Each subnetwork contains an Access Point (AP) which is responsible for managing transmissions with its associated devices. We denote the collection of subnetworks as $\mathcal{N} = 1, \cdots, N$, and the set of devices in the $n^{\text{th}}$ subnetwork as $\mathcal{M}_n = 1, \cdots, M_n$. It is assumed that each subnetwork's AP is equipped with a local Resource Manager (RM), which utilizes data acquired from wireless environment sensing or received from its devices. Within each mobile subnetwork, wireless transmission occurs over one of $K$ shared orthogonal channels. Given the limited availability of resources, the number of bands is typically much less than the number of co-existing subnetworks. We assume that transmissions within each subnetwork are orthogonal, hence there is no intra-subnetwork interference. In practical systems, it may be impossible to make transmissions completely orthogonal due to limitations on the bandwidth of the channel allocated to the subnetwork leading to intra-subnetwork interference. Consideration of the effects of such interference is however left for future work.

The received SINR, denoted $\gamma_{n,m}$, on a link between the $n^{\text{th}}$ AP and $m^{\text{th}}$ device can be expressed

$$\gamma_{n,m}(t) = \frac{g_{n,m}(t)}{\sum_{i \in \mathcal{I}_k} g_{n,i}(t) + \sigma^2}, \quad (1)$$

where $g_{n,m}(t)$ is the received power on the link between the $n$th AP and device $m$, $\mathcal{I}_k$ denotes the APs or devices transmitting on channel $k$ and $\sigma^2$ denotes the noise power and is defined as $\sigma^2 = 10^{(-174+\text{NF}+10\log_{10}(B_k))/10}$, where NF denotes the noise figure and $B_k$ is the bandwidth of the channel. We consider practical 6G in-X subnetworks featuring short packets (in the order of tens of bytes) transmission, making the infinite block length assumption and the Shannon approximation used in most of the existing studies unrealistic. To compute the capacity, we instead use the finite block-length approximation as [10]

$$r_{n,m} = B_k \left[ \log_2\left(1 + \gamma_{n,m}(t)\right) - \sqrt{\frac{V_{n,m}(t)}{l}} Q^{-1}(\epsilon) \log e \right],$$

where $Q$ is the complementary Gaussian cumulative distribution function based on the code-word decoding error probability $\epsilon$, $\log e \approx 0.434$ is a constant constraint of the loss, and $V$ is the channel dispersion defined as [10].

$$V_{n,m}(t) = 1 - \frac{1}{(1 + \gamma_{n,m}(t))^2}. \quad (2)$$

When the length of the code-word block $l$ tends toward infinity, the achieved rate approaches the classical Shannon's approximation.

**Problem Formulation:** We focus on a resource allocation problem that involves the distributed selection of channels to maximize the achieved sum rate while ensuring that devices within each subnetwork achieve a specified minimum rate. This optimization problem can be mathematically expressed as

$$\text{P}: \left\{ \max_{\{\mathbf{c}^t\}} \sum_{m=1}^{M} r_{n,m}\left(\mathbf{c}^t\right) \right\}_{n=1}^{N} \quad \text{st: } r_{n,m} \geq r_{\text{target}}, \ \forall n, m \quad (3)$$

Here, $\mathcal{K} = [1, \cdots, K]; k \in \{1, 2, \cdots, K\}$ represents the vector of channel indices selected by all subnetworks at time t. $r_{\text{target}}$ denotes the target minimum rate, assumed to be the same for all subnetworks. The problem described in (3) involves the joint optimization of $N$ conflicting non-convex objective functions, making it a challenging problem to solve. In this paper, federated MARL methods are proposed as a solution to this problem.

## III. MARL FOR RRM IN SUBNETWORKS

### A. MARL Algorithms

Resource selection in scenarios with multiple in-X subnetworks is cast as a decentralized partially observable Markov decision process (DEC-POMDP). A DEC-POMDP is formally represented as $\left(N, \mathcal{S}, \{\mathcal{A}^n\}_{n=1}^N, \mathcal{P}, \{\mathcal{R}^n\}_{n=1}^N\right)$. The set of all possible states for all agents is defined as the state space $\mathcal{S} = \mathcal{S}_1 \times \cdots \times \mathcal{S}_N$, and the joint action space containing all possible actions for the $n^{\text{th}}$ agent $\mathcal{A}^n$ is denoted $\mathbf{\mathcal{A}} = \mathcal{A}^1 \times \cdots \times \mathcal{A}^N$. The reward signal for the $n^{\text{th}}$ agent is denoted $\mathcal{R}^n : \mathcal{S} \times \mathbf{\mathcal{A}} \times \mathcal{S}$ and the transition probability from state $s \in \mathcal{S}$ to state $s' \in \mathcal{S}$ by joint action $\boldsymbol{a} \in \mathbf{\mathcal{A}}$ is denoted $\mathcal{P} = \mathcal{S} \times \mathbf{\mathcal{A}}$.

*1) MADDQN:* Multi-Agent Double Deep Q Network (MADDQN) [11], [12] extends the Double Deep Q Network (DDQN) algorithm to enable decentralized learning among multiple agents. MADDQN allows each agent to have its own local Q-network and learn from its individual experiences. To train the Q-network, each agent interacts with the environment and collects experiences by executing actions according to its current policy. The experiences, including the state, action, reward, and next state, are stored in either a shared replay buffer, $\mathcal{B}_c$ in case of centralized learning or individual reply buffers, $\{\mathcal{B}_n\}_{n=1}^N$ in distributed learning.

To stabilize the learning process, MADDQN utilizes target networks. These networks are copies of the local Q-networks and their weights are periodically updated by polyak averaging, which involves updating the target network weights with a fraction of the local network weights.

The Q-networks are optimized using mini-batches of experiences sampled from the replay buffer with the goal of learning decentralized policies that maximize the expected cumulative rewards for each agent. For each agent, the Q-network is trained to minimize the mean squared Bellman
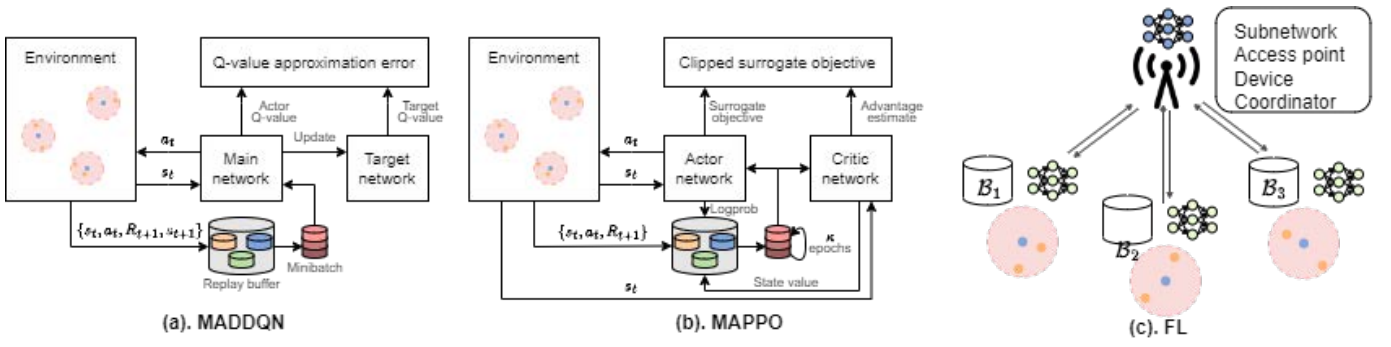
Fig. 1. Illustration of the training procedures for (a) multi-agent DDQN and (b) multi-agent PPO and the federated learning concept (c).

error, i.e., the squared difference between the estimated Q-value and the target Q-value. The loss function is defined as [12]

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^{N} (y_i - Q(s_i, a_i; \theta_i))^2 \qquad (4)$$

where $N$ is the mini-batch size, $s_i$ is the state, $a_i$ is the action, and $Q(s_i, a_i; \theta_i)$ is the estimated Q-value from the local Q-network. The variable, $y_i$ is the target Q-value which is calculated by using the target network to select the best action for the next state and then evaluating it with the local Q-network. The target Q-value is typically defined as

$$y_i = r_i + \gamma Q' \left( s', \max_a Q(s', a'; \theta_i); \theta_i' \right) \qquad (5)$$

where $r_i$ denotes the reward received by agent $i$, $s'$ is the next state, $a'$ is the action selected by the target Q-network, and $Q'$ is the target Q-network.

During training or execution, each agent selects actions using its local Q-network based on the observed state using the well-known $\epsilon$-greedy strategy. The actions can be either random with probability, $\epsilon$, to encourage exploration, or greedy with probability, $1 - \epsilon$, to exploit the learned policies.

*2) MAPPO:* Multi-agent Proximal Policy Optimization (MAPPO) is a state-of-the-art algorithm designed to address challenges in multi-agent environments. It is an extension of the Proximal Policy Optimization algorithm, originally proposed for single-agent reinforcement learning. PPO uses an actor-critic architecture, where each agent has its own policy and value function, allowing them to interact and learn decentralizedly. The core idea is to update the policy while ensuring deviation from the previous policy remains within a certain range. This is achieved by maximizing a surrogate objective function, which approximates the expected improvement of the policy using the current state-action distribution and the ratio between new and old policies. This balances exploration and exploitation, enabling stable learning in multi-agent environments.

In MAPPO, each agent interacts with the environment by executing actions according to its policy. Trajectories of state-action pairs are collected and stored in a replay buffer for efficient sampling. The collected trajectories are used to

estimate the value function for each agent, typically using a separate critic network to predict cumulative rewards from each state. The surrogate objective is computed based on the collected trajectories and the current policy, quantifying the expected improvement and determining the policy update direction. Denoting the ratio of the new policy, $\pi'_{\theta'}$ and the old policy, $\pi(\theta)$ as $r(\theta)$, the surrogate objective is defined as [13]

$$\mathcal{L}(\theta) = \mathbb{E} \left[ \min \left( r(\theta')A, \text{clip} \left( r(\theta', 1 - \eta, 1 + \eta) A \right) \right) \right], \qquad (6)$$

where $A$ denotes the advantage function which provides a measure of the advantage of taking a specific action under the current policy compared to the estimated value function. The parameter $\eta$ is a hyperparameter that controls the range of the policy deviation.

The policy is updated by optimizing the surrogate objective function using stochastic gradient descent (SGD) or a similar optimization technique. The objective is to maximize the surrogate objective with respect to the policy parameters, $\theta$. This is typically done by taking multiple gradient steps on the objective function, considering the trajectories collected from the replay buffer.

The update of the parameters, $\theta$ of the policy can be represented by

$$\theta' = \arg\max_{\theta'} \left\{ \frac{1}{N} \sum_{m=1}^{N} \min \left( r(\theta')A_m, C(\theta')A_m \right) \right\}, \qquad (7)$$

$$C(\theta') = \text{clip} \left( r(\theta', 1 - \eta, 1 + \eta) \right), \qquad (8)$$

where $N$ denotes the number of trajectories and $A_m$ is the advantage of the $m^{\text{th}}$ state-action pair.

The value function is updated to improve its estimation accuracy. This is typically done by minimizing the Mean-Squared Error (MSE) between the predicted value function and the actual cumulative rewards. The value function update is performed using SGD or another optimization algorithm, and the target is to minimize the loss function:

$$\mathcal{L}(\theta_v) = \frac{1}{N} \sum_{i=1}^{N} (V(s_i; \theta_v) - G_i)^2, \qquad (9)$$

where $\theta_v$ represents the value function parameters, $s_i$ is the $i^{\text{th}}$ state, $V_i$ is the predicted value function for state $s_i$, and $G_i$ is the actual cumulative reward obtained from the $i^{\text{th}}$ trajectory.

The steps above are repeated iteratively to improve the policy and value function over multiple iterations. The process continues until a desired level of performance or convergence is achieved.

### B. Proposed Methods

We propose two algorithms combining the MARL techniques, i.e., MADDQN and MAPPO described above with federated learning [14], [15] to solve the resource optimization problem in (3). The methods are referred to as Federated MADDQN (F-MADDQN) and Federated MAPPO (F-MAPPO). The algorithms require the definition of the environment, action space, state space, and reward function. The considered environment is the wireless network with $N$ in-X subnetworks. The other components of the proposed solutions are described below.

*1) Action Space:* As stated earlier, we assume that the available frequency band is divided into a set of $K$ equally sized channels. This leads to a $K$-dimensional action space for each subnetwork. Denoting the $k^{\text{th}}$ channel as $b_k$, the action space for the $n^{\text{th}}$ subnetwork, $\mathcal{A}^n$ is defined as

$$\mathcal{A}^n = \{b_1, \ldots, b_K\} \quad \forall n \in \mathcal{N} \tag{10}$$

*2) Observation Space:* The local observation of subnetwork $n$ at time $t$ is defined as

$$\mathbf{S}_n(t) = \left[ f(\mathbf{s}_n^1(t)), \cdots, f(\mathbf{s}_n^K(t)) \right]^T, \tag{11}$$

where $\mathbf{s}_n^k(t) = [\text{SIR}_{n1}^k(t), \cdots, \text{SIR}_{nM}^k] \in \mathbb{R}^{M \times 1}$ denotes the vector of SIRs measured at all devices in the $n^{\text{th}}$ subnetwork on channel $k$ and $f$ is the Observation Reduction Function (ORF) which is applied to the measurements to obtain the state to be used as input to an RL agent. The choice of ORF affects the input to the RL model and hence plays a significant role in the learning process. Considering the minimum rate constraint in (3), it appears reasonable to optimize for the worst case, i.e., allocate channels based on the minimum SIR over each channel. However, the combinatorial nature of the optimization problem makes it difficult to conclude whether this is the optimal choice. We therefore propose to use two categories of ORF defined as

$$f(\mathbf{x}) : \begin{cases} \mathbf{x} \mapsto \mathbf{x} \in \mathbb{R}^{M \times 1} & \text{Case I:Full state} \\ \mathbf{x} \mapsto f_A(\mathbf{x}) \in \mathbb{R}^1 & \text{Case II:Reduced state} \end{cases}, \tag{12}$$

where $f_A$ is the mean, max, or median aggregation functions.

*3) Reward Signal:* To guide the learning of the agents toward achieving this objective in (3), we define the reward function for the $n$ subnetwork as

$$R_n = \lambda_1 \sum_{m=1}^M r_{n,m} - \lambda_2 \sum_{m=1}^M \mathbb{1}(r_{n,m})(r_{\min} - r_{n,m}) \tag{13}$$

where $\lambda_i; i = 1, 2$ are scaling factors that are selected to create a balance between rate maximization and satisfaction of

the minimum rate constraint. The function, $\mathbb{1}(r)$, is a binary indicator function with a value equal to unity if and only if $r \geq r_{\min}$.

*4) Policy:* The main component of any reinforcement learning solution is the policy that defines how state measurements from the environment are mapped into actions. Representation of this state-action mapping and the associated optimization framework determines the policy. In this paper, we propose two methods viz: F-MADDQN and F-MAPPO. As stated earlier F-MADDQN and F-MAPPO combines federated learning with MADDQN and MAPPO, respectively. In both algorithms, the policy is modeled using a DNN.

*5) Training:* In both F-MADDQN and F-MAPPO, the agents are trained using federated learning (FL) [16] as shown in Figure 1c. Consider $N$ agents $\{\mathcal{F}_i\}_{i=1}^N$, each storing their respective data-sets in experience replay buffer, $\mathcal{B}_i$ from which a set of model parameters $\theta_i$ is learned. Denoting the loss function for the $i^{\text{th}}$ agent as $\mathcal{L}_i(\theta_i)$. the common global model loss can be defined as $\mathcal{L}_g(\theta)$ [16].

$$\mathcal{L}_g(\theta) = \sum_{i=1}^N \eta_i \mathcal{L}_i(\theta), \tag{14}$$

where $|\cdot|$ denotes the size of the set and $\eta_i > 0$ is the relative impact of each agent. Typically, the term $\eta$ is constrained to $\sum_{i=1}^N \eta_i = 1$. To allocate equal priority to all agents, we use $\eta_i = 1/N; \forall i \in \mathcal{N}$. The goal is to find the optimal parameters $w^*$, which minimizes the global loss function [16].

$$\theta^* = \arg \min_\theta \mathcal{L}_g(\theta) \tag{15}$$

One solution for (15) would be a gradient-descent approach, known as the federated averaging algorithm. Each agent uses its local data to perform a number of steps in gradient descent on current model parameters $\bar{\theta}(t)$. This gradient descent step is defined as

$$\theta_i(t+1) = \bar{\theta}(t) - \gamma \nabla \mathcal{L}_i(\bar{\theta}_i(t)) \quad \forall i \in \mathcal{N}, \tag{16}$$

where $\gamma > 0$ is the learning rate, and $\nabla f(\theta)$ for any scalar expression $f(\theta)$ denotes the vector of partial derivatives with respect to the components of the parameters $\theta$. During federated training, the global model is updated as

$$\theta_g(t+1) = \sum_{i=1}^N \frac{1}{N} \theta_i(t+1) \tag{17}$$

Once updated, the global weight is sent back to the agents at specified intervals referred to as *aggregation interval*, $\tau_{\text{agg}}$.

## IV. PERFORMANCE EVALUATION

### A. Simulation Settings

We consider an indoor industrial factory scenario inspired by existing production facilities of manufacturing companies as identified by the industry initiative, 5G alliance for connected industries, and automation [17]. Such a layout is implemented as a 180 m $\times$ 80 m hall containing several separate areas for production, assembly, storage, and human
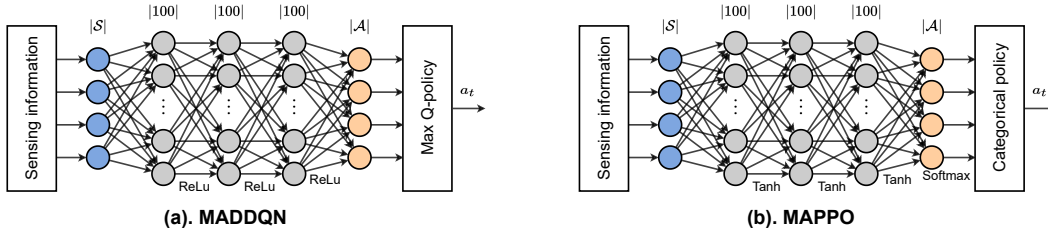
Fig. 2. Illustration of the (a) Deep Q-Network (DQN) architecture for MADDQN and (b) Deep PPO architecture for MAPPO which are used for the simulations.

TABLE I
SIMULATION PARAMETERS.

| Parameter | | Value |
|---|---|---|
| Total factory area | $\mathcal{R}$ | 180 m × 80 m |
| Clutter type table | | Sparse |
| Number of subnetworks | $N$ | 20 |
| Timestep | $t$ | 0.005 s |
| Number of episode | | 2000 |
| Number of steps per episode | $T$ | 200 |
| Subnetwork separation distance | $d_{\min}$ | 1 m |
| Subnetwork radius | $d_r$ | 1 m |
| Subnetwork velocity | | 3 m/s |
| Transmit power | $p_n(t)$ | -10 dBm |
| Number of frequency channels | $K$ | 4 |
| Carrier frequency | $f_c$ | 6 GHz |
| Bandwidth per subnetwork | $BW$ | 10 MHz |
| Noise figure | NF | 10 dB |
| Shadowing decorrelation distance | $d_\delta$ | 10 m |
| Max action switch delay | $\tau_{\max}$ | 10 |

work zones. Multiple in-robot subnetworks can be deployed with the task of transporting materials or tools around the facility. The alleys separating laboring areas are 5 m wide taking up $\sim 1600$ m$^2$ of the factory area and are outlined as two-lane roads in a right-handed traffic setting. In the deployment, the in-robot subnetworks are separated with a minimum distance of $d_{\min} = 1$ m and move with a speed of 3 m/s. The in-robot subnetworks are modeled by a circular coverage area with a radius of $d_r = 1$ m, making it possible for robots to pass each other in the alleys. Collisions are avoided by prioritizing robots with the shortest distances to a common intersection, slowing down any other robot that is within minimum separation distance.

Except where otherwise stated, we consider a total bandwidth of 100 MHz which is partitioned into $K = 4$ channels. Transmission within each subnetwork is then performed over a single channel at each time instant. The channel gain is calculated as $g_{n,m} = h_{n,m} \times \sqrt{10^{(PL_{n,m} + X^{\text{SD}}_{n,m})/10}}$, where $h_{n,m}$ denotes the small scale channel gain which is modeled as a temporally correlated Rayleigh random variable, $PL_{n,m}$ denotes the path loss and $X^{\text{SD}}_{n,m}$ is the shadow fading. We simulate the pathloss using the 3GPP Indoor Factory (InF) channel model for the Dense-clutter Low-antenna (InF-DL) and the Sparse-clutter Low-antenna (InF-SL) scenarios. The

shadow fading $X^{\text{SD}}_{n,m}$ is generated using the spatially correlated shadowing model used in [18]. Other simulation parameters are given in Table I. We consider different values of the aggregation intervals, i.e., $T^{\text{Agg}} = [128, 256, 512, 1024]$ for F-MAPPO and F-MADDQN. The minimum required data rate is set as, $r_{\min} = 11$ bps/Hz.
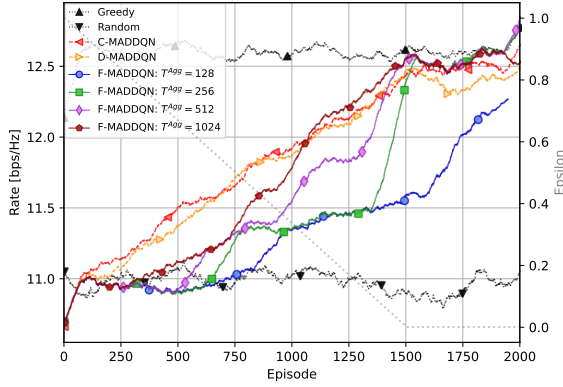
### B. Benchmarks

We benchmark the performance of the proposed F-MADDQN and F-MAPPO algorithms with the following methods.

1. **Centralized Graph Coloring (CGC)**: The CGC algorithm utilizes improper coloring to assign colors equivalent to channels for all subnetworks [19]. At each timestep $t$, the pairwise interference power relationships among subnetworks, $\boldsymbol{I}(t) \in \mathbb{R}^{N \times N}$, are collected, and mapped to a mutual coupling graph $G_t$. Each vertex corresponds to a subnetwork, and edges are created by connecting each vertex to its $K - 1$ nearest neighbors, where the weights of edges are equivalent to the interference power between the connected subnetworks.

2. **Greedy channel selection**: At each switching instant, each subnetwork selects the channel with the highest measured SINR from the previous time step.

3. **Random channel selection**: At the beginning of every episode, a random channel is allocated to each subnetwork.
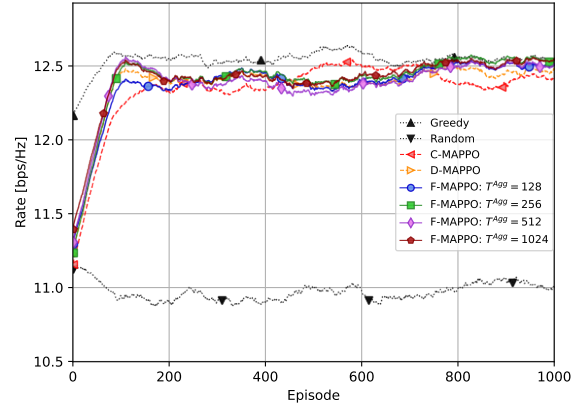
In addition, the federated learning solutions are also compared with previous solutions based on centralized as well as distributed training. We denote the centralized (distributed) MADDQN and MAPPO as C-MADDQN (D-MADDQN) and C-MAPPO (D-MAPPO), respectively.

### C. Training and Convergence

We consider the network architecture in Figure 2a and Figure 2b for F-MADDQN and F-MAPPO, respectively. The networks are trained using the federated learning procedure illustrated in Figure 1c. Figure 3a shows the averaged reward signal achieved at each training episode by the MADDQN-based methods. At convergence approximately after 1500 episodes, the MARL methods approach the distributed greedy selection baseline. At convergence, the F-MADDQN methods with greater aggregation intervals achieve a marginally better performance than the centralized and distributed baselines, indicating a potential advantage of the federated learning

(a) MADDQN.



(b) MAPPO.

Fig. 3. Averaged reward versus number of episodes.
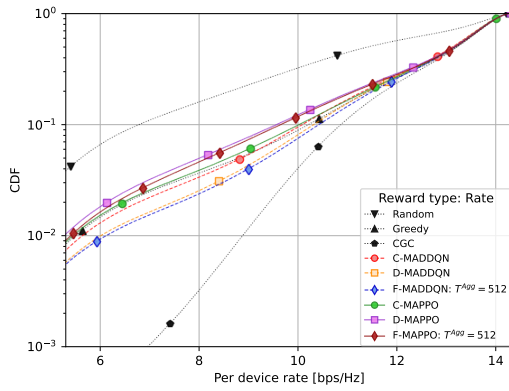


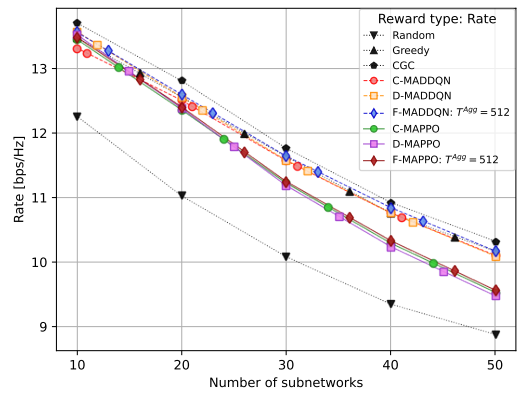Fig. 4. CDF of achieved rate.
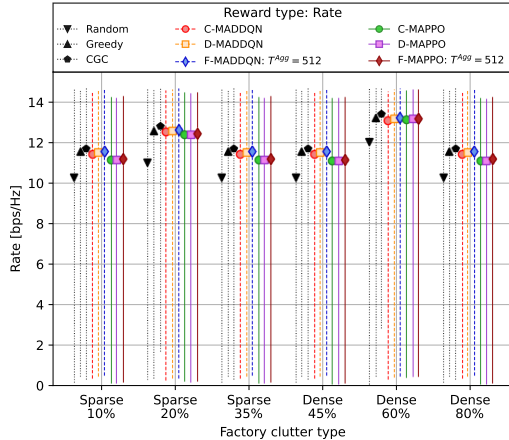


Fig. 6. Rate versus number of subnetworks.



Fig. 5. Sensitivity evaluation.

framework. The figure also shows that too frequent aggregation of DQN weights due to a low value of aggregation interval, $\tau_{\mathrm{agg}}$ can result in a comparatively slower convergence rate. For example, while F-MADDQN with $\tau_{\mathrm{agg}} \geq 256$ converged at $\approx 1500$ episodes, the algorithm failed to converge after 2000 episodes with $\tau_{\mathrm{agg}} = 128$.

We evaluate the convergence performance of F-MAPPO

with the in Figure 3b. At convergence, achieved approximately after $\approx 750$ episodes, the F-MAPPO achieves similar performance to the benchmark for all values of the aggregation interval, $\tau_{\mathrm{agg}}$. Compared to the methods based on MADDQN, the MAPPO-based algorithms give a reduction of about $50\%$ in the number of episodes required for convergence. A plausible explanation for this is the effect of the decay rate of the $\epsilon$-greedy parameter of MADDQN and the iterative training over multiple epochs at each training step of MAPPO.

### D. Performance Comparison

The trained models are deployed in each in-robot subnetwork for distributed channel allocation, and the performance is compared with random, greedy, CGC, and centralized and distributed training frameworks. To evaluate the trained models, the aggregation interval is set to $T^{\mathrm{Agg}} = 512$. Figure 4 shows the CDF of the achieved rate per device. The proposed F-MADDQN performs marginally better than D-MADDQN and C-MADDQN below the $30^{\mathrm{th}}$ percentile. Furthermore, the performance of F-MADDQN is also better than that of the distributed greedy scheme. F-MAPPO similarly achieves marginally better performance than the distributed baseline and approaches the centralized baseline below the $30^{\mathrm{th}}$ percentile.

As expected, the centralized graph coloring baseline shows significant performance superiority to all other schemes which are based on distributed execution. This difference is expected since CGC exploits the global information about the scenario, unlike the MADDQN, MAPPO, and greedy schemes which rely solely on local measurements of only the aggregate interference power. It should, however, be noted that the CGC method has significantly higher sensing, signalling, and computational complexity. Such high complexity makes CGC impractical for dense deployments of subnetworks.

### E. Sensitivity and Robustness Evaluation

We evaluate the ability of the proposed methods to generalize to changes in the deployment density and the wireless environment parameters.

*1) Robustness to changes in deployment density:* The models trained in the factory with $N = 20$ subnetworks are deployed in scenarios with different numbers of subnetworks. We vary the number of subnetworks in the test scenarios between $N = 10$ and $N = 50$ and evaluate the average per-device rate. As shown in Figure 6, when the number of subnetworks increases, an overall decrease in performance with similar proportions across all methods is observed. While the methods based on MADDQN appear to be robust to the changes in deployment density for all values of $N$, there seems to be a degradation in the performance of the methods based on MAPPO. This shows that the proposed F-MADDQN is more robust to changes in deployment density than the F-MAPPO.

*2) Sensitivity to changes in environment parameters:* To study the ability of the proposed methods to generalize to environments with different parameters than those used for the training, we train the models in a factory environment with sparse clutter with an average clutter element size of $d_{\text{clutter}} = 10$ m and a density of $r_{\text{clutter}} = 20\%$. The models are then introduced to environments with different types of clutter. We consider sparse and dense clutter cases with varying clutter density as defined in the 3GPP model for in-factory environment [20]. Figure 5 shows the minimum, maximum, and average values of the achieved rate per device for the different cases using the proposed schemes and the baseline algorithms. The figure shows that the RL-based methods can maintain their performance in comparison to the centralized and distributed baselines. Furthermore, similar to the observation from the sensitivity test with varying numbers of subnetworks, the MADDQN-based solutions appear to be slightly more robust to changes in the environment and maintain performance close to the greedy selection baseline.

## V. CONCLUSION

We proposed two federated reinforcement learning-based schemes named Federated Multi-Agent Double DQN (F-MADDQN) and Federated Multi-Agent Deep Proximal Policy Optimization (F-MADPPO) for dynamic channel allocation in 6G in-X subnetworks which overcome the inherent issues of convergence and high signalling overhead with privacy challenges associated with distributed and centralized multi-agent reinforcement learning techniques. Our performance evaluation results using realistic in-factory models defined by 5G-ACIA and 3GPP have shown that the proposed schemes can achieve similar performance to the best-performing baselines and are robust to changes in the deployment density as well as wireless environment conditions.

## REFERENCES

[1] R. Adeogun, G. Berardinelli, P. E. Mogensen, I. Rodriguez, and M. Razzaghpour, "Towards 6g in-x subnetworks with sub-millisecond communication cycles and extreme reliability," *IEEE Access*, vol. 8, pp. 110 172–110 188, 2020.

[2] H. Sun, X. Chen, Q. Shi, M. Hong, X. Fu, and N. D. Sidiropoulos, "Learning to optimize: Training deep neural networks for interference management," *IEEE Trans on Signal Proc.*, vol. 66, no. 20, pp. 5438–5453, 2018.

[3] R. Adeogun, G. Berardinelli, I. Rodriguez, and P. Mogensen, "Distributed dynamic channel allocation in 6g in-x subnetworks for industrial automation," in *2020 IEEE Globecom Workshops (GC Wkshps*, 2020, pp. 1–6.

[4] R. Adeogun, G. Berardinelli, and P. E. Mogensen, "Enhanced interference management for 6g in-x subnetworks," *IEEE Access*, vol. 10, pp. 45 784–45 798, 2022.

[5] G. Berardinelli and R. Adeogun, "Spectrum assignment for industrial radio cells based on selective subgraph constructions," in *2021 IEEE 94th Vehicular Technology Conference (VTC2021-Fall)*, 2021, pp. 01–05.

[6] R. O. Adeogun, G. Berardinelli, and P. E. Mogensen, "Learning to dynamically allocate radio resources in mobile 6g in-x subnetworks," 2021.

[7] D. Abode, R. Adeogun, and G. Berardinelli, "Power control for 6g industrial wireless subnetworks: A graph neural network approach," in *2023 IEEE WCNC.* IEEE, 2023, pp. 1–6.

[8] R. Adeogun and G. Berardinelli, "Multi-agent dynamic resource allocation in 6g in-x subnetworks with limited sensing information," *Sensors*, vol. 22, no. 13, p. 5062, 2022.

[9] X. Du, T. Wang, Q. Feng, C. Ye, T. Tao, L. Wang, Y. Shi, and M. Chen, "Multi-agent reinforcement learning for dynamic resource management in 6g in-x subnetworks," *IEEE transactions on wireless communications*, vol. 22, no. 3, pp. 1–1, 2023.

[10] B. Lu, H. Zhang, T. Xue, S. Guo, and H. Gai, "Deep reinforcement learning-based power allocation for ultra reliable low latency communications in vehicular networks," in *2021 IEEE/CIC International Conference on Communications in China (ICCC)*, 2021, pp. 1149–1154.

[11] W. Du and S. Ding, "A survey on multi-agent deep reinforcement learning: from the perspective of challenges and applications," *Artificial Intelligence Review*, vol. 54, pp. 3215–3238, 2021.

[12] A. Feriani and E. Hossain, "Single and multi-agent deep reinforcement learning for ai-enabled wireless networks: A tutorial," *IEEE Comm. Surveys & Tutorials*, vol. 23, no. 2, pp. 1226–1252, 2021.

[13] J. G. Kuba, R. Chen, M. Wen, Y. Wen, F. Sun, J. Wang, and Y. Yang, "Trust region policy optimisation in multi-agent reinforcement learning," *arXiv preprint arXiv:2109.11251*, 2021.

[14] L. Collins, H. Hassani, A. Mokhtari, and S. Shakkottai, "Fedavg with fine tuning: Local updates lead to representation learning," 2022.

[15] A. B. Mansour, G. Carenini, A. Duplessis, and D. Naccache, "Federated learning aggregation: New robust algorithms with guarantees," 2022.

[16] J. Qi, Q. Zhou, L. Lei, and K. Zheng, "Federated reinforcement learning: Techniques, applications, and open challenges," *arXiv preprint arXiv:2108.11887*, 2021.

[17] "5G-ACIA - 5G Alliance for Connected Industries and Automation — 5g-acia.org," http://www.5g-acia.org/, [Accessed 16-May-2023].

[18] R. Adeogun and G. Berardinelli, "Multi-agent dynamic resource allocation in 6g in-x subnetworks with limited sensing information," *Sensors*, vol. 22, no. 13, p. 5062, 2022.

[19] R. O. Adeogun, G. Berardinelli, and P. E. Mogensen, "Learning to dynamically allocate radio resources in mobile 6g in-x subnetworks," 2021.

[20] 3GPP, "Study on channel model for frequencies from 0.5 to 100 GHz," 3rd Generation Partnership Project (3GPP), Technical Report (TR) 38.901, 01 2017, version 16.1.0 Release 16.