# A Hybrid Optimization and Deep RL Approach for Resource Allocation in Semi-GF NOMA Networks

Duc-Dung Tran[1], Vu Nguyen Ha[1], Symeon Chatzinotas[1], and Ti Ti Nguyen[2]

[1]*Interdisciplinary Centre for Security, Reliability and Trust (SnT), University of Luxembourg, Luxembourg*
[2]*Université du Québec, Montréal, Canada*
Emails: [1]{duc.tran, vu-nguyen.ha, symeon.chatzinotas}@uni.lu, [2]titi.nguyen@emt.inrs.ca

*Abstract*—Semi-grant-free non-orthogonal multiple access (semi-GF NOMA) has emerged as a promising technology for the fifth-generation new radio (5G-NR) networks supporting the coexistence of a large number of random connections with various quality of service requirements. However, implementing a semi-GF NOMA mechanism in 5G-NR networks with heterogeneous services has raised several resource management problems relating to unpredictable interference caused by the GF access strategy. To cope with this challenge, the paper develops a novel hybrid optimization and multi-agent deep (HOMAD) reinforcement learning-based resource allocation design to maximize the energy efficiency (EE) of semi-GF NOMA 5G-NR systems. In this design, a multi-agent deep Q network (MADQN) approach is employed to conduct the subchannel assignment (SA) among users. While optimization-based methods are utilized to optimize the transmission power for every SA setting. In addition, a full MADQN scheme conducting both SA and power allocation is also considered for comparison purposes. Simulation results show that the HOMAD approach outperforms other benchmarks significantly in terms of the convergence time and average EE.

## I. INTRODUCTION

The future wireless networks are expected to be capable of serving a tremendous number of devices requiring heterogeneous services, e.g., enhanced mobile broadband (eMBB), ultra-reliable low-latency communications (URLLC), and massive machine type communications (mMTC), together with different quality-of-service (QoS) demands [1], [2]. In this context, semi-GF NOMA has been considered as a promising solution for relieving the heavy accessing-process overhead in the dense systems [3]. Following this strategy, the subchannels (SCs) are opened for mMTC users to access freely without waiting for receiving the admission granted, i.e., grant-free (GF) access, while the association process of other users having stringent QoS requirements (e.g., eMBB/URLLC users) are scheduled by the system controllers (such as base stations or access points, etc.), which is also called as grant-based (GB) access. In addition, the NOMA transmission can be exploited when there is more than one user accessing the same SC [4].

However, the without-admission-control property of the GF strategy may result in a serious congestion problem in semi-GF NOMA systems when a tremendously large number of devices tries to access a limited number of SCs. Therefore, GF access needs to be carefully designed to mitigate this problem as well as guarantee the QoS requirements of both GB and GF users in semi-GF NOMA systems. Furthermore, in real-time systems, developing a dynamic resource allocation (RA)

mechanism addressing the congestion problem and fulfilling the various QoS requirements from different services in semi-GF NOMA systems becomes more challenging. In recent years, reinforcement learning (RL) method has been applied to intelligently resolve the RA problem in communications [3]. Its application to GF NOMA and semi-GF NOMA systems has been investigated in [5]–[13]. However, these works have not considered the 5G-NR systems with the coexistence of multiple services. Furthermore, most of them aimed to discretize the continuous power variable to ease the learning process which may result in performance loss.

Regarding the drawback of the existing works, this paper develops two novel learning-based resource allocation designs maximizing EE while guaranteeing heterogeneous requirements relating to various communication services in semi-GF NOMA 5G-NR systems. Both of these proposed algorithms exploit the multi-agent deep RL method where the mMTC users are considered as agents that learn and optimize its SC and transmission power selection. The first algorithm, namely full multi-agent deep Q-network (Full-MAD), aims to set both SC assignment and power allocation (PA) as the action for the learning process, where the transmission power is quantized into a number of discrete levels. On a different method, the second algorithm, namely HOMAD, only considers the SC selection as the action model. In this learning-based solution, the transmission power corresponding to each SC setting can be determined by some efficient optimization-based analysis results. By doing so, the action space size is significantly degraded and the hybrid method can take advantage of both deep Q-network (DQN) and optimization-based approaches to gain better learning performance. The simulation results are then demonstrated to evaluate the performance of our proposed mechanisms in terms of convergence time and the system EE.

## II. SYSTEM MODEL

We investigate an uplink semi-GF NOMA 5G-NR network as shown in Fig. 1. The network consists of one BS located at the center of the cell with a radius of $r$ (m) and a number of users randomly distributed in this cell requiring different services including eMBB/mMTC/URLLC. Let $\mathcal{M}_U$, $\mathcal{M}_E$ and $\mathcal{M}_M$ be the sets of URLLC, eMBB, and mMTC devices, whose cardinalities are $M_U$, $M_E$ and $M_M$, respectively. For convenience, we also denote the set of all users as $\mathcal{M} = \mathcal{M}_U \cup \mathcal{M}_E \cup \mathcal{M}_M$ and $M = M_U + M_E + M_M$. To serve
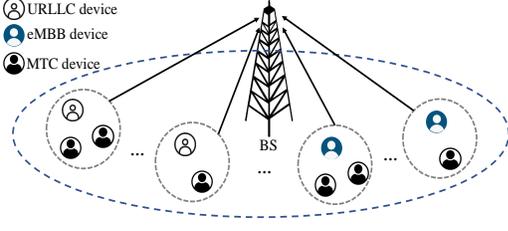
Fig. 1. Illustration of an uplink semi-GF NOMA 5G-NR network.

these users, a total bandwidth of $W$ (Hz) is assumed in the system, which is divided into $K$ SCs. Let $\mathcal{K}$ be the set of all SCs. Furthermore, to guarantee heterogeneous requirements of different services, a semi-GF NOMA transmission scheme is considered for the communication process.

### A. Uplink Semi-GF NOMA 5G-NR Transmission Strategy

*1) 5G-NR Numerology:* Following the 5G-NR standard which introduces various *"numerologies"*, physical-resource-block (PRB) or SC types, supporting different communication requirements, the bandwidth of SC in 5G-NR schemes is defined as $2^\nu$ times SC's bandwidth in 4G systems (i.e., 180 kHz), where $\nu \in \{0; 1; 2; 3; 4\}$ is the numerology index [1], [14]. Herein, PRBs with high SC spacing are arranged for URLLC services while traffic flows from the eMBB service can adopt a numerology with the smaller SC spacing [1]. Therefore, this paper focuses on an SC setting that the whole bandwidth is divided into two sets of SCs, $\mathcal{K}_\mathsf{U}$ and $\mathcal{K}_\mathsf{E}$. Particularly, $\mathcal{K}_\mathsf{U}$ represents the set of SCs serving URLLC users with numerology $\nu_\mathsf{U}$ while $\mathcal{K}_\mathsf{E}$ is the set of eMMB-service SCs with numerology $\nu_\mathsf{E}$. Herein, $\mathcal{K}_\mathsf{U} \cup \mathcal{K}_\mathsf{E} = \mathcal{K}$. One assumes that $\nu_\mathsf{E} < \nu_\mathsf{U}$ and denotes $W_\mathsf{E} = 2^{\nu_\mathsf{E}} \times 180$ (kHz) and $W_\mathsf{U} = 2^{\nu_\mathsf{U}} \times 180$ (kHz) as the bandwidth of SCs corresponding to eMBB and URLLC services, respectively.

*2) Semi-Grant-Free Radio Access Strategy:* In this system, the radio access of eMBB and URLLC users is managed by BS under the GB access scheme due to their requirements for high reliability, latency, and achievable rate. Specifically, each of these users is granted several distinct SCs for its transmission. In contrast, the mMTC users can access the network based on GF access method to improve connection density due to the massive access requirement of mMTC service. Herein, the mMTC users can access any SCs freely without a scheduling process to increase the access rate and the number of active mMTC users. In this context, many mMTC users can access the same SC; furthermore, they can use the SCs which are already granted to the eMBB and URLLC users.

Considering the transmission over SC $k$ ($k \in \mathcal{K}$), we denote $b_z^{(k)}(t)$ ($z \in \mathcal{M}$) as a binary SC allocation variable at time-slot (TS) $t$, where $b_z^{(k)}(t) = 1$ if device $z$ occupies SC $k$ and $b_z^{(k)}(t) = 0$ otherwise. In our scheme, we assume the orthogonal SC scheduled for URLLC/eMBB services and one-SC freely access strategy for mMTC users where each mMTC device can select only one arbitrary SC for its transmission. This assumption yields the following conditions,

$$(C1): \quad \sum_{z \in \mathcal{M}_\mathsf{U} \cup \mathcal{M}_\mathsf{E}} b_z^{(k)}(t) \leq 1, \quad \forall k \in \mathcal{K}. \quad (1)$$

$$(C2): \quad \sum_{k \in \mathcal{K}} b_z^{(k)}(t) = 1, \quad \forall z \in \mathcal{M}_\mathsf{M}. \quad (2)$$

In addition, the set of devices occupying SC $k$ in TS $t$ can be described as $\mathcal{Z}^{(k)}(t) = \{z | b_z^{(k)} = 1, z \in \mathcal{M}\}$.

*3) NOMA Transmission Mechanism:* In uplink NOMA, the decoding order of the multi-user data stream is affected by various different factors. Specifically, a decoding order can be formulated based on channel gain conditions [15], received power levels [5], or QoS constraints of users [16], [17]. In this paper, the messages of the users over each SC can be decoded at the BS as follows:

- Due to strict QoS requirements on reliability and latency, the URLLC user's signal needs to be decoded first.
- The symbols belonging to eMBB and mMTC users will be decoded in the order of the corresponding channel gains. In particular, the user having the higher channel gain will be decoded earlier at the BS.
- After decoding the message of a user with higher channel gain, the BS removes this component from its observation to decode the remaining users' messages by using the successive interference cancellation (SIC) technique.

Without loss of generality, one assumes there are $Z^k$ users accessing SC $k$ in TS, then they are arranged in the decoding order discussed above as $\mathcal{Z}^{(k)}(t) = \left\{ z_1^{(k)}, ..., z_{Z^k}^{(k)} \right\}$. Accordingly, the received signal-to-interference-plus-noise ratio (SINR) of user $z_\ell^{(k)}$ is expressed as

$$\gamma_{z_\ell^{(k)}}^{(k)}(t) = \mathcal{Y}_{z_l^{(k)}}^{(k)}(t) / \Big( \sum_{j > \ell} \mathcal{Y}_{z_j^{(k)}}^{(k)}(t) + \sigma_k^2 \Big), \quad (3)$$

where $\mathcal{Y}_z^{(k)}(t) = P_z^{(k)}(t) g_z^{(k)}(t)$ is the power of signal due to user $z$'s data over SC $k$ in TS $t$; $P_z^{(k)}(t)$ is the transmission power of user $z$ over SC $k$, in which $P_z^{(k)}(t) = 0$ if $b_z^{(k)}(t) = 0$ and $P_z^{(k)}(t) \neq 0$, otherwise; $g_z^{(k)}(t)$ denote the corresponding channel gain and $\sigma_k^2$ represents the noise power over SC $k$.

### B. Achievable Rate of Users

*1) URLLC Communication:* Regarding the transmission of URLLC user $u$ over SC $k$ in $\mathcal{K}_\mathsf{U}$, which happens when $b_u^{(k)} = 1$. Based on the NOMA transmission mechanism given in Section II-A, one must have $u \equiv z_1^{(k)}$. Moreover, the SINR of URLLC device $u$ over SC $k$ is expressed as

$$\gamma_u^{(k)}(t) = \mathcal{Y}_u^{(k)}(t) / \Big( \mathcal{I}_u^{(k)}(t) + \sigma_u^2 \Big), \quad (4)$$

where $\mathcal{I}_u^{(k)}(t) = \sum_{j=2}^{Z^k} \mathcal{Y}_{z_j^{(k)}}^{(k)}(t)$ represents the interference caused by mMTC users over SC $k$. Furthermore, bandwidth of SC $k$ in $\mathcal{K}_\mathsf{U}$ is $W_\mathsf{U}$ and $\sigma_u^2 = FN_0 W_\mathsf{U}$ denotes the noise power, where $F$ is the noise figure, $N_0$ is the noise power spectral density (PSD). Accordingly, the achievable rate of URLLC user $u$ over SC $k$ in finite blocklength regime for a quasi-static flat fading channel can be approximated as [18]

$$R_u^{(k)}(t) = W_\mathsf{U}[\log_2(1 + \gamma_u^{(k)}(t)) - \Phi_u^{(k)}(t)], \quad (5)$$

where $\Phi_u^{(k)}(t) = \sqrt{\frac{V_u^{(k)}(t)}{D_u W_\mathsf{U}}} \frac{Q^{-1}(\varepsilon_u)}{\ln 2}$, $V_u^{(k)}(t) = 1 - \left( 1 + \gamma_u^{(k)}(t) \right)^{-2} \approx 1$ [18] is the channel dispersion, $\varepsilon_u$ is

the decoding error probability (DEP) which can be used to evaluate the transmission reliability, $D_u$ is the transmission latency threshold, and $Q^{-1}(x)$ is the inverse of the Gaussian Q-function. Here, we define a data-rate demand for URLLC $u$ to satisfy the URLLC requirements (i.e., $\varepsilon_u$ and $D_u$) when transmitting one packet over one SC in each TS as $R_u^{\mathsf{tar}} = W_{\mathsf{U}} \left[ \log_2 \left( 1 + \gamma_u^{\mathsf{tar}} \right) - \Phi_u^{\mathsf{tar}} \right]$, where $\gamma_u^{\mathsf{tar}} = 2^{\frac{n_b}{D_u W_{\mathsf{U}}} + \frac{Q^{-1}(\varepsilon_u)}{\ln 2 \sqrt{D_{max} W_{\mathsf{U}}}}} - 1$ is the target SNR for user $u$ [18], $n_u$ is the packet size, and $\Phi_u^{\mathsf{tar}}$ is defined similarly as in (5). This demand yields the following constraints,

$$(C3): \quad b_u^{(k)}(t) R_u^{(k)}(t) \geq R_u^{\mathsf{tar}}, \quad \forall k \in \mathcal{K}. \quad (6)$$

*2) eMBB Communication:* Assume that eMBB user $e$ access SC $k$ in $\mathcal{K}_{\mathsf{E}}$ which implies $b_e^{(k)} = 1$. Due to its order in the NOMA-based decoding process, its SINR denoted as $\gamma_e^{(k)}(t)$, can be defined as in (3) with noting that $\sigma_k^2 = F N_0 W_{\mathsf{E}}$. Then, the achievable rate of eMBB device $e$ is given by

$$R_e^{(k)}(t) = W_{\mathsf{E}} \log_2 \left( 1 + \gamma_e^{(k)}(t) \right). \quad (7)$$

Herein, one addresses a predetermined target transmission rate, $R_e^{\mathsf{tar}}$, for each eMBB user $e$ in every TS as

$$(C4): \quad \sum_{k \in \mathcal{K}} b_e^{(k)}(t) R_e^{(k)}(t) \geq R_e^{\mathsf{tar}}, \quad \forall e \in \mathcal{M}_{\mathsf{E}}. \quad (8)$$

*3) mMTC Communication:* Based on the NOMA transmission strategy mentioned earlier in Section II-A, mMTC devices can select a free SC or the one occupied by either URLLC or eMBB device. When $b_m^{(k)}(t) = 1$, mMTC user $m$ utilize SC $k$ in TS $t$. In such case, the SINR of this device, denoted as $\gamma_m^{(k)}(t)$, can be calculated as in (3) with noting that $\sigma_k^2 = F N_0 W_k$ where $W_k = W_{\mathsf{E}}$ if $k \in \mathcal{K}_{\mathsf{E}}$, and $W_k = W_{\mathsf{U}}$, otherwise. Similar to URLLC devices, the achievable rate of mMTC device $m$ is given by $R_m^{(k)}(t) = W_k \left[ \log_2 (1 + \gamma_m^{(k)}(t)) - \Phi_m^{(k)}(t) \right]$, where $\Phi_m^{(k)}(t) = \sqrt{\frac{V_m^{(k)}(t)}{D_m W_k}} \frac{Q^{-1}(\varepsilon_m)}{\ln 2}$, $V_m^{(k)}(t) = 1 - \left( 1 + \gamma_m^{(k)}(t) \right)^{-2} \approx 1$ [18]. Furthermore, the target SNR of mMTC device $m$ can be defined as $\gamma_m^{\mathsf{tar}} = 2^{\frac{n_m}{D_m W_k} + \frac{Q^{-1}(\varepsilon_m)}{\ln 2 \sqrt{D_m W_k}}} - 1$, where $n_m$, $D_m$, and $\varepsilon_m$ denote the packet size, transmission latency, and DEP of mMTC device $m$. Then, the SINR of mMTC user should be greater than a threshold for successful decoding, i.e.,

$$(C5): \quad b_m^{(k)}(t) \gamma_m^{(k)}(t) \geq \gamma_m^{\mathsf{tar}}, \quad \forall m \in \mathcal{M}_{\mathsf{M}}. \quad (9)$$

*C. Energy Efficiency Maximization Problem*

In this paper, we aim to design an effective SC and power allocation strategy to maximize the network EE while guaranteeing the different requirements of all services. To do so, we first define an EE factor as follows:

$$\zeta(t) = R^{\mathsf{tot}}(t) / (P^{\mathsf{Tx}}(t) + M P_{\mathsf{c}}), \quad (10)$$

where $R^{\mathsf{tot}}(t) = \sum_{k \in \mathcal{K}} \sum_{z \in \mathcal{M}} b_z^{(k)}(t) R_z^{(k)}(t)$, $P^{\mathsf{Tx}}(t) = \sum_{k \in \mathcal{K}} \sum_{z \in \mathcal{M}} P_z^{(k)}(t)$, and $P_{\mathsf{c}}$ denotes the circuit power consumption. Then, the design problem can be formulated as

$$\max_{\mathbf{b}, \mathbf{P}} \mathbb{E}_t \left[ \zeta(t) \right] \quad \text{s.t.} \quad \text{constraints } (C1) - (C5), \quad (11a)$$

$$(C6): \sum_{k \in \mathcal{K}} P_z^{(k)}(t) \leq P_z^{\mathsf{max}}, \; \forall (z, t), \quad (11b)$$

where $\mathbf{b}$ and $\mathbf{P}$ denote the SC assignment and power control strategies, respectively; and constraint $(C6)$ stands for the power budget of devices.

## III. TWO PROPOSED MULTI-AGENT DEEP RL SOLUTIONS

*A. Full multi-agent DQN Approach*

A full multi-agent DQN approach, named Full-MAD, is first studied in this section. Herein, all mMTC users are considered agents. Employing a multi-agent deep RL mechanism, they separately learn and define optimal policies for selecting SC and PA. In addition, the multi-level quantization strategy is exploited to deal with the continuous characteristic of power variables in the similar approach introduced in [5], [10]. Herein, the power is quantized into $L$ levels to build the action sets for the RL process. Particularly, the state, action, and reward of each agent ( e.g., $m \in \mathcal{M}_{\mathsf{M}}$) in TS $t$ are defined as follows. The state of agent $m$ is defined as

$$s_m(t) = \left\{ g_m^{(1)}(t), \ldots, g_m^{(K)}(t), \mathbf{c}_m(t-1) \right\}, \quad (12)$$

where $\mathbf{c}_m(t-1) = \{ b_m^{(1)}(t), \ldots, b_m^{(K)}(t), P_m^{(1)}(t), \ldots, P_m^{(K)} \}$ is the SC and transmission power level (TPL) selection status of agent $m$ in TS $t-1$. Additionally, one defines the action of agent $m$ as its SC and TPL selection, expressed as

$$a_m(t) = \{ 1, \ldots, kl, \ldots, KL \}, \quad (13)$$

where $a_m(t) = kl$ indicates that agent $m$ select SC $k$ and the $l$-th TPL in TS $t$. Let $\mathcal{A}_m$ be the set of all actions of mMTC user $m$, we have $|\mathcal{A}_m| = KL$. For action selection strategy, the $\epsilon$-greedy policy can be exploited where the random action is taken with the probability of $\epsilon$, and the action with the highest Q-value, i.e., $a_m^{max} = \arg \max_{a \in \mathcal{A}_m} \{ Q_m (s_m(t), a; \boldsymbol{\theta}_m) \}$ is employed for the remaining probability. Herein, $Q_m (s_m(t), a_m(t); \boldsymbol{\theta}_m)$ is the Q-value corresponding to action $a_m(t)$. Regarding the EE factor, we can define the reward function in TS $t$ as

$$r(t) = \begin{cases} \zeta(t), & \text{if all constraints are satisfied,} \\ 0, & \text{otherwise.} \end{cases} \quad (14)$$

Based on the actions and rewards obtained from trials, each agent builds its own DQN model consisting of two deep neural networks (DNNs), namely online and target networks corresponding to weight vectors $\boldsymbol{\theta}_m$ and $\boldsymbol{\theta}_m'$, respectively. Herein, the online network is used to select an action. Meanwhile, the target network is applied to evaluate the online network-based action. Thus, the objective is to reduce the loss function as [10]

$$\hat{L}(\boldsymbol{\theta}_m) = \left[ y_m(t) - Q_m (s_m(t), a_m(t); \boldsymbol{\theta}_m) \right]^2, \quad (15)$$

where $y_m(t)$ denotes the target Q-value determined by the target network as $y_m(t) = r(t) + \max_{a \in \mathcal{A}_m} Q_m (s_m(t+1), a; \boldsymbol{\theta}_m')$. Given the MDP principle and DQN model of each agent mentioned above, the Full-MAD learning-based solution approach is summarized in Algorithm 1.

**Algorithm 1** MULTI-AGENT DRL-BASED ENERGY EFFICIENCY MAXIMIZATION ALGORITHM

1: Initialize the weight vectors of the online and target networks, i.e., $\boldsymbol{\theta}_m$ and $\boldsymbol{\theta}'_m$, $\forall m \in \mathcal{M}_\mathsf{M}$.
2: **for** $i = 1, \ldots, E_p$ **do**
3:  Initialize the state $s_m(t)$, $\forall m$.
4:  **for** $t = 1, \ldots, T$ **do**
5:   All agents take actions as
   - **For Full-MAD approach:** *Both SC and power-level selection* (as in (13)) following the $\epsilon$-greedy policy.
   - **For HOMAD approach:** *Only SC selection* (as in (16)) following the $\epsilon$-greedy policy. *Optimizing transmission power as in Section III-B.*
6:   All agents observe the reward in (14) and move to the next states.
7:   **for** $m = 1, \ldots, M_\mathsf{M}$ **do**
8:    Store an experience tuple $(s_m(t), a_m(t), r(t), s_m(t+1))$ to the memory of agent $m$.
9:    Randomly sample a mini-batch of experiences from the memory to train the online network.
10:    Update $\boldsymbol{\theta}_m$ by using gradient descent to minimize the loss function in (15).
11:    Update $\boldsymbol{\theta}'_m$ as $\boldsymbol{\theta}'_m = \boldsymbol{\theta}_m$ after every $B$ TSs.
12:   **end for**
13:  **end for**
14: **end for**

### B. Hybrid Optimization and Multi-Agent DQN Approach

In this section, we developed the HOMAD approach to speed up the learning process and eliminate the power-quantization loss. Specifically, the HOMAD mechanism employs a MADQN-based method for SC selection according to which the transmission power is optimized efficiently. In this approach, the states and rewards are defined similarly to those presented in the Full-MAD method while the action is simplified to SC selection as

$$a_m(t) = \{1, \ldots, k, \ldots, K\}. \quad (16)$$

Once the action is taken, the transmission power is optimized as follows. Firstly, employing the Dinkelbach algorithm [19], [20], we aim to solve problem (11) by iteratively solving,

$$\max_{\mathbf{P}(t)} R^{\mathsf{tot}}(t) - \zeta P^{\mathsf{Tx}}(t) \text{ s.t. } (C1) - (C6), \quad (17)$$

and adjusting $\zeta$ until an optimal $\zeta^\star \geq 0$ satisfying $R^{\mathsf{tot}}(t) = \zeta^\star \left( P^{\mathsf{Tx}}(t) + M P_\mathsf{c} \right)$ is found. To cope with (17), let's regard the following remark based on which we propose efficient approaches to optimize the power for each SC setting.

**Remark 1.** *The formula given in* (3) *demonstrates that there is no interference suffering the decoding process due to user $z_{Z^k}^{(k)}$. Moreover, once the power of all users in set $\left\{\ell+1, \ldots, z_{Z^k}^{(k)}\right\}$ is defined, the transmission power of user $z_\ell^{(k)}$, i.e., $P_{z_\ell^{(k)}}^{(k)}$, can be optimized without coupling to other users.*

*1) Power Allocation for mMTC Service:* Once user $z_\ell^{(k)}$ is an mMTC user, $P_\ell^{(k)}$ can be optimized for given $\{P_m^{(k)}\}_{m=\ell+1}^{Z_k}$ by solving the following sub-problem.

$$\max_{p_\ell} R_m^{(k)}(t) - \zeta p_\ell \text{ s.t. } \gamma_{z_\ell^{(k)}}^{\mathsf{tar}}/A_\ell^{(k)} \leq p_\ell \leq P_{z_\ell^{(k)}}^{\mathsf{max}}, \quad (18)$$

where $A_\ell^{(k)} = g_{z_\ell^{(k)}}^{(k)}(t)/(\sum_{j>\ell} g_{z_j^{(k)}}^{(k)}(t) P_{z_j^{(k)}}^{(k)} + \sigma_k^2)$.

**Proposition 1.** *The solution of problem* (18) *is given as,*

$$P_\ell^{(k)\star} = \min(\max(W_k/(\zeta \ln 2) - 1/A_\ell^{(k)}, \gamma_{z_\ell^{(k)}}^{\mathsf{tar}}/A_\ell^{(k)}), P_{z_\ell^{(k)}}^{\mathsf{max}}). (19)$$

*Proof:* The proof is described simply as follows. As can be seen, problem (18) is convex due to the concave objective function and the convex feasible set. Then, the optimal solution can be obtained by setting the derivative of the objective function to zero and solving it with the feasible set [21]. ■

*2) Power Allocation for eMBB Service:* Assume that user $e$ is assigned $n$ SCs named as $\{k_1^e, \ldots, k_n^e\} \subset \mathcal{K}_\mathsf{E}$, and it is denoted as user $z_{\ell_j}^{(k_j^e)}$ over SC $k_j^e$ ($j = 1, \ldots, n$). Then, problem (17) can be decomposed for user $e$ as

$$\max_{\mathbf{P}^e} \sum_{j=1}^n C_j^e - \zeta p_j^e \text{ s.t.} \sum_{j=1}^n C_j^e \geq R_e^{\mathsf{tar}}, \sum_{j=1}^n p_j^e \leq P_e^{\mathsf{max}}, \quad (20)$$

where $C_j^e = W_\mathsf{E} \log_2 \left(1 + A_j^e p_j^e\right)$, $\mathbf{P}^e = [p_1^e, \ldots, p_n^e]$, $p_j^e$ is the transmission power variable of eMBB user $e$ over SC $k_j^e$, and $A_j^e$ is defined similarly as in (18). Since problem (20) is convex, its solution can be obtained by using the duality method. In particular, the Lagrangian of (20) is described as $\mathcal{L}(\mathbf{P}^e, \mu, \nu) = \sum_{j=1}^n \left[(1+\mu)C_j^e - (\zeta + \nu)p_j^e\right] - \mu R_e^{\mathsf{tar}} + \nu P_e^{\mathsf{max}}$, where $\mu$ and $\nu$ are the Lagrangian multipliers corresponding to the constrains of (20). Then, the dual function is defined as $\mathsf{g}(\mu, \nu) = \max_{\mathbf{P}^e} \mathcal{L}(\mathbf{P}^e, \mu, \nu)$.

**Proposition 2.** *The solution of dual function is defined as*

$$p_j^e = \max \left((1+\mu)W_\mathsf{E}/[(\nu + \zeta) \ln 2] - 1/A_j^e, 0\right). \quad (21)$$

*Proof:* The proof of this proposition can be obtained easily by solving the equation $\partial \mathcal{L}(\mathbf{P}^e, \mu, \nu)/\partial p_j^e = 0$. ■

The dual problem can be rewritten as $\max_{\mu, \nu} \mathsf{g}(\mu, \nu)$ s.t. $\mu, \nu \geq 0$. Since problem (20) is convex and the dual gap is zero, the optimal solution of the dual problem can be found by iteratively updating the dual variables $\mu$ and $\nu$ as $\mu^{[v+1]} = \left[\mu^{[v]} - \delta^{[v]}(\sum_{j=1}^n C_j^e - R_e^{\mathsf{tar}})\right]^+$ and $\nu^{[v+1]} = \left[\nu^{[v]} + \delta^{[v]}(\sum_{j=1}^n p_j^e - P_e^{\mathsf{max}})\right]^+$, where the suffix $[v]$ represents the iteration index, $\delta_{[v]}$ is the step size. This sub-gradient method guarantees the convergence if $\delta_{[v]} \overset{v \to \infty}{\longrightarrow} 0$ [22].

*3) Power Allocation for URLLC Service:* Similar to the previous section, one assumes that there are $l$ SCs assigned to URLLC user $u$, namely $\{k_1^u, \ldots, k_l^u\} \subset \mathcal{K}_\mathsf{U}$. Then, if the power of all mMTC users on SCs $\{k_1^u, \ldots, k_l^u\}$ are determined, the power transmission over all SCs can be determined by solving the following problem

$$\max_{\mathbf{P}^u} \sum_{j=1}^l \left(C_j^u - \Phi_j^u\right) - \zeta p_j^u \text{ s.t. } p_j^u \geq \gamma_u^{\mathsf{tar}}/A_j^u, \ \forall j, \quad (22a)$$

$$\sum_{j=1}^l p_j^u \leq P_u^{\mathsf{max}}, \quad (22b)$$

where $C_j^u = W_\mathsf{U} \log_2 \left(1 + A_j^u p_j^u\right)$, $\Phi_j^u = \sqrt{\frac{V_j^u}{D_u W_\mathsf{U}}} \frac{Q^{-1}(\varepsilon_j^u)}{\ln 2}$, $V_j^u = 1 - \left(1 + A_j^u p_j^u\right)^{-2} \approx 1$ [18], $\mathbf{P}^u = [p_1^u, \ldots, p_l^e]$ and $p_j^u$ denotes the transmission power variable corresponding to URLLC user $u$ over SC $k_j^u$. Similar to the approach employed for solving problem (20), the transmission power of URLLC users can be determined in the following proposition.

**Algorithm 2** ENERGY-EFFICIENCY POWER ALLOCATION ALGORITHM

1: Initialize $\zeta^{(0)} = 0$, set $q = 0$, and choose predetermined tolerate $\tau$.
2: **repeat**
3:   The power allocation can be optimized in a parallel manner over all SCs for mMTC users, but the process over an SC involved in an eMBB user process can stop and then continue when the power of that eMBB user is updated. The process is described as
   a. **The power of every mMTC user** are defined as in (19).
   b. **The power of every eMBB user** is optimized as described in Section III-B2 when all mMTC ordered before it over the corresponding SCs having their transmission power optimized.
   c. **The power of every URLLC user** is optimized as in Section III-B3 when all mMTC users have their power transmission determined.
4:   Update $\zeta^{(q+1)} = \frac{R^{\text{tot}}(t)}{P^{\text{Tx}}(t) + MP_\text{c}}$.
5:   Set $q := q + 1$.
6: **until** $|\zeta^{(q)} - \zeta^{(q-1)}| \leq \tau$.

TABLE I
EXPERIMENTAL PARAMETERS

| Parameters | Value |
|---|---|
| Cell radius ($r$) | 500 m |
| Channel model | Rician |
| eMBB, URLLC numerology indices ($\nu_\text{E}, \nu_\text{U}$) | 1, 4 |
| eMBB Data-rate demand $\left(R_e^{\text{tar}}\right)$ | {2; 4; 6; 8} bps/Hz |
| Latency threshold ($D_u = D_m = D_{max}$) | {2; 1; 0.5; 0.4} ms |
| Reliability threshold $\left(\varepsilon_u = \varepsilon_m = \varepsilon_{th} = 10^{-x}\right)$ | $x = \{2; 4; 5, 6, 7\}$ |
| Maximum transmission power | 23 dBm |
| Circuit power consumption ($P_c$) | 0.05 W |
| Noise figure and PSD ($F$ and $N_0$) | 6 dB and -174 dBm/Hz |
| Packet length ($n_u = n_m = n_b$) | 32 bytes |
| Number of hidden layers, neurons per hidden layers | 3, {256, 128, 64} |
| Learning rate ($\alpha$) and discount factor ($\gamma$) | 0.001 and 0.9 |
| Optimizer | Adam |

**Proposition 3.** *The transmission power of URLLC users $u$ over SCs $\{k_1^u, ..., k_l^u\}$, can be defined as*

$$p_j^u = \max\left(W_\text{U}/\left[(\theta + \zeta)\ln 2\right] - 1/A_j^u, \bar{\gamma}_u/A_j^u\right), \forall j, \quad (23)$$

*where $\theta$ is iteratively updated as $\theta^{[v+1]} = \left[\theta^{[v]} + \delta^{[v]}(\sum_{j=1}^{l} p_j^u - P_u^{\text{max}})\right]^+$.*

*Proof:* The proof can be obtained by employing the similar duality method presented in Section III-B2. ∎

In summary, the proposed energy-efficiency power allocation algorithm is described in Algorithm 2.

*4) Proposed HOMAD Algorithm:* As mentioned above, the HOMAD approach allows agents to select SCs by using a MADQN scheme similar to Full-MAD approach whereas the transmission power for each SC setting is optimized by applying Algorithm 2. The summary of this algorithm is also provided in Algorithm 1 with **"HOMAD" remark** in **Step 5**.

## IV. SIMULATION RESULTS

This section provides the simulation results to evaluate our proposed algorithms' performance. The simulations were performed on a PC equipped an Intel Xeon W-11855M CPU with 3.2 GHz frequency, 64-GB RAM, and 64-bit Windows 10 operating system. The DQN model consists of three fully-connected hidden layers including 256, 128, and 64 neurons. The experimental parameters are provided in Table I.

Fig. 2 depicts the convergence trend during the learning process of full multi-agent Q-learning (Full-MAQL), Full-
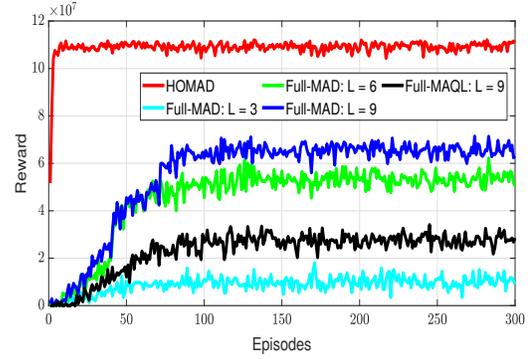


Fig. 2. Convergence comparison of different learning methods, where $M_\text{U} = M_\text{E} = 1$, $M_\text{M} = 4$, $K = 2$, $\nu_\text{E} = 1$, $\nu_\text{U} = 4$, $D_{max} = 2$ (ms), $\varepsilon_{th} = 10^{-5}$, and $R_e^{tar} = 4$ (bps/Hz).

MAD, and HOMAD approach by illustrating the variation of the reward achieved by all agents versus the various number of episodes. In the Full-MAQL scheme, each agent needs to build its own Q-table including all possible sate-action combinations. Fig. 2 shows that Full-MAQL returns the lowest reward (i.e., worst performance) in comparison to other schemes. This demonstrates the limitation of the Q-learning approach in a very large-space environment. Considering our proposed schemes, the HOMAD algorithm can achieve the highest rewards in this simulation with the significant gaps between its corresponding curve and the others thanks to the power-optimization process. Due to the discrete power level of the quantization process, which is widely used in literature [5], [10], the Full-MAD scheme obtains a lower reward than that due to the HOMAD. Interestingly, the Full-MAD scheme can improve its performance by increasing the number of TPLs ($L$) as shown in Fig. 2 which however leads to a larger action space together with a higher complexity level of the learning process. The convergence is further clarified in Table II where the number of episodes required for convergence, the average implementing time per episode, and the convergence time corresponding to three schemes are provided. As given in this table, HOMAD scheme converges with the lowest number of episodes but also requires the highest implementing time per episode because of its smallest action space size and also the power-optimization process. Inversely, Full-MAQL needs the largest number of episodes for convergence while its average time for each episode is the shortest. In summary, the HOMAD scheme again shows its superiority to the others when requires the shortest time for convergence.

Next Figs. 3 and 4 illustrate the variation of the average EE versus the different value sets of $(D_{max}, \varepsilon_u, R_e^{\text{tar}})$ and number of mMTC users, respectively. In Fig. 3, we aim to consider a scenario, where mMTC users can use the same SCs assigned to URLLC and eMBB users to improve spectral efficiency and connectivity density, as long as URLLC and eMBB requirements, i.e., $(D_{max}, \varepsilon_u, R_e^{\text{tar}})$, are still guaranteed. As expected, the higher stringent requirements of URLLC and eMBB users (i.e., lower $D_{max}$ and $\varepsilon_u$, and higher $R_e^{\text{tar}}$) result in the lower EE achieved by all schemes. In Fig. 4, one can observe that increasing $M_\text{M}$ will degrade the system EE. This

TABLE II
CONVERGENCE TIME COMPARISON

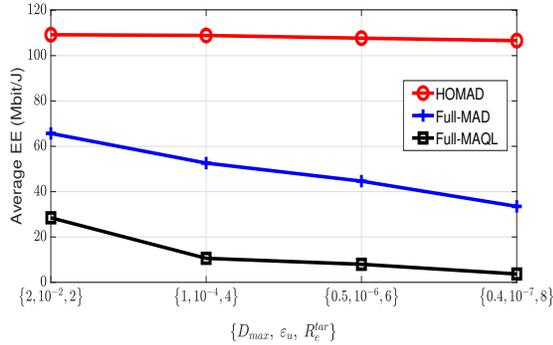| Methods | Avg. time per episode | No. of episodes | Conv. time |
|---------|----------------------|-----------------|------------|
| Full-MAQL | 0.36 sec. | 85 | 31.05 sec. |
| Full-MAD | 1.58 sec. | 85 | 134.46 sec. |
| HOMAD | 2.49 sec. | 6 | 14.93 sec. |



Fig. 3. Effect of URLLC and eMBB requirements $\{D_{max}, \varepsilon_u, R_e^{tar}\}$, where $M_U = M_E = 1$, $M_M = 4$, $K = 2$, $\nu_E = 1$, and $\nu_U = 4$.
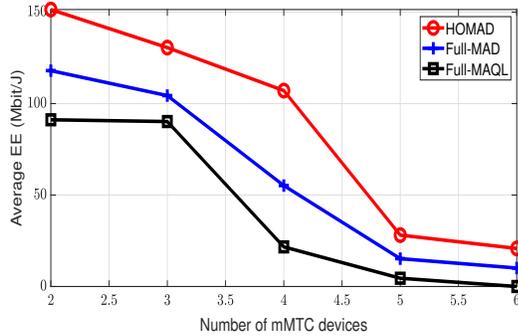


Fig. 4. Effect of $M_M$, where $M_U = M_E = 1$, $K = 2$, $\nu_E = 1$, $\nu_U = 4$, $D_{max} = 2$ (ms), $\varepsilon_{th} = 10^{-5}$, and $R_e^{tar} = 4$ (bps/Hz).

is because the number of mMTC users using the same SC gets higher as $M_M$ increases which results in the higher interference suffering the URLLC and eMBB users. In addition, these figures again confirm the superiority of the proposed HOMAD algorithm in all simulation scenarios, while the Full-MAD algorithm outperforms the Full-MAQL scheme.

## V. CONCLUSION

We have proposed two multi-agent Deep RL-based resource allocation mechanisms, HOMAD and Full-MAD algorithms, for maximizing the system EE of eMBB/mMTC/URLLC-coexistence 5G-NR networks using semi-GF NOMA transmission strategy. In particular, the Full-MAD approach addresses the EE maximization problem by employing the MADQN method to conduct both SC and PA selection. Furthermore, the HOMAD approach aims to use the MADQN method to only select SC solution while the power corresponding to a given SC setting can be optimized effectively. Simulation results have shown that the Full-MAD method outperforms the conventional Full-MAQL mechanism, while the HOMAD algorithm can return a higher EE and converge faster than other benchmark schemes.

## REFERENCES

[1] S.-Y. Lien *et al.*, "5G new radio: Waveform, frame structure, multiple access, and initial access," *IEEE Commun. Mag.*, vol. 55, no. 6, pp. 64–71, 2017.

[2] T. T. Nguyen, V. N. Ha, and L. B. Le, "Wireless scheduling for heterogeneous services with mixed numerology in 5G wireless networks," *IEEE Commun. Lett.*, vol. 24, no. 2, pp. 410–413, 2020.

[3] M. B. Shahab, R. Abbas, M. Shirvanimoghaddam, and S. J. Johnson, "Grant-free non-orthogonal multiple access for IoT: A survey," *IEEE Commun. Surveys Tuts.*, vol. 22, no. 3, pp. 1805–1838, 2020.

[4] D.-D. Tran, H.-V. Tran, D.-B. Ha, and G. Kaddoum, "Secure transmit antenna selection protocol for MIMO NOMA networks over Nakagami-m channels," *IEEE Syst. J.*, vol. 14, no. 1, pp. 253–264, 2020.

[5] M. V. da Silva, R. D. Souza, H. Alves, and T. Abrão, "A NOMA-based Q-learning random access method for machine type communications," *IEEE Wireless Commun. Lett.*, vol. 9, no. 10, pp. 1720–1724, 2020.

[6] D.-D. Tran, S. K. Sharma, S. Chatzinotas, and I. Woungang, "Learning-based multiplexing of grant-based and grant-free heterogeneous services with short packets," in *Proc. IEEE Global Commun. Conf.*, 2021.

[7] ——, "Q-learning-based SCMA for efficient random access in mMTC networks with short packets," in *Proc. IEEE Int. Symp. Pers., Indoor, Mobile Radio Commun.*, 2021, pp. 1334–1338.

[8] D.-D. Tran, S. K. Sharma, and S. Chatzinotas, "BLER-based adaptive Q-learning for efficient random access in NOMA-based mMTC networks," in *Proc. IEEE Veh. Technol. Conf.*, 2021, pp. 1–5.

[9] D.-D. Tran, V. N. Ha, and S. Chatzinotas, "Novel reinforcement learning based power control and subchannel selection mechanism for grant-free NOMA URLLC-enabled systems," in *IEEE VTC-Spring*, 2022.

[10] M. Fayaz, W. Yi, Y. Liu, and A. Nallanathan, "Transmit power pool design for grant-free NOMA-IoT networks via deep reinforcement learning," *IEEE Trans. Wireless Commun.*, vol. 20, no. 11, 2021.

[11] ——, "A power-pool-based power control in semi-grant-free NOMA transmission," *arXiv preprint arXiv:2106.11190v2*, pp. 1–14, 2022.

[12] Y. Liu, Y. Deng, H. Zhou, M. Elkashlan, and A. Nallanathan, "Deep reinforcement learning-based grant-free NOMA optimization for mURLLC," *IEEE Trans. Commun.*, vol. 71, no. 3, pp. 1475–1490, 2023.

[13] D.-D. Tran, S. K. Sharma, V. N. Ha, S. Chatzinotas, and I. Woungang, "Multi-agent DRL approach for energy-efficient resource allocation in URLLC-enabled grant-free NOMA systems," *IEEE Open J. Commun. Soc.*, pp. 1–17, 2023, early access.

[14] V. N. Ha, T. T. Nguyen, L. B. Le, and J.-F. Frigon, "Admission control and network slicing for multi-numerology 5G wireless networks," *IEEE Netw. Lett.*, vol. 2, no. 1, pp. 5–9, 2020.

[15] S. Han *et al.*, "Energy-efficient short packet communications for uplink NOMA-based massive MTC networks," *IEEE Trans. Veh. Technol.*, vol. 68, no. 12, pp. 12066–12078, December 2019.

[16] Z. Ding, R. Schober, and H. V. Poor, "Unveiling the importance of SIC in NOMA systems-part 1: State of the art and recent findings," *IEEE Commun. Lett.*, vol. 24, no. 11, pp. 2373–2377, 2020.

[17] D.-D. Tran, S. K. Sharma, S. Chatzinotas, I. Woungang, and B. Ottersten, "Short-packet communications for MIMO NOMA systems over Nakagami-m fading: BLER and minimum blocklength analysis," *IEEE Trans. Veh. Technol.*, vol. 70, no. 4, pp. 3583–3598, 2021.

[18] C. Sun *et al.*, "Optimizing resource allocation in the short blocklength regime for ultra-reliable and low-latency communications," *IEEE Trans. Wireless Commun.*, vol. 18, no. 1, pp. 402–415, Jan. 2019.

[19] V. N. Ha, D. H. N. Nguyen, and J.-F. Frigon, "System energy-efficient hybrid beamforming for mmwave multi-user systems," *IEEE Trans. Green Commun. Netw.*, vol. 4, no. 4, pp. 1010–1023, 2020.

[20] ——, "Energy-efficient hybrid precoding for mmwave multi-user systems," in *Proc. IEEE Int. Conf. Commun.*, 2018, pp. 1–6.

[21] V. N. Ha, G. Kaddoum, and G. Poitau, "Joint radio resource management and link adaptation for multicasting 802.11ax-based WLAN systems," *IEEE Trans. Wireless Commun.*, vol. 20, no. 9, pp. 6122–6138, 2021.

[22] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge University Press, 2004.