# Hierarchical Multi-Agent DRL-Based Framework for Joint Multi-RAT Assignment and Dynamic Resource Allocation in Next-Generation HetNets

Abdulmalik Alwarafy[1], Bekir Sait Ciftler[1], *Member, IEEE,* Mohamed Abdallah[1], *Senior Member, IEEE,*
Mounir Hamdi[1], *Fellow Member, IEEE,* and Naofal Al-Dhahir[2], *Fellow Member, IEEE*

[1]Division of Information and Computing Technology, College of Science and Engineering,
Hamad Bin Khalifa University, Doha, Qatar
[2]Electrical and Computer Engineering Department, Erik Jonsson School of Engineering and Computer Science,
The University of Texas at Dallas, USA

*Abstract*—This paper considers the problem of cost-aware downlink sum-rate maximization via joint optimal radio access technologies (RATs) assignment and power allocation in next-generation heterogeneous wireless networks (HetNets). We consider a future HetNet comprised of multi-RATs and serving multi-connectivity edge devices (EDs), and we formulate the problem as a mixed-integer non-linear programming (MINP) problem. Due to the high complexity and combinatorial nature of this problem and the difficulty to solve it using conventional methods, we propose a hierarchical multi-agent deep reinforcement learning (DRL)-based framework, called DeepRAT, to solve it efficiently and learn system dynamics. In particular, the DeepRAT framework decomposes the problem into two main stages; the *RATs-EDs assignment stage*, which implements a single-agent Deep $Q$ Network (DQN) algorithm, and the *power allocation stage*, which utilizes a multi-agent Deep Deterministic Policy Gradient (DDPG) algorithm. Using simulations, we demonstrate how the various DRL agents efficiently interact to learn system dynamics and derive the global optimal policy. Furthermore, our simulation results show that the proposed DeepRAT algorithm outperforms existing state-of-the-art heuristic approaches in terms of network utility. Finally, we quantitatively show the ability of the DeepRAT model to quickly and dynamically adapt to abrupt changes in network dynamics, such as EDs' mobility.

*Index Terms*—Deep Reinforcement Learning, Deep Q Network, Deep Deterministic Policy Gradient, Resource Allocation, Multi-RAT Assignment, Power Allocation, Heterogeneous Networks.

## I. INTRODUCTION

Heterogeneous wireless networks (HetNets) are expected to be one of the key enablers for next-generation wireless communication networks [2]. In such networks, a massive number of multi-radio access technologies (multi-RATs) in the licensed and unlicensed frequency bands across the ground, space, and underwater coexist to enhance the network's quality of service (QoS). The main goal of next-generation HetNets is to support the stringent QoS requirements of the emerging disruptive wireless applications in terms of rate, coverage, and reliability with ubiquitous connectivity. On the other hand, emerging user edge devices (EDs) are equipped with advanced multi-access physical capabilities that enable them to simultaneously aggregate radio resources from various RATs (i.e.,

multi-homing mode of operation) to guarantee an enhanced reliable communication for their running applications [3]. The number of such EDs is expected to be around 30 Billion by 2023, including smartphones, IoT devices, and sensors [4].

Radio resource allocation is crucial in the planning, orchestration, and resource optimization of next-generation HetNets. It is mainly used to guarantee enhanced system efficiency, increased network connectivity, and reduced energy consumption. However, allocating and managing radio resources, such as power, spectrum, and rate, in the next-generation HetNets is a persistent challenge. In particular, the multi-RAT assignment for EDs (i.e., RATs-EDs associations) and RATs' power allocation are among the key issues in the emerging HeNets. Various methods have been proposed to address these two radio resource allocation issues, such as optimization theory, ranking-based, and game theory [5]–[7]. However, most of these conventional resource allocation methods generally suffer from the following shortcomings. They require full and real-time knowledge of the network dynamics. Unfortunately, obtaining such information is not possible as most real-world networks are dynamic and change over time causing a rapid variation of the wireless channel. In addition, most of these conventional resource allocation methods suffer from high computational complexity, lack of scalability, and do not always guarantee convergence. These issues are even exacerbated in next-generation HetNets due to the large-scale and massive heterogeneous nature of these networks in terms of the underlying RATs and QoS demands of supported applications as well as the explosive number and type of emerging EDs. All of the above mentioned issues and challenges render the use of existing resource allocation techniques for future wireless HetNets quite difficult if not even impossible [8]. Hence, it is of a paramount importance to develop alternative resource allocation solutions that can overcome these challenges while quickly adapting to the varying systems' dynamics.

Deep reinforcement learning (DRL) has emerged recently as one of the most promising branches of the artificial intelligence (AI) field. In DRL, intelligent agents are trained to make autonomous decisions and observe their results in order to

A conference version of this work was published in the IEEE International Conference on Communications Workshops (ICC Workshops) 2021 Proceedings [1].

learn optimal control policies. In the context of radio resource allocation, DRL methods possess several advantages over state-of-the-art techniques. First, they provide autonomous and real-time decision-making even for highly complex and large-scale HetNets. Second, DRL provides efficient solutions for complex and high-dimensional wireless radio resource allocation optimization problems with limited channel state information (CSI) knowledge. These unique features make DRL techniques one of the key enabling technologies that can be utilized to address the radio resource allocation in next-generation HetNets.

In light of the imperative need for radio resource allocation in next-generation HetNets, the shortcomings of traditional radio resource allocation approaches, and the efficiency of DRL techniques in solving complex radio resource allocation optimization problems, this paper proposes a DRL-based framework for radio resource allocation in next-generation HetNets. In particular, we propose a hierarchical DQN and DDPG-based scheme called DeepRAT to study the problem of adaptive multi-RAT assignment and continuous downlink power allocation for multi-homing EDs in next-generation HetNets. Our aim is to jointly optimize the sum rate of the entire network, monetary cost, and power allocation while satisfying the QoS demands of EDs and the constraints of the RATs' power resources. Our simulation results demonstrate that the DeepRAT algorithm can efficiently learn the optimal policy, and it outperforms the greedy, random, and fixed algorithms in terms of satisfying the objective of the optimization problem. Moreover, after training, our proposed algorithm converges 2.5 times faster than the case with the initial training. In general, the main contributions of this paper are summarized as follows:

- We formulate an optimization problem whose objective is to cost-effectively maximize the downlink sum-rate of the multi-RAT HetNet via jointly optimizing the RATs-EDs assignment and RATs' power allocations while considering the limited RATs' power resources, multi-homing capabilities of EDs, and QoS data requirements of EDs.
- Due to the extensive computational complexity and combinatorial nature of the formulated problem, as well as the difficulty of applying conventional approaches to solve it, we propose a DRL-based algorithm called DeepRAT to solve the problem hierarchically and learn the system dynamics using a mix of value-based and policy-based DRL algorithms.
- Using simulations, we show how the various agents of the DeepRAT model interact in order to learn the global optimal policy and solve our proposed optimization problem, relying only on limited information about the network dynamics and CSI.
- We demonstrate quantitatively that our proposed algorithm outperforms existing heuristic-based methods in terms of satisfying the objective of the optimization problem. Also, we show how the DeepRAT algorithm can quickly adapt to the abrupt changes in network dynamics,

such as EDs' mobility.

The rest of this paper is organized as follows. Table I defines the main acronyms used in this paper. Section II describes some related work that implements DRL methods for radio resource allocation. Section III presents the proposed system model and formulates the optimization problem. Section IV shows the architecture of the proposed DeepRAT framework and explains its underlying DRL-based models along with a brief mathematical background. Section V explains the simulation setup and discusses the corresponding numerical results. Finally, Section VI concludes the paper.

## II. RELATED WORK

DRL-based techniques have attracted considerable research lately in the context of radio resource allocation for wireless networks [6], [8]. The authors in [9] proposed DRL models based on the single and multi-agent actor-critic algorithms to address the problem of total sum-rate maximization via power allocation for cellular networks. In [10], the authors used DRL methods to study the joint optimization of user association and power allocation in orthogonal frequency division multiple access (OFDMA)-based HetNets. Ye *et al.* [11] presented a DRL-based mechanism to study the problem of resource allocation for unicast and broadcast scenarios in vehicle-to-vehicle (V2V) networks. The work in [12] investigated the power control problem of device-to-device (D2D)-enabled networks in time-varying environments using a centralized DRL algorithm. Zhang *et al.* [13] presented a DRL algorithm to study the problem of energy-efficient resource allocation in ultra-dense cellular networks. The authors in [14] presented a multi-agent DQN-based model to study the problem of joint power, bandwidth, and throughput allocation in unmanned aerial vehicle (UAV)-assisted IoT networks. In [3], the authors proposed a non-cooperative multi-agent DQN-based method to study the problem of power allocation in hybrid RF/VLC networks. The authors show via simulation that the convergence rate of the DQN-based model is 96.1% compared to that of the $Q$-learning-based algorithm, which is 72.3%.

On the other hand, using deep deterministic policy gradient (DDPG) models has also gained increasing interest recently. They have shown superior performance in addressing the radio resource allocation problems in continuous and high dimensionality environments compared to the vanilla DQN algorithms [15]. The authors in [16] presented a comparative study for the applications of three DRL algorithms, namely the DDPG, Neural Episodic Control (NEC), and Variance Based Control, in the optimization of wireless networks. The authors concluded that the DDPG and VBC methods achieve better performance than the NEC-based algorithm. In [7], the authors presented a single-agent DDPG algorithm to address the problem of network selection in heterogeneous health systems. Their goal was to optimize the medical data delivery from Patient Edge Nodes (PENs) via multi-radio access networks (RANs) to the core network. Nasir *et al.* [17] presented a multi-agent DDPG-based algorithm to study the problem of joint power and spectrum allocation in wireless networks.

Table I: Definitions of Main Acronyms Used in this Paper.

| Acronym | Definition | Acronym | Definition | Acronym | Definition |
|---|---|---|---|---|---|
| RAT | Radio Access Technology | HetNet | Heterogeneous Network | ED | Edge Device |
| MINLP | Mixed-Integer Non-Linear Programming | DRL | Deep Reinforcement Learning | DQN | Deep Q Network |
| DDPG | Deep Deterministic Policy Gradient | QoS | Quality of Service | AI | Artificial Intelligence |
| OFDMA | Orthogonal Frequency Division Multiple Access | CSI | Channel State Information | V2V | Vehicle-to-Vehicle |
| NEC | Neural Episodic Control | VLC | Visible Light Communication | D2D | Device-to-Device |
| DSRC | Dedicated Short-Range Communication | PEN | Patient Edge Node | RF | Radio Frequency |
| SDN | Software-Defined Networking | DNN | Deep Neural Network | ES | Edge Server |
| AWGN | Additive White Gaussian Noise | SNR | Signal to Noise Ratio | AP | Access Point |
| C-RAN | Cloud/Centralized Radio Access Network | LTE | Long-Term Evolution | BBU | Baseband Unit |
| CDF | Cumulative Distribution Function | UAV | Unmanned Aerial Vehicle | OU | Ornstein-Uhlenbeck |
| OFDM | Orthogonal Frequency Division Multiplexing | IoT | Internet of Things | NR | New Radio |

Based on simulation results, the authors demonstrated how their proposed technique outperforms the conventional fractional programming algorithm. In [18], the authors investigated the problem of rate resource allocation for 5G network slices. The authors decomposed the problem into a master-slave, and proposed a multi-agent DDPG-based algorithm to solve it. Experimental results showed that their proposed algorithm performs better than some baseline approaches and provides a near-optimal solution.

In this paper, we present a multi-agent algorithm based on DRL called the DeepRAT to study the problem of cost-effective sum-rate maximization of HetNets via dynamic multi-RAT assignment and continuous power allocation for multi-connectivity multi-homing EDs. Towards this end, we formulate this problem as a mixed-integer non-linear programming (MINLP) problem and, due to the high complexity and combinatorial nature of the problem, we propose the DeepRAT algorithm to solve it efficiently and learn system dynamics, relying only on limited information about the network and CSI.

## III. System Model and Problem Formulation

This section describes our proposed system model and formulates the optimization problem.

### A. System Model

We consider a next-generation HetNet as depicted in Fig. 1. It consists of various RAT access points (APs), such as sub 6GHz, dedicated short-range communication (DSRC) for vehicular networks, 5G NR, 4G long-term evolution (4G LTE), and WiFi. It is assumed that the RATs have different operating characteristics, such as carrier frequency, spectrum, data rate, energy consumption, the monetary cost for using RAT services, and transmission delay. In order to guarantee judicious and efficient management of networks radio resources, we assume that the RATs are controlled by a cloud-based edge server (ES). The RATs are assumed to serve multi-connectivity (i.e., multi-access), multi-homing EDs. Note that, unlike our previous work in [1] in which each ED can be assigned one RAT at any time in a greedy fashion, i.e., multi-mode, the EDs in this paper are assumed to have the ability to connect to multiple RATs at any time to aggregate RATs' radio resources.

### B. Optimization Problem Formulation

Our objective is to study the problem of downlink sum-rate maximization of the multi-RAT HetNet while;

Table II: Definitions of Key Symbols Used in this Paper.

| Symbol | Definition |
|---|---|
| $L$ | The total number of RATs. |
| $U$ | The total number of EDs. |
| $L_u$ | The subset of RATs assigned to the $u$th ED. |
| $U_l$ | The subset of EDs assigned to the $l$th RAT. |
| $x_{lu}$ | The assignment indicator for $u$th ED over $l$th RAT. |
| $\mathbb{U}_{lu}$ | The utility function of $u$th ED over $l$th RAT |
| $R_{lu}$ | The data rate for $u$th ED over $l$th RAT. |
| $R_u$ | The achieved rate by $u$th ED from its assigned RATs. |
| $R_u^{min}$ | The minimum required rate for $u$th ED. |
| $C_{lu}$ | The monetary cost per second for $u$th ED over $l$th RAT. |
| $g_{lu}$ | The channel gain for $u$th ED over $l$th RAT. |
| $\Gamma_{lu}$ | The SNR for $u$th ED over $l$th RAT. |
| $p_{lu}$ | The allocated power for $u$th ED over $l$th RAT. |
| $P_l^{max}/W_l$ | The total power/bandwidth of $l$th RAT. |
| $h_{lu}$ | The small scale fading for $u$th ED over $l$th RAT. |
| $\varepsilon_l$ | The monetary cost per bit for $l$th RAT. |
| $\alpha_u$ & $\gamma_u$ | The weighting coefficients representing relative importance of the two metrics of utility function. |
| $\zeta_{ES}, \eta_{ES}, \eta_1, \eta_2, \zeta_l$ | The weighting factors used in rewards to show if satisfying the constraints has priority over maximizing objective. |

1) reducing the monetary cost of connection,
2) assigning each ED to the optimal set of RAT(s),
3) allocating the optimal downlink power for each RAT-ED communication link,
4) meeting the QoS requirements of EDs in terms of the minimum required data rate.

We denote by $\mathcal{L} \triangleq \{1, 2, \cdots, l, \cdots, L\}$ the set of all RATs indices, where $L$ represents the total number of RAT APs. We also denote by $\mathfrak{U} \triangleq \{1, 2, \cdots, u, \cdots, U\}$, the set of all EDs indices, where $U$ represents the total number of multi-homing EDs. The subset of EDs assigned to the $l$th RAT is denoted by $U_l$ such that $U_l \subseteq \mathfrak{U} \ \forall l \in \mathcal{L}$. In addition, the subset of RATs assigned to the $u$th ED is denoted by $L_u$ such that $L_u \subseteq \mathcal{L} \ \forall u \in \mathfrak{U}$.

Assuming an OFDM-based system with flat fading for each ED and the total bandwidth of each RAT is equally divided between the RATs' assigned EDs, i.e., $W_l/U_l$, the upper bound of downlink data rate for the $u$th ED over the $l$th RAT at time slot $t$ is expressed as [1], [3], [19]:

$$R_{lu}(t) = \frac{W_l}{U_l} \log_2 \left( 1 + \frac{g_{lu}(t) p_{lu}(t)}{\frac{W_l}{U_l} \sigma_l^2} \right), \quad (1)$$

where $W_l$ is the total bandwidth of the $l$th RAT and $\sigma_l^2$ is the $l$th RAT power spectral density of the additive white Gaussian noise (AWGN). The parameter $g_{lu}(t)$ in (1) represents the
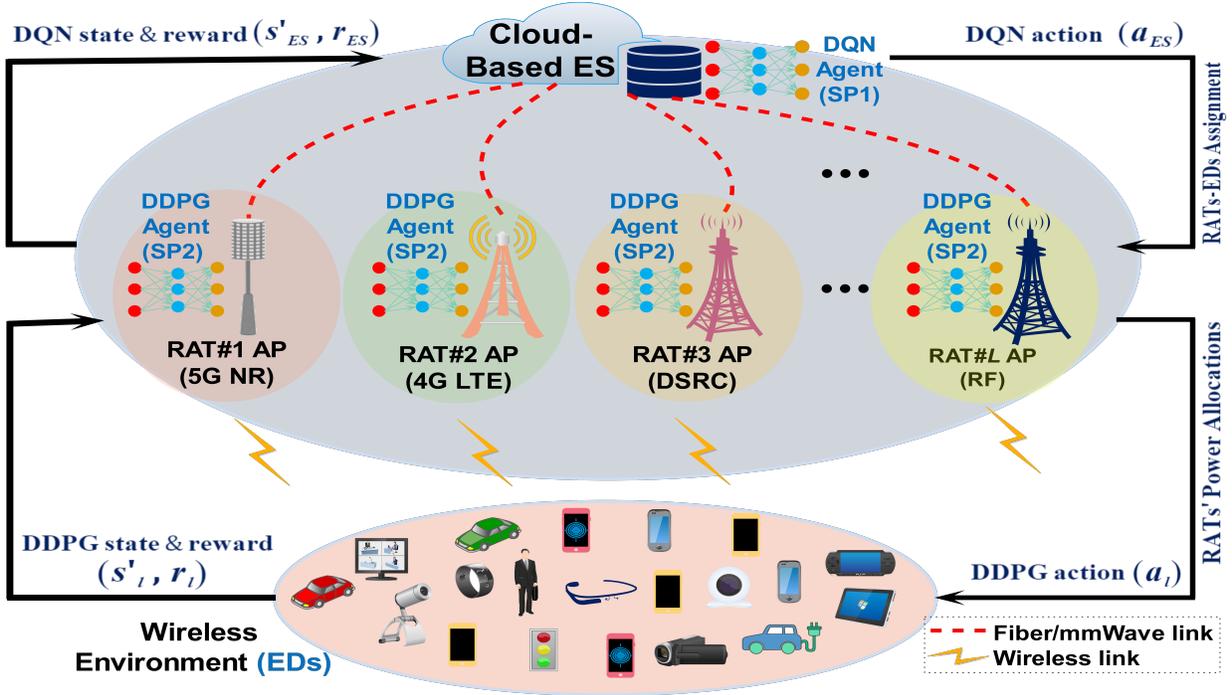
Fig. 1: System model of our proposed DeepRAT framework for next-generation HetNets. It comprises $L$ radio access technology (RAT) access points (APs) controlled by a cloud-based edge server (ES) and serves $U$ multi-homing edge devices (EDs).

channel gain for the $u$th ED over the $l$th RAT at time slot $t$, which is defined based on the channel model of the RAT used. In Section V, we will simulate three types of channel models, namely, the mmWave with beamforming for 5G NR [20], the COST 231 for 4G-LTE [20], and the exponential for 3G [10], [19]. Also, we will simulate a dynamic wireless system, as we will discuss later in Section V.

Now we define $\varepsilon_l$ as the monetary cost of the $l$th RAT, which is expressed in Euro per bit [7] and can be obtained from e.g., the IEEE 802.21 standard [21]. In our system model, such information can be easily gathered by the ES and stored in advance and only updated if there are changes in RATs' pricing. The monetary cost, expressed in Euro per second, resulting from using the $l$th RAT by the $u$th ED to receive the data rate $R_{lu}(t)$ at time slot $t$ is expressed as $C_{lu}(t) = \varepsilon_l R_{lu}(t)$. Table II summarizes the definitions of the symbols used in this paper.

In this paper, we are optimizing the downlink sum-rate of the HetNet in a cost-effective way, such that the QoS requirements of EDs are guaranteed. This is achieved via jointly : 1) assigning the optimal set of RAT(s) to each ED, i.e., $x_{lu} \forall u \in \mathfrak{U}, l \in \mathcal{L}$, and 2) allocating the optimal downlink power to each active RAT-ED link, i.e., $p_{lu}(t), \forall l \in \mathcal{L}, u \in U_l$. Consequently, our optimization problem is formulated

as follows:

$$\textbf{P1:} \quad \max_{x_{lu}, p_{lu}} \quad \sum_{l=1}^{L} \sum_{u=1}^{U} x_{lu} \mathbb{U}_{lu}(t)$$

$$\text{subject to} \quad C_1: \sum_{l=1}^{L} x_{lu} \geq 1, \; x_{lu} \in \{0,1\} \quad \forall u \in \mathfrak{U},$$

$$C_2: \sum_{u=1}^{U} x_{lu} p_{lu}(t) \leq P_l^{max} \quad \forall l \in \mathcal{L},$$

$$C_3: \sum_{l=1}^{L} x_{lu} R_{lu}(t) \geq R_u^{min} \quad \forall u \in \mathfrak{U},$$

$$C_4: p_{lu}(t) \geq 0, \quad \forall l \in \mathcal{L}, u \in \mathfrak{U}.$$

(2)

where $x_{lu}$ is the RTAs-EDs assignment indicator, such that $x_{lu} = 1$ if the ES assigns the $l$th RAT to the $u$th ED, and $x_{lu} = 0$ otherwise. To achieve our goal, we combine the data rate and monetary cost as a weighted sum objective utility function. In the weighted sum method, the Pareto optimal values can be achieved by adjusting the weighting parameters [22], and thus there is no optimality loss in the problem formulation. Therefore, we define $\mathbb{U}_{lu}(t)$ in **P1** as the utility function of the $u$th ED over the $l$th RAT at time slot $t$, which is given by:

$$\mathbb{U}_{lu}(t) = \alpha_u R_{lu}(t) - \gamma_u C_{lu}(t) \quad (3)$$

where $\alpha_u$ and $\gamma_u$ are weighting parameters representing the relative importance of the objectives of jointly maximizing data rate and reducing the cost at each ED, such that $\alpha_u + \gamma_u = 1$. In other words, $\alpha_u$ and $\gamma_u$ are EDs'-defined

parameters used to show if the ED cares more about getting a higher data rate over the monetary cost or getting a lower monetary cost over the data rate. Our objective is to find the optimum data rates $R_{lu}(t)$ that maximize our utility function via jointly controlling the RATs-EDs assignment and links' power allocations while considering the EDs preferences in terms of $\alpha_u$, $\gamma_u$, and $R_u^{min}$ and the RATs' rate prices $\epsilon_l$. Note that all quantities in (3) are normalized to their maximum values in order to make them comparable.

The optimization problem (2) is over the two unknowns: $x_{lu}$ and $p_{lu}$ subject to four constraints. $C_1$ ensures that each ED can be connected to multiple RATs simultaneously, and it reflects the multi-homing capabilities of EDs. $C_2$ and $C_4$ ensure that the power allocations from RATs to their assigned EDs do not violate the RATs' available power resources. $C_3$ ensures that the achievable data rates for EDs from their assigned RATs are greater than the minimum QoS requirements.

### C. Why DRL?

Problem (2) is a combinatorial Mixed-Integer Non-Linear Programming (MINLP) [23], which is highly complex to solve using traditional approaches. In particular, applying the exhaustive search algorithm to find $x_{lu}$ followed by optimization approaches to find the corresponding $p_{lu}$ is not practical as the search space will grow exponentially. For example, as we will show later in Section V, we simulate a scenario with $L = 3$ and $U = 10$. This means that applying the exhaustive method requires a full search over $(2^L - 1)^U = 282,475,249$ possible combinations, each of which is followed by a constrained optimization process to find $p_{lu}$. This is quite difficult and impractical. In addition, transforming the problem into a geometric program is not possible due to constraint $C_1$ with the nonlinearity of the objective [23], [24]. Note that compared to the problem formulation in [25], in **P1** we add the monetary cost aspect to the utility function, the multi-homing capabilities of EDs, and the power allocation issue. These aspects added new dimensions to the formulated problem and increased its complexity. In addition, compared to the problem formulation in [26], we added the following four additional dimensions in **P1**: the monetary cost issue to the objective function, the multi-homing constraint, the RATs' power allocation constraint, and the QoS requirements constraint of EDs. These new dimensions enriched the problem while making it more computationally expensive and difficult to solve using conventional approaches (e.g., optimization, ranking-based, and game theory methods discussed in Section I). The same observations are made when comparing **P1** with the problems formulated in [27] and [28]. Furthermore, unlike the previous works in [29], [30], we added the monetary cost issue and the multi-homing constraint in **P1**. Note that compared to our problem formulation in [1], we add the monetary cost and multi-homing dimensions. In addition, this paper carefully addresses the scalability issue of the proposed DeepRAT algorithm for the increasing number of EDs. Finally, we demonstrate how our proposed DRL-based

algorithm can quickly adapt, in terms of convergence speed, to the abrupt changes of the network, such as EDs' mobility.

Hence, and due to the high complexity of our formulated optimization problem, we propose to solve it using emerging DRL techniques instead. In particular, we hierarchically decompose **P1** into two optimization sub-problems, such that each sub-problem is a function of only one decision variable and, hence, can be solved separately and independently of the other sub-problem. The first sub-problem is to find the optimal RAT-EDs assignment $x_{lu}$, which depends on the parameters of the EDs. It is considered a global variable relevant to the overall HetNet system, and it can be solved at the ES level. The second sub-problem is to find the optimal power allocation for each RAT-ED $p_{lu}(t)$. It is considered a local variable that depends only on the parameters of RATs and can be solved at the RATs level. Consequently, **P1** is decomposed into the following two optimization sub-problems:

$$\textbf{SP1:} \quad \max_{x_{lu}} \quad \sum_{l=1}^{L}\sum_{u=1}^{U} x_{lu}\mathbb{U}_{lu}(t) \qquad (4)$$

subject to $\quad C_1 \quad$ and $\quad C_3$.

and

$$\textbf{SP2:} \quad \max_{p_{lu}} \quad \sum_{l=1}^{L}\sum_{u=1}^{U} x_{lu}\mathbb{U}_{lu}(t) \qquad (5)$$

subject to $\quad C_2 - C_4$.

Next, we proceed with our proposed methodology to solve sub-problems **SP1** and **SP2** using DRL.

## IV. DRL FOR DYNAMIC MULTI-RAT ASSIGNMENT AND POWER ALLOCATION

In this section, we first explain our proposed DeepRAT algorithm to solve sub-problems **SP1** and **SP2**. Then, we provide a detailed description of the elements of the DQN and DDPG models used in the proposed DeepRAT model.

### A. The DeepRAT Framework

We propose a multi-agent DRL-based framework called the DeepRAT, which hierarchically solves **SP1** and **SP2** iteratively and interactively in two stages; *RATs-EDs assignment* and continuous *power allocation*. DeepRAT employs two types of DRL algorithms, a single-agent DQN at the ES and multi-agent DDPG at the RATs level, as depicted in Fig. 1. The methodology of the DeepRAT algorithm to solve the two sub-problems is explained as follows. For sub-problem **SP1**, Deep-RAT utilizes a single-agent DQN algorithm to optimize the RATs-EDs assignment $x_{lu}$ at the ES level, while considering $p_{lu}$ as constants, which are passed by the RATs. Note that this RATs-EDs assignment is initially performed randomly by the ES without prior knowledge of whether it would be optimal or not. Then, the ES broadcasts $x_{lu}$ to the multi-agent DDPG algorithms of each RAT in order to optimize their power allocation $p_{lu}$ according to sub-problem **SP2** while considering $x_{lu}$ as constants. The ES then receives feedback ACK signals from all RATs indicating whether the objective of **P1** has

been successfully solved for the current RATs-EDs assignment $x_{lu}(t)$ (i.e., the RATs are in good status) or not (i.e., the RATs are in bad status). Based on these ACK signals, the ES starts learning to make better assignments in the future time slots. These two stages are iteratively executed until all DeepRAT's agents learn the global policy that solves our main problem in **P1**, i.e., the single-agent DQN learns the optimal RATs-EDs assignment policy, and the multi-agents DDPG learn the optimal power allocation policy. The main elements of these two types of DRL models are defined next.

### B. DeepRAT Stage 1: DQN Algorithm for RATs-EDs Assignment

Due to the discrete nature of the RATs-EDs assignment problem, we adopt the DQN algorithm to act as an ES agent to learn the optimal policy for this problem. Below, we define the state space, action space, and reward function for the single-agent ES DQN model.

*1) DQN action space:*
At each time step $t$, the main role of the DQN ES agent is to take an action $a_t^{ES}$ that optimally assigns each of the EDs to the optimal set of RAT(s), i.e., obtaining $x_{lu}$ in **SP1**. As shown in the optimization problem **SP1**, the DQN assignment action should: 1) maximize the objective function, and 2) satisfies $C_1$ and $C_3$. These conditions can be achieved iteratively by the design of the reward function. This action is a combinatorial problem that scales exponentially with the number of EDs and RATs, causing degradation in both system scalability and convergence speed. Unlike our previous work [1], in which the size of action space was proportional to both the number of EDs and RATs, i.e., $L^U$, in this paper we carefully address this scalability issue. Specifically, we make the size of the action space proportional to only the number of RATs, i.e., $2^L$. The scalability enhancement achieved is evident. As an example, for $L = 3$ and $U = 5$, the size of the action space proposed in this paper is around 30 times less than the one proposed in [1]. Similarly, when $L = 4$ and $U = 10$, the size reduction of the action space is around 65536. This will greatly enhance system scalability and reduce convergence time. Therefore, the action space of the DQN ES agent is discrete, corresponding to assigning the optimal set of RAT(s) $L_u$ to the $u$th ED $\forall u \in \mathfrak{U}$, where $L_u \in \mathcal{L}$, which is expressed as:

$$\mathcal{A}_{ES}(t) = [a_1(t), a_2(t), \cdots, a_u(t), \cdots, a_U(t)], \quad (6)$$
$$a_u(t) = \{L_1, L_2, \cdots, L_{2^L}\},$$

*2) DQN state space:*
Due to the holistic view of the DQN ES agent, its state space must include effective and rich information about RATs and EDs to help the DQN agent in taking optimal assignment actions. To address the scalability issue mentioned previously, the RATs-EDs assignment is done by the ES iteratively for each ED, i.e., the ES assigns the $u$th ED to the best set of RATs $L_u$, while considering the assignment of the remaining EDs constant. Therefore, the ES DQN is configured to run on episodes of $U$ EDs time steps, and the DQN state must indicate the ED investigated [7]. In addition, unlike our previous work

in [1] in which we assumed that the CSI is available at the agents, in this work we consider a more practical and challenging scenario by assuming that the agents have limited information about network dynamics and CSI. In particular, we assume that only historical information about the achieved data rates is available to the agents. With this in mind, the state space of ES is comprised of two main components, global information related to all RATs and EDs in the network and local information related to the $u$th ED investigated. The global information contains three elements; the matrix of all $x_{lu}$ at the previous time step (i.e., $x_{lu}(t-1) \in \mathbb{C}^{U \times L} \ \forall l, u$), the matrix of all $R_{lu}$ at the previous time step (i.e., $R_{lu}(t-1) \in \mathbb{C}^{U \times L} \ \forall l, u$), and the vector of all $R_u^{\min} \in \mathbb{C}^{U \times 1} \ \forall u$. The local information is only related to the $u$th ED under investigation, which has three *scalar* elements; the index of ED under investigation at the current time step $uth(t) \in \mathbb{C}^{1 \times 1}$, the minimum required date of this ED $R_{uth}^{\min} \in \mathbb{C}^{1 \times 1}$, and the achieved data rate of this ED from *its assigned RATs* at the previous time step $R_{uth}(t-1) \in \mathbb{C}^{1 \times 1}$. Consequently, the state space of the DQN ES agent is expressed as:

$$\mathcal{S}_{ES}(t) = [\text{Global information}, \ \text{Local information}], \quad (7)$$
$$\mathcal{S}_{ES}(t) = [\underbrace{x_{lu}(t-1), R_{lu}(t-1), R_u^{min}}_{\text{Global information}},$$
$$\underbrace{uth(t), R_{uth}^{\min}, R_{uth}(t-1)}_{\text{Local information}}].$$

*3) DQN reward function:*
The reward function is designed to incorporate the objective of our optimization problem in **P1** on one hand and the constraints $C_1$ and $C_3$ on the other hand. Hence, the agent will receive a negative punishment if the constraints are violated. The instantaneous reward $r_{ES}$ that the ES agent receives when taking action $a_{ES}$ given state $s_{ES}$ for the $u$th ED is expressed as:

$$r_{ES}(t) \triangleq \eta_{ES} constraints + \zeta_{ES} objective, \quad (8)$$
$$r_{ES}(t) = \eta_{ES} \sum_{u}^{U} \left( \sum_{l \in L_u} R_{lu}(t) - R_u^{min} \right) +$$
$$\zeta_{ES} \sum_{l=1}^{L} \sum_{u=1}^{U} x_{lu} \mathbb{U}_{lu}(t).$$

where $\eta_{ES}$ and $\zeta_{ES}$ are weighting factors used to indicate whether satisfying the constraints has priority over maximizing the objective or not. These factors are manually tuned during simulation.

The RATs-EDs assignment problem can be formulated based on the immediate rewards achieved. Towards this goal, the expected accumulated discounted instantaneous reward over the time horizon $T$ is defined as $\mathcal{R}_{ES} = \mathbb{E}\left[\sum_{t=1}^{T} \gamma^{t-1} r_{ES}(t)\right]$, where $0 \leq \gamma \leq 1$ is a discounted factor [15]. The objective of the DQN ES agent is to obtain the optimal decision policy $\pi_{ES}^*$ (i.e., selecting the optimal RATs-EDs assignment $x_{lu}$) that maximizes $\mathcal{R}_{ES}$. This is expressed as $\pi_{ES}^* = \underset{\pi}{\operatorname{argmax}} \ \mathcal{R}_{ES}$.

However, as we discussed earlier in Subsection III-C, this RATs-EDs assignment problem is space-hard and it is quite difficult for traditional resource allocation techniques to solve it [6], [8]. Therefore, the DQN algorithm can be leveraged instead to learn $\pi^*_{ES}$. In DQN, the optimal policy is expressed as $\pi^*_{ES} = \arg\max_a Q^{\pi^*_{ES}}_{ES}(s^{ES}_t, a^{ES}_t)$, where the function $Q^{\pi^*_{ES}}_{ES}(s^{ES}_t, a^{ES}_t)$ is called the state-action value function. This value function defines the expected accumulated discounted instantaneous reward achieved when executing action $a^{ES}_t$ in state $s^{ES}_t$ and then following the policy $\pi_{ES}$ thereafter. The value-function is defined as $Q^{\pi^*_{ES}}_{ES}(s^{ES}_t, a^{ES}_t) = \mathbb{E}_{\tau \sim \pi_{ES}}[\mathcal{R}_{ES}|s^{ES}_t, a^{ES}_t]$, and the DQN algorithm utilizes the following iterative Bellman equation to compute it:

$$Q^*_{ES}(s^{ES}_t, a^{ES}_t) = r_{ES}(s^{ES}_t, a^{ES}_t) + \gamma \max_{a^{ES}_t} Q^*_{ES}(s^{ES}_{t+1}, a^{ES}_{t+1}),$$
(9)

At each decision time step $t$, the deep neural network (DNN) of the ES DQN model iteratively updates its weights $\theta_{ES}$ to minimize the following loss function:

$$L_{ES}(\theta_t) = \mathbb{E}_{s,a,r,s' \in \mathcal{D}_{ES}}[(y_{ES}(t) - Q_{ES}(s^{ES}_t, a^{ES}_t|\theta_{ES}))^2],$$
(10)

where $y_{ES}(t) = r_{ES}(s^{ES}_t, a^{ES}_t) + \gamma \max_{a^{ES}_{t+1}} Q_{ES}\left(s^{ES}_{t+1}, a^{ES}_{t+1}|\theta'_{ES}\right)$ is the target value, which is obtained from the target network with old weights $\theta'_{ES}$, and $\mathcal{D}_{ES}$ represents the DQN replay buffer.

## C. DeepRAT Stage 2: DDPG Algorithm for Power Allocation

The DDPG is an efficient DRL algorithm developed to learn policies for continuous-based problems with high dimensionality in state and action spaces [16], [18], [31]. The DDPG algorithm will be leveraged in our second stage, i.e., solving the power optimization problem. In particular, a multi-agent deployment is considered in which each RAT employs a DDPG agent, whose main goal is solving its own objective function in **SP2** for its assigned EDs, $U_l$. The main elements of the multi-agent DDPG algorithm used are defined below.

*1) DDPG action space:*

Each DDPG RAT agent takes action independently and uncooperatively from the other agents. At time slot $t$, once the DQN ES agent executes the assignment action for the $u$th ED ($L_u$), the main goal of each DDPG agent is to optimize the power allocation $p_{lu}$ in **SP2** $\forall u \in \mathcal{U}$ and $\forall l \in \mathcal{L}$. As shown in the optimization problem **SP1**, the power allocation action of the $l$th DDPG RAT agent should: 1) maximize the objective function, and 2) satisfy $C_2$, $C_3$, and $C_4$. These conditions can be achieved by incorporating them into the reward function. Consequently, the action space of the $l$th DDPG RAT agent is continuous with a size of $U$, corresponding to deciding the optimal power allocation for each of the RATs-EDs communication links (i.e., $p_{lu}(t)$). The action space of the $l$th DDPG RAT agent is defined as:

$$\mathcal{A}_l(t) = [p_{l1}(t), p_{l2}(t), \cdots, p_{lu}(t), \cdots, p_{lU}(t)]$$
$$\text{where} \quad p_{lu}(t) \in [0, P^{max}_l] \ Watt,$$
(11)

The action exploration-exploitation problem in the DDPG algorithm is addressed via adding some Gaussian or Ornstein-Uhlenbeck (OU) noise $n_t$ to the selected action $a^l_t$ [31].

*2) DDPG state space:*

We assume that the DDPG RAT agents do not cooperate, and there is no direct communication between them. The agents, however, have direct communication with the DQN ES, which has a holistic view of all DDPG agents and can coordinate them. This means that each DDPG agent can acquire information about the other agents via the ES. The state space of the $l$th DDPG RAT agent is designed to contain useful information on the underlying HetNet. Four main types of representative information are incorporated in each $l$th agent state space. The first type of information is discrete and is directly related to the current RATs-EDs assignment action of the ES agent, which is a vector of the set of EDs assigned to the $l$th RAT at the current time step ($U_l(t) \in \mathbb{C}^{U \times 1}$). The second type of information is continuous, which is the vector of $R^{min}_u \in \mathbb{C}^{U \times 1}$, $\forall u \in U_l$. The third information is the vector of $R_{lu}$ for the $l$th RAT at the previous time slot (i.e., $R_{lu}(t - 1) \in \mathbb{C}^{U \times 1}$, $\forall u \in U_l$). The fourth information is the vector of $R_u$ at the previous time slot (i.e., $R_u(t - 1) \in \mathbb{C}^{U \times 1}$, $\forall u \in U_l$). Note that while $R_{lu}(t - 1)$ denotes the downlink rate for a *single* link between the $u$th ED and $l$th RAT at the previous time slot, the notation $R_u(t-1)$ denotes the rate achieved by the $u$th ED from *its assigned RATs* $L_u$ at the previous time slot, i.e., $R_u(t-1) = \sum_{l \in L_u} R_{lu}(t-1)$. Consequently, the state space of the $l$th DDPG RAT agent is represented as:

$$\mathcal{S}_l(t) = [U_l(t), R^{min}_u, R_{lu}(t-1), R_u(t-1)].$$
(12)

*3) DDPG reward function:*

The reward of the $l$th DDPG RAT agent is expressed as a continuous function that is governed by the RAT's achieved constrained objective function. It is quantified by including the optimization constraints $C_2$, $C_3$, and $C_4$ of (2) into the reward function so that the instantaneous reward reflects whether the constraints are satisfied or not [7], [16]. The reward function is given by (13). In (13), $\eta_1, \eta_2$, and $\zeta_l$ are also weighting factors used to indicate whether satisfying the constraints has priority over maximizing the objective or not. These factors are manually tuned during simulation.

In this second stage of our problem **SP2**, the main objective is to derive the optimal power allocation policy $\pi^*_l$ that maximizes the long-term reward of the $l$th agent $\mathcal{R}_l$, i.e., $\pi^*_l = \underset{\pi}{\text{argmax}} \ \mathcal{R}_l = \underset{\pi}{\text{argmax}} \sum_{t=0}^T \gamma^{t-1} r_l(s^l_t, a^l_t)$. Towards this goal, we implement the DDPG algorithm to derive this policy $\pi^*_l$. The DDPG algorithm integrates the DQN and actor-critic algorithms [16], [31], and will be utilized to perform the training of the RATs' DNNs. The DDPG has one parameterized actor function and one parameterized critic function represented by $\mu(s^l_t|\theta^\mu_l)$ and $Q_l(s^l_t, a^l_t|\theta^Q_l)$, respectively, where $\theta^\mu_l$ and $\theta^Q_l$ denote the weights of the actor and critic networks, respectively. The parameterized actor function is used to derive the policy, and it is implemented with a DNN trained

$$r_l(t) \triangleq constraints + \zeta_l objective,$$
$$r_l(t) = \eta_1 \left( P_l^{max} - \sum_{u \in U_l} p_{lu}(t) \right) + \eta_2 \sum_{u \in U_l} \left( \sum_{l \in L_u} R_{lu}(t) - R_u^{min} \right) + \zeta_l \mathbb{U}_{lu}(t). \qquad (13)$$

based on the iterative Bellman equation. On the other hand, the parameterized critic function is used to derive the value function $Q_l^\mu(s_t^l, a_t^l)$, defined as $Q_l^\mu(s_t^l, a_t^l) = \mathbb{E}_{\tau \sim \pi}[\mathcal{R}_l | s_t^l, a_t^l]$, and it is implemented using a DQN. The goal of each DDPG is to find the optimal policy $\pi_l^*$ that maximizes the long-term reward $\mathcal{R}_l$ using $\pi_l^* = \arg\max_a Q_l^{\mu^*}(s_t^l, a_t^l)$. The value function $Q_l^{\mu^*}(s_t^l, a_t^l)$ is derived iteratively using the Bellman equation similar to (9), and the policy $\pi_l^*$ is found via training the DNN of the $l$th DDPG RAT agent to minimize the Bellman loss function given by the following formula:

$$L_l(\theta_l^Q) = \mathbb{E}_{s,a,r,s' \in \mathcal{D}_l}[(y_l(t) - Q_l(s_t^l, a_t^l | \theta_l^Q))^2] \qquad (14)$$

where $\mathcal{D}_l$ denotes the $l$th DDPG agent's replay memory and $y_l(t)$ represents the target value, which is derived from the target network and obtained from the following equation:

$$y_l(t) = r_l(s_t^l, a_t^l) + \gamma \max_{a_{t+1}} Q_l \left( s_{t+1}^l, \mu_l(s_{t+1}^l | \theta_l^{\pi_l'}) | \theta_l^{Q'} \right) \quad (15)$$

where $\theta_l^{Q'}$ denotes the weights of target critic network, which has the same architecture as the main $Q$-network. These weights are mainly used to make the training more stable, and they are periodically updated based on the weights of the main $Q$-network $\theta_l^Q$.

The actor network of the $l$th DDPG agent is trained via applying the chain rule to the expected return from the cumulative reward distribution $J$ with respect to $\theta_l^\mu$ [31]:

$$\nabla_{\theta_l^\mu} J = \mathbb{E} \left[ \nabla_{\theta_l^\mu} Q_l(s_t^l, a_t^l | \theta_l^Q) |_{s=s_t, a=\mu(s_t, \theta_t^\mu)} \right] \qquad (16)$$

The detailed pseudo-code of our proposed multi-homing DeepRAT algorithm for solving (2) is given in Algorithm 1 and explained next. Lines 1 to 3 initialize the network parameters, ES DQN model, RATs DDPG models, and initial states. The episode begins with initial states for all the agents and iterates over all EDs. In Lines 4 to 12, the DQN ES agent observes the state for each ED and takes the corresponding assignment action to the RATs, $L_u$. In Lines 13 to 23, each DDPG RAT agent observes the state space (including the current assignment action of the DQN agent) and takes the corresponding power allocation action to the EDs $p_{lu} \forall u \in U_l$. In Lines 24 to 27, the DQN ES agent receives the reward for each ED assignment and learns to take a better assignment in future episodes. This process is repeated until the DeepRAT converges to the optimal policy that solves our main problem **P1** in (2), i.e., the DQN ES agent learns the optimal RATs-EDs assignment policy and all the DDPG RAT agents learn the optimal power allocation policy.

---

**Algorithm 1** The Multi-Homing DeepRAT Algorithm

**Input:** $L, U, R_u^{min}, \alpha_u, \gamma_u$.
**Output:** Optimal RATs-EDs assignment $x_{lu}$ & power allocation $p_{lu}$.

1: **Initialization**: Set $t = 0$ and initialize $\mathcal{D}_{ES}$ of ES DQN agent and $\mathcal{D}_l$ of RATs' DDPG agents, $\forall l \in \mathcal{L}$.
2: Randomly initialize weights of ES's DQN ($\theta_{ES}$ & $\theta'_{ES}$), and RATs' DDPG ($\theta_l^\mu, \theta_l^Q, \theta_l^{\mu'}, \& \theta_l^{Q'}$).
3: Initialize states $s_0$ of ES and RATs with initial observations.
4: **for** episode = 1 to $M$ **do**
5:    **for** ED = 1 to $U$ **do**
6:      Generate a random number $x$ from $[0, 1]$
7:      **if** $x \geq \epsilon(t)$ **then**
8:        Choose action $a_t^{ES}$ from ES's action space $\mathcal{A}_{ES}(t)$ according to $\max_{a \in \mathcal{A}_{ES}} Q_{ES}(s_t^{ES}, a_t^{ES} | \theta_{ES})$.
9:      **else**
10:        Choose a random action $a_t^{ES}$ from $\mathcal{A}_{ES}(t)$.
11:      **end if**
12:      Observe state for ES $s_t^{ES}$ and perform action $a_t^{ES}$.
13:      **for** $t_{RAT} = 0$ to $K$ **do**
14:        Observe state for each RAT $s_t^l$ using (12) and take action $a_t^l$ from $\mathcal{A}_l^p$ with OU noise, i.e., $a_t^l = \mu_l(s_t^l | \theta_l^{\mu_l}) + n_t^l$.
15:        Receive reward $r_l(s_t^l, a_t^l)$ using (13), observe $s_{t+1}^l$ using (12), and store transitions $(s_t^l, a_t^l, r_l, s_{t+1}^l)$ in $\mathcal{D}_l$.
16:        Sample random mini-batch of $N$ transitions $(s_i^l, a_i^l, r_i^l, s_{i+1}^l)$ from $\mathcal{D}_l$.
17:        Set $y_l(i)$ based on (15).
18:        Update $\theta_l^Q$ by minimizing loss (from (14)): $L_l(\theta_l^Q) = \frac{1}{N} \sum_{i=1}^N (y_l(i) - Q_l(s_i^l, a_i^l | \theta_l^Q))^2$.
19:        Update $\theta_l^\mu$ using sampled policy gradient (from (16)): $\nabla_{\theta_l^\mu} J = \frac{1}{N} \sum_{i=1}^N \nabla_{\theta_l^\mu} Q_l(s_i^l, a_i^l | \theta_l^Q) |_{s=s_i, a=\mu(s_i, \theta_l^\mu)}$
20:        Update weights of all RATs' DDPG target networks:
21:        $\theta_l^{Q'} \leftarrow \tau \theta_l^Q + (1 - \tau) \theta_l^{Q'}$,
22:        $\theta_l^\mu \leftarrow \tau \theta_l^\mu + (1 - \tau) \theta_l^\mu$
23:      **end for**
24:      Receive reward $r_{ES}(s_t^{ES}, a_t^{ES})$ using (9), observe $s_{t+1}^{ES}$ using (8), & store $(s_t^{ES}, a_t^{ES}, r_{ES}, s_{t+1}^{ES})$ in $\mathcal{D}_{ES}$,
25:      Sample random mini-batch of $M$ transitions $(s_i^{ES}, a_i^{ES}, r_{ES,i}, s_{i+1}^{ES})$ from $\mathcal{D}_{ES}$.
26:      Update weights of ES $\theta_{ES}$ to minimize loss in (10).
27:    **end for**
28: **end for**

---

### D. Deployment Scenario of the DeepRAT Framework

Our proposed multi-homing DeepRAT algorithm is simple yet practical for implementation using simple software-defined radios (SDRs). During the training phase of the DeepRAT model, the expensive computations are conducted offline on quasi-centralized hardware, such as GPUs and/or tensor processing units. Once the DeepRAT algorithm learns the optimal global policy, it can be deployed online to perform optimal decisions autonomously by relying only on its learned policies without inducing any extra delay. This aspect will be quantified in the next section. It is noteworthy that updating DeepRAT's

DNNs is only required if the characteristics of the wireless environment have changed significantly, and they are no longer reflecting the training experiences. Such a case occurs once per several weeks or even months.

In addition, the design principle of the DeepRAT framework is quite practical for the modern AI-driven wireless networks. Here, we provide three practical deployment scenarios of the proposed DeepRAT ecosystem. First, the DeepRAT can be deployed in the Cloud/Centralized Radio Access Networks (C-RANs) architecture, in which the cloud-based ES can be allocated at the baseband unit (BBU) pool of the C-RANs. Second, the DeepRAT framework can be deployed in the Software-Defined Networking (SDN) architecture, in which the ES can be placed at the control-plane side of the architecture. Third, the DeepRAT framework can be deployed in the future self-organizing/sustaining networks [2], in which the ES can be placed at the self-organizing/sustaining server.

## V. PERFORMANCE EVALUATION

This section presents a detailed description of the simulation setup we used to evaluate the performance of our proposed multi-homing DeepRAT algorithm. We first discuss the specifications of the HetNet under investigation, the DRL models used, and the EDs requirements. Then, we present and discuss the simulation results. Also, we present a practical scenario where there are abrupt changes in EDs mobility to demonstrate the ability of our proposed DeepRAT model to dynamically adapt to these varying system dynamics.

### A. Simulation Setup

We consider a practical scenario of a next-generation cellular network comprised of three multi-RAT OFDM-based systems, i.e., $L = 3$, specifically, 5G NR, 4G LTE, and 3G, as shown in Fig. 1. The specifications of these three systems are shown in Table III. Note that during our simulation, we considered three practical channel models for each RAT with flat fading, namely, the mmWave with beamforming for 5G NR [20], the COST 231 for 4G-LTE [20], and the exponential for 3G [10], [19]. The specifications of these channel models are listed in Table III. The RATs are 100 meters apart. Ten single-antenna EDs are assumed, i.e., $U = 10$, requesting services from the ES with random QoS requirements i.e., $R_u^{min}$, $\alpha_u$, and $\gamma_u$ as shown in Table IV. In order to model a dynamic wireless system, we assume a time-varying network where the mobility of EDs varies over time with random speeds ranging from 2 to 6 km/h. This means that the CSI, in terms of channel gain, of all links will dynamically change over time.

The number of agents is four; one DQN-based located at the ES side and three DDPG-based located at each RAT, as depicted in Fig. 1. These DRL models are simulated in Python using the Pytorch library, with architectures as shown in Table V. Relu activation functions are used at the output layers of all NNs, and the weights are updated using the Adam optimizer [32]. Also, in order to satisfy $C_4$ in **P1**, we employ the sigmoid function at the output layers of the DDPG actor networks.

Table III: SIMULATION PARAMETERS [10], [20].

| Parameter | RAT1 (5G) | RAT2 (4G LTE) | RAT3 (3G) |
|---|---|---|---|
| Frequency (GHz) | 28 | 6 | 2.4 |
| Bandwidth (MHz) | 200 | 40 | 27 |
| Max power (dBm) | 43 | 40 | 42 |
| Noise spectral density (dBm/MHz) | -57 | -57 | -57 |
| Channel model | Directional | COST 231 (Urban) | Exponential |
| Path loss exponent | 2(LOS), 4(NLOS) | - | 2 (LOS) |
| Number of uniform linear array antennas | 4 | 4 | 1 |
| Number of multipaths | 4 | 4 | - |
| Antenna gain (dBi) | 3 | 11 | - |
| Shadowing (dB) | 3.1 | 3 | 1.8 |
| $\epsilon_l$ (Euro/bit) | 9e-6 | 6e-6 | 1e-6 |

Table IV: QoS REQUIREMENTS OF EDs.

| ED ID | $R_u^{min}$ (bps) | $\alpha_u$ | $\gamma_u$ |
|---|---|---|---|
| ED1 | $8.3 \times 10^4$ | 0.4 | 0.6 |
| ED2 | $8.49 \times 10^4$ | 0.3 | 0.7 |
| ED3 | $1.17 \times 10^4$ | 0.2 | 0.8 |
| ED4 | $4.78 \times 10^4$ | 0.2 | 0.8 |
| ED5 | $1.37 \times 10^4$ | 0 | 1 |
| ED6 | $1.43 \times 10^4$ | 0.5 | 0.5 |
| ED7 | $6.1 \times 10^4$ | 0.4 | 0.6 |
| ED8 | $1.58 \times 10^4$ | 0.6 | 0.4 |
| ED9 | $8.93 \times 10^4$ | 0.6 | 0.4 |
| ED10 | $7.24 \times 10^4$ | 0.1 | 0.9 |

Table V: HYPERPARAMETERS OF THE DEEPRAT MODEL.

| Hyperparameter | Value | |
|---|---|---|
| | Single-Agent DQN | Multi-Agent DDPG |
| Buffer size | 1000 | 500 |
| Batch size | 64 | 16 |
| $\gamma$ | 0.99 | 0.99 |
| Number of layers | 2 | actor=2 critic=2 |
| Number of neurons | (256, 128) | (16, 16) |
| Learning rate | $8 \times 10^{-4}$ | actor=$5 \times 10^{-4}$ critic=$5 \times 10^{-4}$ |
| Exploration/ Exploitation noise | $\epsilon_{start} = 1$ $\epsilon_{end} = 0.005$ $\epsilon_{decay} = 5 \times 10^{-4}$ | OU ($\theta = 0.15$, $\sigma = 0.03$) |
| Rewards weighting factors | $\eta_{ES} = 1 \times 10^3$ & $\zeta_{ES} = 8 \times 10^{-4}$ | $\eta_1 = 1$, $\eta_2 = 1 \times 10^3$, & $\zeta_l = 5 \times 10^{-3}$ |

### B. Numerical Results

In this subsection, we present simulation results to evaluate the performance of our proposed DeepRAT algorithm when deployed in an online fashion. Also, in order to evaluate the performance of our proposed multi-homing DeepRAT algorithm, we compare it against four state-of-the-art benchmarks. 1) The multi-mode method (i.e., maximum or greedy method), in which the ES will greedily assign each ED to only *one* RAT that gives the maximum utility after the convergence of the DeepRAT algorithm, and the DDPG agents are utilized for power allocation, i.e., similar to our conference version in [1] and the work in [29]. It is noteworthy that the core difference between our implementation in this paper and the previous implementation in [1] is the limited information about the multi-RAT network in this paper. In particular, our implementation of the multi-mode algorithm in [1] assumes that the agents have full knowledge of network dynamics and CSI, i.e., channel gains $g_{lu}$, and we compared our approach against the CVXPY solver's solution [33]. However, our proposed multi-homing DeepRAT approach presented in this paper assumes that the agents have limited information about system dynamics and CSI). 2) The random approach, in which the ES will randomly assign each ED to only one RAT after the convergence of the

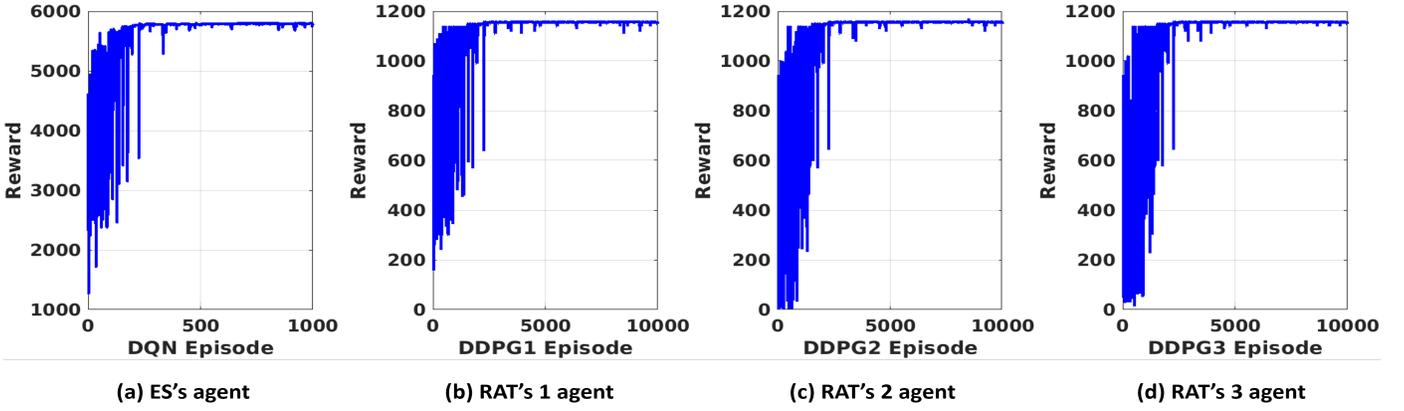(a) ES's agent      (b) RAT's 1 agent      (c) RAT's 2 agent      (d) RAT's 3 agent

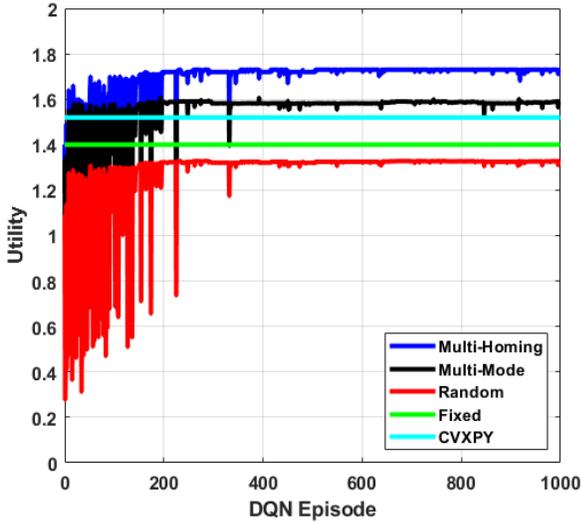Fig. 2: Reward training convergence of the DeepRAT's DQN and multi-agents DDPG models.



Fig. 3: Utility of the proposed multi-homing DeepRAT algorithm compared to the multi-mode, random, fixed, and CVXPY schemes. Note that, unlike the proposed multi-homing DeepRAT technique, these four heuristic-based methods are not guaranteed to meet the EDs' QoS requirements in terms of rates, $\alpha_u$, and $\gamma_u$ preferences.

DeepRAT algorithm, and the DDPG agents are used for power allocation. 3) The fixed approach, in which the ES will assign EDs to *all* existing RATs, and the RATs will allocate their power *equally* to all EDs. 4) The CVXPY solver's solution with MOSEK sub-solver [33], in which the ES will assign EDs to *all* existing RATs, and the CVXPY solver is utilized to solve the power allocation optimization problem assuming full knowledge of system dynamics and instantaneous CSI, similar to the works in [1], [18]. However, we should emphasize the following. 1) Although these conventional approaches show good results compared to our proposed multi-homing approach, they do not always guarantee an optimal solution, i.e., the QoS requirements of EDs are not guaranteed. 2) Unlike our proposed multi-homing DeepRAT approach, which works based on limited information about system dynamics and CSI, the fixed and CVXPY methods require perfect knowledge of

the multi-RATs and instantaneous CSI. 3) Additionally, our DeepRAT algorithm adapts to abrupt network changes such as EDs' mobility. However, these conventional approaches are not adaptable, which severely degrades performance, accuracy, and reliability of the learned policies [6], [8].

Fig. 2 shows the training rewards of all DeepRAT's DQN and DDPG agents. The DQN converges to the steady-state of optimal RATs-EDs assignment policy after 226 DQN episodes, while all three DDPG agents converge to the optimal power allocation policy after 2252 DDPG episodes. These results clearly show how our various value-based (i.e., DQN) and policy-based (i.e., DDPG) DRL agents efficiently interact with each other in order to learn a unified global optimal policy that solves our optimization problem in **P1**.

Fig. 3 shows the utility function, i.e., the objective of the optimization problem **P1**, which clearly demonstrates that the proposed multi-homing DeepRAT algorithm converges to the optimum solution after 226 episodes. Also, Fig. 3 shows that the utility of the proposed multi-homing DeepRAT algorithm outperforms the ones achieved by the state-of-the-art approaches mentioned above, i.e., the multi-mode, random, fixed, and CVXPY approaches. In particular, the steady-state utility function of the proposed DeepRAT algorithm is 1.73 compared to 1.58, 1.52, 1.33, and 1.4, respectively, for the multi-mode, CVXPY, random, and fixed methods. Recall that the multi-mode, random, and fixed methods do not guarantee that they satisfy the EDs' QoS requirements. In addition, note that the multi-mode outperforms the CVXPY as the latter does not assign EDs to the optimal set of RATs leading to a lower utility value. Fig. 4 also shows the corresponding cumulative distribution function (CDF) of the utility function for the proposed multi-homing DeepRAT algorithm and these four state-of-the-art approaches. It clearly shows that the median of the utility for multi-homing DeepRAT is 1.73 compared to 1.58, 1.52, 1.33, and 1.4, respectively, for the multi-mode, CVXPY, random, and fixed methods.

In Fig. 5, we present the total sum-rate achieved by our proposed multi-homing algorithm compared to the four conventional approaches. The steady-state of total sum-rate of our proposed DeepRAT approach is 4.2 Gbps compared to 4.33 Gbps, 3.86 Gbps, 3.23 Gbps, and 3.42 Gbps for the CVXPY,
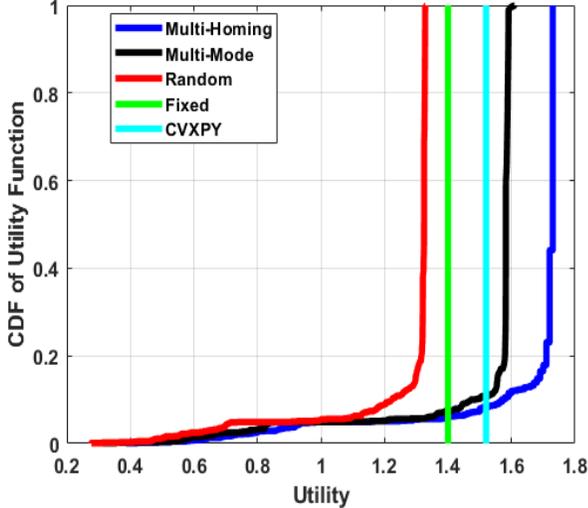
Fig. 4: CDF of the utility function of the proposed multi-homing DeepRAT algorithm compared with the other heuristic-based schemes.

multi-mode, random, and fixed methods, respectively. Hence, the percentage increase in total sum-rate achieved by the multi-homing DeepRAT over the multi-mode, random, and fixed is 8.81%, 30.03%, and 22.81%, respectively. Note that although the rate achieved by the CVXPY solver is slightly higher than the proposed DeepRAT algorithm, it has the following shortcomings. 1) Unlike the multi-homing DeepRAT method, the CVXPY-based approach does not guarantee that the EDs are assigned to the optimal set of RATs as it assigns all EDs to all existing RATs. 2) The CVXPY solver requires full and instantaneous knowledge of CSI to solve the power allocation problem, whereas the proposed DeepRAT method does not. This aspect is important when we discuss the fast adaptivity of our proposed DeepRAT algorithm in the next subsection. 3) The utility values obtained by the DeepRAT method are higher than those obtained by the CVXPY solver, as observed from Figs. 3 and 4.

The optimal RATs-EDs assignment process after the convergence of all DeepRAT's models is shown in Fig. 6. It shows which EDs have been assigned to RATs and the corresponding percentages of downlink data rates achieved. The top plot shows the percentage of downlink data rates delivered by each RAT to its assigned EDs, while the bottom plot shows the percentage of data rates delivered to each ED from its assigned set of RATs. As an example, the top plot shows that the ES assigned RAT1 to six EDs, namely ED2, ED3, ED4, ED6, ED8, and ED9 (i.e., $U_1 = \{2, 3, 4, 6, 8, 9\}$). The percentages of data rates delivered to these EDs from RAT1 are 52.9%, 4.42%, 10%, 0.08%, 14.7%, and 17.9%, respectively. Also, the bottom plot shows that the ES assigned ED3 to all three RATs (i.e., $L_3 = \{1, 2, 3\}$), and the percentages of data rates delivered from RATs 1, 2, 3 are 53%, 30.7%, and 16.3%, respectively. For convenience, Table VI shows the results presented at the bottom of Fig. 6 in a tabular form.

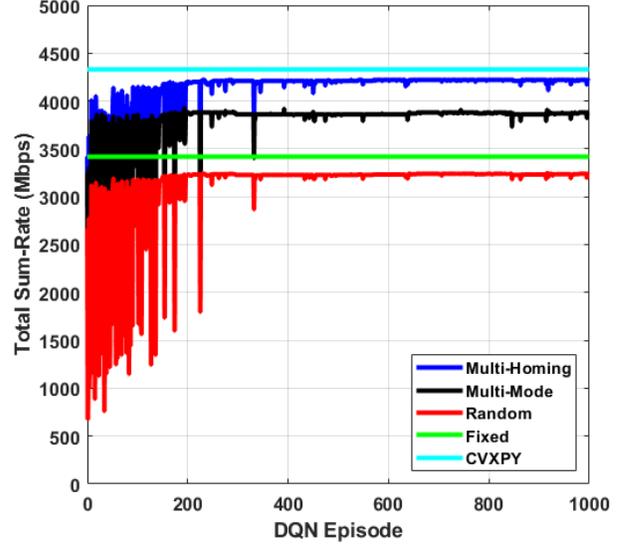Fig. 7 shows the achieved data rate for each multi-homing



Fig. 5: Total sum-rate achieved by the proposed multi-homing DeepRAT algorithm compared to the multi-mode, random, fixed, and CVXPY schemes. Although the CVXPY scheme provides a slightly better rate than DeepRAT, it does not assign EDS to the optimal set of RATs, requires full knowledge of CSI, and has inferior utility values as shown in Figs. 3 and 4.

Table VI: PERCENTAGES OF DATA RATES DELIVERED BY RATS TO THEIR ASSIGNED EDs, AS SHOWN AT THE BOTTOM OF FIG. 6.

| ED ID | RAT1 (%) | RAT2 (%) | RAT3 (%) |
|-------|----------|----------|----------|
| ED1 | 0 | 76.2 | 23.1 |
| ED2 | 99.1 | 0.36 | 0.54 |
| ED3 | 53 | 30.7 | 16.3 |
| ED4 | 76.2 | 1.3 | 22.5 |
| ED5 | 0 | 0 | 100 |
| ED6 | 60.1 | 39.9 | 0 |
| ED7 | 0 | 0 | 100 |
| ED8 | 91.8 | 0 | 8.2 |
| ED9 | 100 | 0 | 0 |
| ED10 | 0 | 0.7 | 99.3 |

ED from its assigned RAT(s). All EDs converge to the optimal rate and reach steady-state after 226 episodes. Also, when comparing these results with the minimum data rate requirements of EDs $R_u^{min}$ in Table IV, we can see that the ES assignment guarantees that all EDs satisfy their QoS data rate requirements in a manner that cost-effectively maximizes the network sum-rate. For example, the data rate requirement of ED4 is 47.8 kbps, whereas the achievable rate after connecting ED6 to RATs 1 and 2 is 503 Mbps bps, which is much greater than the required rate. Indeed, our proposed scheme can significantly enhance the performance of the HetNet by enabling data transfer from RATs to EDs in a cost-effective manner while guaranteeing satisfactory QoS.

### C. DeepRAT's Adaptivity to the Mobility of EDs

In this subsection, we demonstrate the efficiency of the DeepRAT algorithm to quickly and dynamically adapt, in terms of the convergence speed, to abrupt network changes.
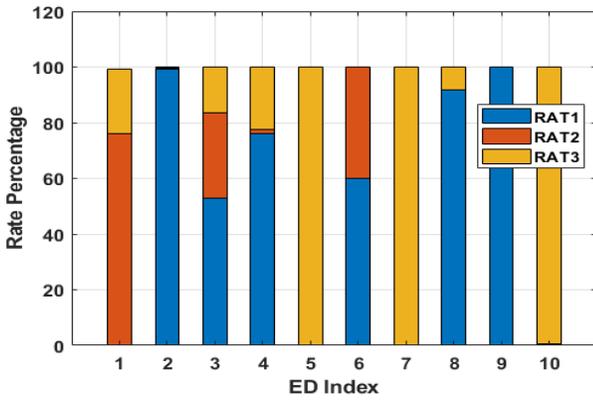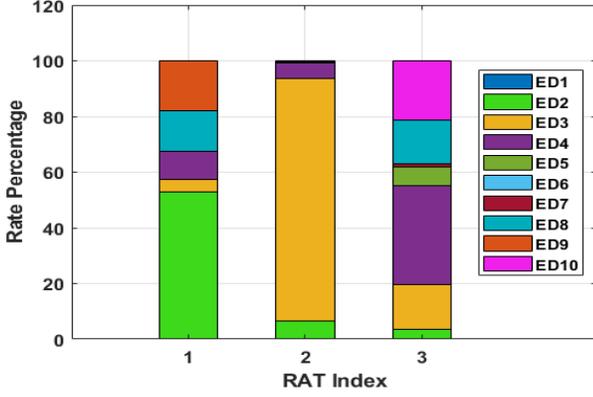
Fig. 6: Percentage of data rate delivered from each RAT to its assigned EDs (top) and achieved by each ED from its assigned RATs (bottom).
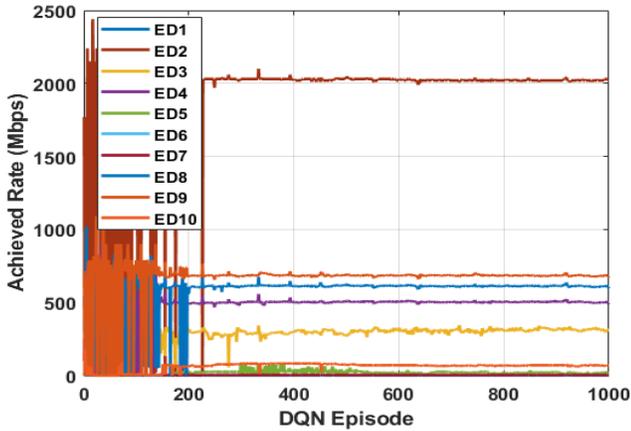


Fig. 7: Total achieved rate by each multi-homing ED from its assigned RATs. All rates are greater than the minimum required data rate requested by EDs as shown in Table IV.

During the simulation, we define the convergence speed as the number of episodes required to reach the steady-state, which
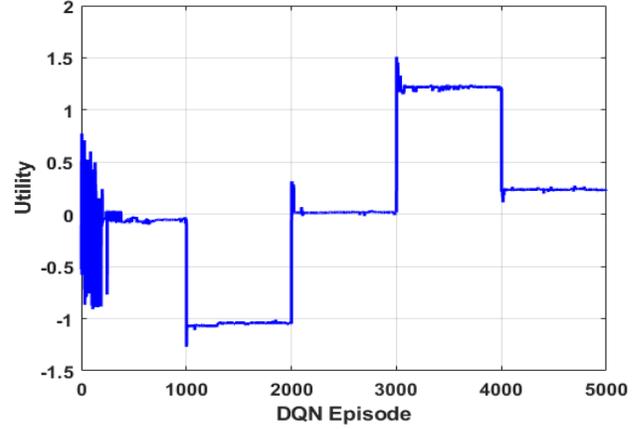


Fig. 8: Adaptivity of the DeepRAT algorithm to the dynamic mobility of EDs. The EDs randomly move at every 1000 episodes, and the DeepRAT quickly adapts to these changes. The convergence speed after training is around 2.5 faster than the one at the initial training.

is defined as having the utility values constant for the last 200 episodes [3].

We investigate a practical scenario where the EDs move randomly during each 1000 episodes. Fig. 8 shows the corresponding simulation results for the utility function. We notice that the DeepRAT algorithm adapts very quickly, in terms of the convergence speed, to the abrupt system dynamics, i.e., EDs mobility, and it dynamically finds the optimal solution of the problem in **P1**. These trends are clear at episodes 1000, 2000, 3000, and 4000, where the DeepRAT algorithm converges and reaches steady-states after 246, 1081, 2097, 3078, and 4034 episodes, respectively. The initial training phase takes around 246 episodes to converge, whereas the worst-case after the training takes only less than 97 episodes to converge, i.e., to re-solve the optimization problem **P1**. This means that the convergence speed after training is around 2.5 faster than the convergence speed at the initial training, which quantifies the dynamic adaption performance of the DeepRAT algorithm for the random changes in network dynamics.

## VI. CONCLUSION

This paper investigated the problem of cost-effective downlink sum-rate maximization in multi-RAT multi-homing HetNets. The problem was formulated as a MINLP whose objective is to cost-effectively maximize network sum-rate via jointly assigning EDs to the optimal set of RATs and allocating the optimal RATs' power levels. Due to the high complexity and combinatorial nature of the problem on the one hand and the limited knowledge of network statistics on the other hand, we proposed to solve the problem using DRL methods. Towards this goal, we proposed a multi-agent DQN and DDPG-based DRL algorithm, called DeepRAT, which solved the problem hierarchically in two stages; dynamic RATs-EDs assignment and power allocation. Our simulation results showed that the proposed multi-homing DeepRAT algorithm outperforms four benchmark heuristic algorithms in terms

of utility value. In addition, our simulation results showed the ability of the DeepRAT algorithm to quickly adapt to abrupt network changes, such as EDs' mobility, and that its convergence speed after training is around 2.5 faster than the initial training. As future work, we will extend the multi-homing DeepRAT algorithm to address the problem of joint optimization of both power and spectrum.

## ACKNOWLEDGMENT

## REFERENCES

[1] A. Alwarafy, B. S. Ciftler, M. Abdallah, and M. Hamdi, "DeepRAT: A DRL-based framework for multi-RAT assignment and power allocation in hetnets," in *2021 IEEE International Conference on Communications Workshops (ICC Workshops)*. IEEE, 2021, pp. 1–6.

[2] W. Saad, M. Bennis, and M. Chen, "A vision of 6g wireless systems: Applications, trends, technologies, and open research problems," *IEEE network*, vol. 34, no. 3, pp. 134–142, 2019.

[3] B. S. Ciftler, M. Abdallah, A. Alwarafy, and M. Hamdi, "DQN-based multi-user power allocation for hybrid RF/VLC networks," in *ICC 2021-IEEE International Conference on Communications*. IEEE, 2021, pp. 1–6.

[4] CISCO (2020) Cisco Annual Internet Report (2018-2023). White Paper. [Online]. Available: https://www.cisco.com/c/en/us/solutions/collateral/executive-perspectives/annual-internet-report/white-paper-c11-741490.pdf

[5] F. Hussain, S. A. Hassan, R. Hussain, and E. Hossain, "Machine learning for resource management in cellular and IoT networks: Potentials, current solutions, and open challenges," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 2, pp. 1251–1275, 2020.

[6] A. Alwarafy, M. Abdallah, B. S. Ciftler, A. Al-Fuqaha, and M. Hamdi, "Deep reinforcement learning for radio resource allocation and management in next generation heterogeneous wireless networks: A survey," *arXiv preprint arXiv:2106.00574*, 2021.

[7] Z. Chkirbene, A. Awad, A. Mohamed, A. Erbad, and M. Guizani, "Deep reinforcement learning for network selection over heterogeneous health systems," *IEEE Transactions on Network Science and Engineering*, 2021.

[8] N. C. Luong, D. T. Hoang, S. Gong, D. Niyato, P. Wang, Y.-C. Liang, and D. I. Kim, "Applications of deep reinforcement learning in communications and networking: A survey," *IEEE Communications Surveys & Tutorials*, vol. 21, no. 4, pp. 3133–3174, 2019.

[9] A. A. Khan and R. S. Adve, "Centralized and distributed deep reinforcement learning methods for downlink sum-rate optimization," *IEEE Transactions on Wireless Communications*, vol. 19, no. 12, pp. 8410–8426, 2020.

[10] H. Ding, F. Zhao, J. Tian, D. Li, and H. Zhang, "A deep reinforcement learning for user association and power control in heterogeneous networks," *Ad Hoc Networks*, vol. 102, p. 102069, 2020.

[11] H. Ye, G. Y. Li, and B.-H. F. Juang, "Deep reinforcement learning based resource allocation for V2V communications," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 4, pp. 3163–3173, 2019.

[12] Z. Bi and W. Zhou, "Deep reinforcement learning based power allocation for D2D network," in *2020 IEEE 91st Vehicular Technology Conference (VTC2020-Spring)*. IEEE, 2020, pp. 1–5.

[13] Z. Zhang, H. Qu, J. Zhao, and W. Wang, "Deep reinforcement learning method for energy efficient resource allocation in next generation wireless networks," in *Proceedings of the 2020 International Conference on Computing, Networks and Internet of Things*, 2020, pp. 18–24.

[14] Y. Y. Munaye, R.-T. Juang, H.-P. Lin, G. B. Tarekegn, and D.-B. Lin, "Deep reinforcement learning based resource management in UAV-assisted IoT networks," *Applied Sciences*, vol. 11, no. 5, p. 2163, 2021.

[15] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.

[16] K. Yang, C. Shen, and T. Liu, "Deep reinforcement learning based wireless network optimization: A comparative study," in *IEEE INFOCOM 2020-IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS)*. IEEE, Conference Proceedings, pp. 1248–1253.

[17] Y. S. Nasir and D. Guo, "Deep reinforcement learning for joint spectrum and power allocation in cellular networks," *arXiv preprint arXiv:2012.10682*, 2020.

[18] Q. Liu, T. Han, N. Zhang, and Y. Wang, "Deepslicing: Deep reinforcement learning assisted resource allocation for network slicing," in *GLOBECOM 2020-2020 IEEE Global Communications Conference*. IEEE, 2020, pp. 1–6.

[19] J. Kong, Z.-Y. Wu, M. Ismail, E. Serpedin, and K. A. Qaraqe, "Q-learning based two-timescale power allocation for multi-homing hybrid rf/vlc networks," *IEEE Wireless Communications Letters*, vol. 9, no. 4, pp. 443–447, 2020.

[20] F. B. Mismar, B. L. Evans, and A. Alkhateeb, "Deep reinforcement learning for 5G networks: Joint beamforming, power control, and interference coordination," *IEEE Transactions on Communications*, vol. 68, no. 3, pp. 1581–1592, 2020.

[21] I. . W. Group *et al.*, "IEEE standard for local and metropolitan area networks-part 21: Media independent handover," *IEEE Std 802.21-2008*, pp. c1–c301, 2009.

[22] A. Zhou, B.-Y. Qu, H. Li, S.-Z. Zhao, P. N. Suganthan, and Q. Zhang, "Multiobjective evolutionary algorithms: A survey of the state of the art," *Swarm and evolutionary computation*, vol. 1, no. 1, pp. 32–49, 2011.

[23] S. Boyd, S. P. Boyd, and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.

[24] A. A. Abdellatif, M. S. Allahham, A. Mohamed, A. Erbad, and M. Guizani, "ONSRA: an optimal network selection and resource allocation framework in multi-RAT systems," in *ICC 2021-IEEE International Conference on Communications*. IEEE, 2021, pp. 1–6.

[25] N. Zhao, Y.-C. Liang, D. Niyato, Y. Pei, M. Wu, and Y. Jiang, "Deep reinforcement learning for user association and resource allocation in heterogeneous cellular networks," *IEEE Transactions on Wireless Communications*, vol. 18, no. 11, pp. 5141–5152, 2019.

[26] Q. Zhang, Y.-C. Liang, and H. V. Poor, "Intelligent user association for symbiotic radio networks using deep reinforcement learning," *IEEE Transactions on Wireless Communications*, vol. 19, no. 7, pp. 4535–4548, 2020.

[27] Z. Li, C. Wang, and C.-J. Jiang, "User association for load balancing in vehicular networks: An online reinforcement learning approach," *IEEE Transactions on Intelligent Transportation Systems*, vol. 18, no. 8, pp. 2217–2228, 2017.

[28] Y. Chen, J. Li, W. Chen, Z. Lin, and B. Vucetic, "Joint user association and resource allocation in the downlink of heterogeneous networks," *IEEE Transactions on Vehicular Technology*, vol. 65, no. 7, pp. 5701–5706, 2016.

[29] Y. Choi, H. Kim, S.-w. Han, and Y. Han, "Joint resource allocation for parallel multi-radio access in heterogeneous wireless networks," *IEEE Transactions on Wireless Communications*, vol. 9, no. 11, pp. 3324–3329, 2010.

[30] J. Miao, Z. Hu, C. Wang, R. Lian, and H. Tian, "Optimal resource allocation for multi-access in heterogeneous wireless networks," in *2012 IEEE 75th Vehicular Technology Conference (VTC Spring)*. IEEE, Conference Proceedings, pp. 1–5.

[31] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, 2015.

[32] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[33] S. Diamond and S. Boyd, "CVXPY: A Python-embedded modeling language for convex optimization," *Journal of Machine Learning Research*, vol. 17, no. 83, pp. 1–5, 2016.