

Resource Allocation in an Open RAN System using Network Slicing

Mojdeh Karbalaee Motaleb*, Vahid Shah-Mansouri*, Saeede Parsaeefard[†], and Onel Luis Alcaraz López[‡]

Email: {mojdeh.karbalaee, vmansouri}@ut.ac.ir, {saeede.parsaeefard}@gmail.com, {onel.alcarazlopez}@oulu.fi

*School of ECE, University of Tehran, Tehran, Iran, [†]University of Toronto, Toronto, Canada, [‡]University of Oulu, Oulu, Finland

Abstract—The next radio access network (RAN) generation, open RAN (O-RAN), aims to enable more flexibility and openness, including efficient service slicing, and to lower the operational costs in 5G and beyond wireless networks. Nevertheless, strictly satisfying quality-of-service requirements while establishing priorities and promoting balance between the significantly heterogeneous services remains a key research problem. In this paper, we use network slicing to study the service-aware baseband resource allocation and virtual network function (VNF) activation in O-RAN systems. The limited fronthaul capacity and end-to-end delay constraints are simultaneously considered. Optimizing baseband resources includes O-RAN radio unit (O-RU), physical resource block (PRB) assignment, and power allocation. The main problem is a mixed-integer non-linear programming problem that is non-trivial to solve. Consequently, we break it down into two different steps and propose an iterative algorithm that finds a near-optimal solution. In the first step, we reformulate and simplify the problem to find the power allocation, PRB assignment, and the number of VNFs. In the second step, the O-RU association is resolved. The proposed method is validated via simulations, which achieve a higher data rate and lower end-to-end delay than existing methods.

Index Terms—open radio access network (O-RAN), virtual network function (VNF), network slicing, knapsack problem, greedy algorithm, Karush-Kuhn-Tucker (KKT) Conditions.

I. INTRODUCTION

Network slicing is a key technology in 5G wireless systems. Network slicing isolates network resources into slices, e.g., via core slicing and/or radio access network (RAN) slicing, for serving various service classes. [1]–[3].

There are three main service classes in 5G, namely enhanced mobile broadband (eMBB), ultra-reliable low latency communications (URLLC), and massive machine-to-machine communications (mMTC). Each service is assigned to a network slice depending on its corresponding quality of service (QoS) requirements. For instance, the eMBB service demands high capacity and throughput, e.g., 8K video streaming and immersive gaming. Meanwhile, the URLLC service provides ultra-reliable and low-latency connectivity, e.g., for autonomous vehicles, Tactile Internet, and remote surgeries. Finally, mMTC services require connectivity for a large number of Internet of Things (IoT) devices that transmit small payloads [4]–[6].

Radio access networks (RANs) currently lack adequate flexibility and openness to handle these simultaneous service demands. Given this, a new RAN paradigm, called open RAN (O-RAN), schematically represented in Fig. 1, is introduced to deal with these issues.

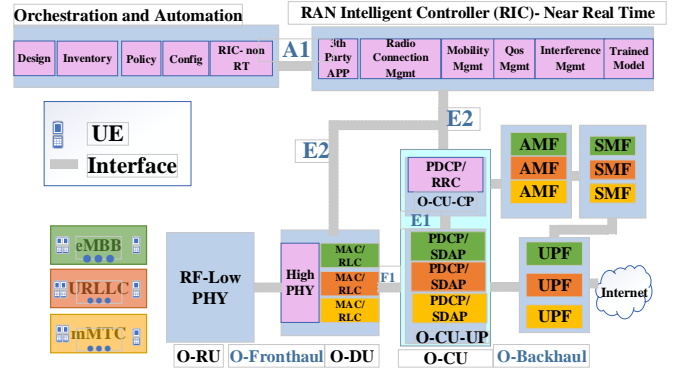


Fig. 1: Network sliced O-RAN system

The O-RAN architecture separates hardware and software, enabling network function virtualization (NFV). Additionally, each component is implemented as a virtual network function (VNF), the system function block in NFV, that can be deployed on a virtual machine (VM) or container [7]. As a result, some O-RAN components defined in Section III are virtualized and implemented as VNFs.

A. Motivation

The optimal resource allocation for the 5G and 6G systems is crucial to reducing costs and improving user equipment (UE). Significant challenges face these systems, including interference alignment, limited capacity of the fronthaul links, energy restrictions on VMs, etc [8], [9].

Many studies have investigated resource allocation in cloud RAN (C-RAN) by considering a single service's power, data rate, and delay limitations. This architectural model is inefficient in 5G and 6G since it needs to serve services with different QoS simultaneously. Therefore, O-RAN can simultaneously support multiple services at a lower cost by being flexible, layered, and modular. Balancing different services with different QoS, resource requirements, and priorities is a fundamental problem [1], [10]–[12].

This paper aims to design an O-RAN architecture to support the three types of 5G services, namely, eMBB, URLLC, and mMTC, via network slicing and resource allocation. Specifically, we intend to maximize the total achievable data rate and satisfy the minimum achievable data rate for eMBB services while meeting the URLLC requirements in the presence of numerous low-power IoT devices.

B. Main Contributions

This paper studies the resource utilization of a downlink O-RAN system to develop an isolated network slicing outline for the three 5G services. We use mathematical methods to decompose and convexify the problem and solve it using hierarchical algorithms. The main contributions of this paper are summarized as follows:

- We formulate a problem for allocating baseband resources such as power, physical resource blocks (PRBs), O-RUs, and activating VNFs to maximize the weighted throughput of the O-RAN architecture. The three types of 5G service classes, i.e., eMBB, URLLC, and mMTC, are considered in this system. We take into account their corresponding QoS requirements and service priorities.
- We propose a two-step resource management algorithm for solving the optimization problem. In the first step, we reformulate and simplify the problem so as to find an upper and lower bound for the number of activated VNFs. Moreover, we use the Lagrangian function and Karush-Kuhn-Tucker (KKT) conditions to obtain the optimal power and PRB allocation. In the second step, the problem of O-RU association is converted to a multiple knapsack problem and solved by a greedy algorithm.
- We provide insights into the complexity of the proposed algorithms and demonstrate their convergence. Additionally, we analyze the feasibility region of the problem and introduce a fast algorithm to check it numerically.
- We show via numerical results that the proposed algorithm outperforms two baseline schemes in terms of achievable data rate and mean total delay. Remarkably, the proposed algorithm performs close to the optimal solution in low-interference conditions.

Table I lists the acronyms used throughout this paper, which is organized as follows. Relevant literature related to our work is discussed in Section II, while Section III briefly overviews the O-RAN architecture. The system model and the problem formulation are described in Section IV, V, respectively. The details of our proposed resource management algorithm are introduced in Section VI. In Section VII, numerical results are provided to evaluate the performance of the proposed algorithm. Finally, Section VIII concludes the paper.

II. RELATED LITERATURE

The network slicing problem in multi-tenant cellular networks has received significant attention recently, e.g., [8], [13], [14]. Two levels of dynamic network slicing in heterogeneous C-RAN (H-CRAN) are examined in [8]. The higher level manages user acceptance control, RRH association, and the allocation of BBU capacity. Meanwhile, PRB and power are allocated at lower levels. In [15], RAN slicing is considered for the fog RAN (F-RAN) system and executed using deep reinforcement learning. In [16], [17],

TABLE I: List of Acronyms

Acronym	Definition
VNF	virtual network function
VM	virtual machine
RAN	radio access network
O-RAN	open RAN
vRAN	virtual RAN
CRAN	cloud RAN
RRH	radio remote head
BBU	baseband unit
QoS	quality of service
MIMO	multiple input multiple output
PRB	physical resource block
eMBB	enhanced mobile broadBand
URLLC	ultra-reliable low latency communication
mMTC	massive machine-to-machine communications
O-RU	O-RAN radio unit
O-DU	O-RAN distributed unit
O-CU	O-RAN central unit
UPF	user plane function
UE	user equipment
SINR	signal-to-noise-plus-interference ratio
CAPEX	capital expenditures
OPEX	operating expenses
KKT	Karush-Kuhn-Tucker

the implementation of RAN level slicing is discussed among multiple mobile network operators.

Multiplexing eMBB and URLLC services on the same RAN is particularly challenging, which has motivated some recent work. In [18], a RAN slicing is considered in a coordinated multipoint system to ensure eMBB and URLLC QoS requirements. Moreover, [19] investigates the minimization of the system's power for the RAN slicing of eMBB and URLLC downlink services using non-orthogonal multiple access techniques. In [5], a deep reinforcement learning algorithm is used to solve the resource allocation problem for eMBB and URLLC services.

In [9], [20], VMs activation and beamforming allocation are discussed in C-RAN systems. Paper [9] minimizes energy cost with system delay, fronthaul capacity, and rate constraints. To guarantee UE delays, M/M/1 queueing theory is used for transmission and processing delays. In [21], [22], the problem of joint virtual computing resource allocation with beamforming is formulated. Also, the association of RRH to the UE is considered and solved using innovative methods.

In [23], [24], the problem of joint power allocation and RRH association in a H-CRAN system is considered to maximize the energy efficiency. Finally, in [25], the optimum power and association of RRH to BBU and RRH to UE are obtained in the massive MIMO-aided C-RAN system.

III. BACKGROUND

The O-RAN architecture is a flexible, open, low-cost, and intelligent alternative to future RANs. O-RAN combines

features of virtual RAN (vRAN) and cloud RAN (C-RAN). By virtualizing RANs, operators can improve flexibility, reduce CAPEX and OPEX, and add new capabilities to their networks more quickly. The C-RAN architecture consists of two parts: the radio remote head (RRH) and the baseband unit (BBU). Several distributed RRHs can be connected to a centralized BBU, called BBU-pool [26]. Unlike C-RAN, O-RAN separates RAN into three different units, namely Radio Unit (O-RU), Distributed Unit (O-DU), and Central Unit (O-CU). **Mostly non-real-time baseband processing occurs in the O-CU layer, while real-time baseband processing occurs in the O-DU layer.**

In the O-RAN architecture, the PHY is divided into low and high PHY, unlike C-RAN. As shown in Fig. 1, O-RU is a logical node that contains RF and low PHY. The former transmits or receives radio signals, while the latter includes digital beamforming. Typically, the O-DU constitutes a logical node with high PHY, MAC, and RLC. It contains a subfunction of the eNodeB and is deployed near the O-RU. Moreover, O-DU is connected to an O-RU with an open fronthaul interface. In addition to supporting the lower layers of the protocol stack, O-CU also provides support for the higher layers. The O-CU contains two parts: the O-CU user plane (O-CU-UP) and the O-CU control plane (O-CU-CP). The former hosts the packet data convergence protocol (PDCP)-UP and the service data adaptation protocol (SDAP), while the latter hosts PDCP-CP and radio resource control (RRC). O-DU and O-CU are connected via an open and well-defined interface F₁. Moreover, O-CU-UP is connected to user plane function (UPF) via O-backhaul link. The O-RAN architecture contains other principal logical nodes called Orchestration and Automation, RAN Intelligent Controller (RIC)- Near Real-Time, and O-Cloud. Orchestration and Automation include functions such as RIC Non-Real-Time. RIC is responsible for machine learning methods and making the system more intelligent. **A key feature of the O-RAN architecture is that the hardware is disaggregated from the software, leading to network function virtualization (NFV).** Additionally, each component is implemented as a virtual network function (VNF), the system function block in NFV, that can be deployed on a virtual machine (VM) or container [7]. As a result, some O-RAN components such as user plane function (UPF), O-CU, O-DU, and RIC-near real-time, are virtualized and implemented as VNFs. [10]–[12], [27]–[30].

IV. SYSTEM MODEL

We consider a downlink (DL) system, and an O-RAN architecture using RAN slicing as depicted in Fig. 1. In this section, we present the system and signal model, derive the achievable data rates, power of O-RU, and the fronthaul capacity of the O-RAN system. Moreover, we discuss the mean delay and the power of VNFs.

A. System Model

Assume, there are three service types: eMBB, URLLC, and mMTC, which support different applications. Accordingly, there are S_1 slices for the first service type

(eMBB), S_2 slices for the second service type (URLLC), and S_3 slices for the third service type (mMTC). Therefore, there are $S = S_1 + S_2 + S_3$ pre-allocated slices serving these services. Moreover, each service request $s \in \{1, \dots, S\}$ is served by its corresponding slice. So we have the set $\{1, 2, \dots, S_1\}$ of eMBB service instances, the set $\{1, 2, \dots, S_2\}$ of URLLC service instances, and the set $\{1, 2, \dots, S_3\}$ of mMTC service instances. Each service $s_j \in \{1, 2, \dots, S_j\}$ consists of U_{s_j} requests from single-antenna UEs requiring certain level of QoS. Notice that $j \in \{1, 2, 3\}$ indicates the service type. Based on the application and QoS request, UE may be admitted and allocated to the resources.

Each pre-allocated slice contains reserved VNFs for the three logical nodes:

- MAC/RLC functions in the O-DU
- PDCP/SDAP functions in the O-CU-UP
- UPF which is a functional layer

Each slice $s \in \{1, 2, \dots, S\}$, consists of M_s^d VNFs for the processing of O-DU, M_s^c VNFs for the processing of O-CU-UP, and M_s^u VNFs for the processing of UPF. The VNFs of O-DU, O-CU-UP, and UPF are interconnected, which is defined as the service function chain in the O-RAN system. Also, each VNF instance runs on a VM that uses resources from the data centers.

Assume there are K PRBs in this system. Suppose each slice s consists of \bar{K}_s pre-allocated virtual resource blocks that are mapped to PRBs. Therefore, we have $\sum_s \bar{K}_s \leq K$. In addition, there are R multi-antenna O-RUs that are shared between the slices. Specifically, the O-RU $r \in \mathcal{R} = \{1, 2, \dots, R\}$ has J antennas for transmitting and receiving data. Moreover, all O-RUs have access to all PRBs.

B. Signal Model

Let $y_{u(s,i)}$ be the received signal of UE i in the s^{th} service such that

$$y_{u(s,i)} = \sum_{r=1}^R \sum_{k=1}^{K_s} \mathbf{h}_{r,u(s,i)}^H g_{u(s,i)}^r e_{r,u(s,i)}^k x_{Q,r,u(s,i)}^k + z_{u(s,i)}, \quad (1)$$

where $x_{Q,r,u(s,i)}^k = x_{P,r,u(s,i)}^k + \mathbf{q}_r \cdot x_{P,r,u(s,i)}^k = \mathbf{w}_{r,u(s,i)}^k \sqrt{p_{r,u(s,i)}^k} x_{u(s,i)}^k$, and $x_{u(s,i)}^k$ depicts the transmitted symbol vector, $z_{u(s,i)} \sim \mathcal{CN}(0, BN_0)$ is the receive additive Gaussian noise and BN_0 is the noise power in a given bandwidth B . Here, x_P denotes the precoded message before compression, while x_Q illustrates the precoded message after compression. In addition, $\mathbf{q}_r \sim \mathcal{CN}(0, \sigma_q^2 \mathbf{I}_R)$ indicates the quantization Gaussian noise which comes from the signal compression in O-DU. Furthermore, $g_{u(s,i)}^r \in \{0, 1\}$ is a binary variable that illustrates whether O-RU r serves the i^{th} UE that is allocated to the s^{th} slice or not. Furthermore, $p_{r,u(s,i)}^k$ represents the transmission power of the O-RU r serve the i^{th} UE in slice s and PRB k , while $\mathbf{h}_{r,u(s,i)}^k \in \mathbb{C}^J$ corresponding channel vector. In addition, $\mathbf{w}_{r,u(s,i)}^k \in \mathbb{C}^J$ depicts the associated transmit beamforming

vector. Therefore, the SINR of the i^{th} UE served at slice s on PRB k is given by

$$\rho_{r,u(s,i)}^k = \frac{p_{r,u(s,i)}^k |\mathbf{h}_{r,u(s,i)}^k H \mathbf{w}_{r,u(s,i)}^k|^2}{BN_0 + I_{r,u(s,i)}^k}, \quad (2)$$

A UE in an O-RU r using PRB k receives interference from other O-RUs in the set $\mathcal{R} \setminus r$ that are using the same PRB k . Two types of interference occur between UEs in each slice: i) inter-slice interference between signals transmitted over different slices, and ii) intra-slice interference between signals transmitted over the same slice.

Network slicing techniques significantly reduce inter-service (inter-slice) interference. One way to leverage a two-time scale PRB scheduling is to isolate PRBs in slices (in the first time scale) and schedule the PRBs to the UEs of the slices (in the second time scale). Another method consists of allocating part of the PRBs of eMBB services to URLLC and mMTC [2], [5], [31]. In this paper, we assume that the PRB scheduling is performed. Also, in Section V-A, we briefly study the PRB scheduling between slices. Since there are limited resources, inter-service interference cannot be eliminated entirely. Nevertheless, isolating the slices reduces inter-service interference considerably and allows us to ignore it mathematically.

Back to (2), $I_{r,u(s,i)}^k$ is the sum of the power of interfering signals and quantization noise, and can be represented as

$$I_{r,u(s,i)}^k = \underbrace{\sum_{j=1}^R \sigma_q^2 |\mathbf{h}_{r,u(s,i)}^k|^2}_{\text{(quantization noise)}} + \underbrace{\sum_{\substack{l=1 \\ l \neq i}}^{U_s} e_{u(s,i)}^k e_{u(s,l)}^k p_{u(s,l)}^k \sum_{\substack{r'=1 \\ r' \neq r}}^R |\mathbf{h}_{r',u(s,i)}^k H \mathbf{w}_{r',u(s,i)}^k g_{u(s,l)}^{r'}|^2}_{\text{(intra-slice interference)}}, \quad (3)$$

where $e_{u(s,i)}^k$ is a binary variable that indicates whether the k^{th} PRB is allocated to the UE i in slice s , assigned to r^{th} O-RU, or not. Furthermore, there is no inter-slice interference, only intra-slice interference, since slices are assumed to be isolated.

Herein, we consider a zero forcing beamforming vector, which minimizes the experienced intra-slice interference, and is given by [32]

$$\mathbf{w}_{r,u(s,i)}^k = \hat{\mathbf{h}}_{r,u(s,i)}^k (\hat{\mathbf{h}}_{r,u(s,i)}^k H \hat{\mathbf{h}}_{r,u(s,i)}^k)^{-1}. \quad (4)$$

where $\mathbf{h}_{r,u(s,i)}^k$ is the channel estimate, which is assumed imperfect. Mathematically, $\hat{\mathbf{h}}_{r,u(s,i)}^k = \mathbf{h}_{r,u(s,i)}^k + \Delta \mathbf{h}_{r,u(s,i)}^k$, where $\Delta \mathbf{h}_{r,u(s,i)}^k \sim \mathcal{N}(0, \phi_{r,u(s,i)}^2)$ indicates the estimating error vector with a Gaussian distribution and $\phi_{r,u(s,i)}^2 = \text{diag}(\phi_{r,u(s,i)}^2, \dots, \phi_{r,u(s,i)}^2)$.

C. Achievable Data Rate

The achievable data rate for the i^{th} UE request in the s_1^{th} application of service type 1 (eMBB) can be written as

$$\mathcal{R}_{u(s_1,i)} = \sum_{r=1}^R \mathcal{R}_{r,u(s_1,i)} g_{u(s_1,i)}^r, \quad (5)$$

where

$$\mathcal{R}_{r,u(s_1,i)} = \sum_{k=1}^K \mathcal{R}_{r,u(s_1,i)}^k e_{r,u(s_1,i)}^k, \quad (6)$$

is the achievable data rate of RU r to UE i in slice s_1 , which depends on the achievable data rate per PRB, i.e.,

$$\mathcal{R}_{r,u(s_1,i)}^k = B \log_2(1 + \rho_{r,u(s_1,i)}^k), \quad (7)$$

Since the blocklength in URLLC and mMTC is finite, the achievable data rate for the i^{th} UE request in the application of service type 2 (URLLC) and 3 (mMTC) is not achieved from the Shannon capacity formula. Instead, in a short packet transmission, the achievable data rate is approximated as [2]

$$\mathcal{R}_{u(s_j,i)} = \sum_{r=1}^R \mathcal{R}_{u(s_j,i)}^r g_{u(s_j,i)}^r, \quad (8)$$

where

$$\mathcal{R}_{r,u(s_j,i)} = \mathcal{R}_{r,u(s_j,i)}^k e_{u(s_j,i)}^k, \quad (9)$$

is the achievable data rate of RU r to UE i in slice s_1 , which depends on the achievable data rate per PRB, i.e.,

$$\mathcal{R}_{r,u(s_j,i)}^k = B \log_2(1 + \rho_{r,u(s_j,i)}^k - \zeta_{u(s_j,i)}^k) e_{u(s_j,i)}^k, \quad (10)$$

where

$$\zeta_{u(s_j,i)}^k = \log_2(e) Q^{-1}(\epsilon) \sqrt{\frac{\mathfrak{C}_{u(s_j,i)}^k}{N_{u(s_j,i)}^k}}. \quad (11)$$

Here, ϵ is the transmission error probability, Q^{-1} is the inverse of the Q function, $\mathfrak{C}_{u(s_j,i)}^k = 1 - \frac{1}{(1 + \rho_{u(s_j,i)}^k)^2}$ depicts the channel dispersion of UE i at slice s_j and PRB k , while $N_{u(s_j,i)}^k$ represents the corresponding transmit blocklength. $\mathcal{R}_{r,u(s_j,i)}$ is the achievable data rate that is transmitted by O-RU r to UE i requesting service s_j .

If we replace $p_{u(s,l)}^k$ and $p_{u(n,l)}^k$ in (3) by P_s^{\max} , an upper bound $\bar{I}_{r,u(s,i)}^k$ is obtained for $I_{r,u(s,i)}^k$. Therefore, $\bar{\mathcal{R}}_{u(s,i)} \forall s, i$ is derived by using $\bar{I}_{r,u(s,i)}^k$ instead of $I_{r,u(s,i)}^k$ in (8) and (5).

D. Power of the O-RU and the Fronthaul Capacity

Let P_r denote the power of the transmitted signal from the r^{th} O-RU to all the UEs served by it. From (1), the power of each O-RU r is obtained as follows,

$$P_r = \sum_{s=1}^S \sum_{k=1}^K \sum_{i=1}^{U_s} |\mathbf{w}_{r,u(s,i)}^k|^2 \alpha_{r,u(s,i)}^k + \sigma_q^2, \quad (12)$$

where $\alpha_{r,u(s,i)}^k = p_{r,u(s,i)}^k g_{u(s,i)}^r e_{r,u(s,i)}^k$. Since we have a fiber link between O-RU and O-DU, the rate of users

on the fronthaul link between O-DU and the r^{th} O-RU is formulated as

$$C_r = \log \left(1 + \frac{\sum_{s=1}^S \sum_{k=1}^{K_s} \sum_{i=1}^{U_s} |\mathbf{w}_{r,u(s,i)}^k|^2 \alpha_{r,u(s,i)}^k}{\sigma_q^2} \right),$$

$$C_r = \log_2 \left(\frac{P_r}{\sigma_q^2} \right). \quad (13)$$

E. Mean Delay

In this part, the end-to-end mean delay for each service is obtained. The total delay (T^{tot}), is the sum of the processing delay (T^{proc}), the transmission delay (T^{tr}), and the total propagation delay (T^{pro}).

$$T^{\text{tot}} = T^{\text{proc}} + T^{\text{tr}} + T^{\text{pro}}, \quad (14a)$$

$$T^{\text{proc}} = T^{\text{RU}} + T^{\text{DU}} + T^{\text{CU}} + T^{\text{UPF}}, \quad (14b)$$

$$T^{\text{tr}} = T^{\text{fr},t} + T^{\text{mid},t} + T^{\text{b},t}, \quad (14c)$$

$$T^{\text{pro}} = T^{\text{fr},p} + T^{\text{mid},p} + T^{\text{b},p}. \quad (14d)$$

Mathematically, the total propagation delay (T^{pro}) is the sum of the propagation delay in the fronthaul link $T^{\text{fr},p}$, the midhaul link $T^{\text{mid},p}$, and the backhaul link $T^{\text{b},p}$. In each link, the propagation delay is the time a signal to reach its destination. It is obtained based on the length of the fiber link and the capacity of the link (as $T = L/c$, where L is the length of the link and c is the propagation speed of the medium). Meanwhile, the total transmission delay (T^{tr}) is the sum of the transmission delay in the fronthaul $T^{\text{fr},t}$, the midhaul $T^{\text{mid},t}$, and the backhaul $T^{\text{b},t}$. In each link, the transmission delay is the amount of time required to push all the packets into the transmission medium, and can be formulated as $T = \frac{\alpha}{R}$, where R is the data-rate of the packet and α is the mean packet size. Notice that taking the propagation and transmission delays into account in the formulation is straightforward, but we have avoided it for the sake of succinctness and simplicity. Therefore, the propagation delay is fixed and does not affect the optimization problem.

Next, we present a brief calculation of propagation delay. Assume a distance between the O-RU and O-DU around 10 km, the distance between O-DU and O-CU around 80 km, not greater than the distance from O-CU to the network around 200 km [33]. Then, assuming the fronthaul, midhaul and backhaul are connected with fiber optics and c is the speed of light, the propagation delay is about $T^{\text{pro}} = (10 + 80 + 200) \times 10^3 / (3 \times 10^8) < 1$ ms.

The following is a brief calculation of transmission delay to show that it does not affect also the optimization since its contribution to the total delay is negligible. In URLLC and mMTC, the mean packet size can be between 20 to 32 bytes; Also, the minimum data rate is assume to be $1\text{bps}/\text{Hz} \times BW (180\text{KHz})$. So the transmission delay from O-RU to O-DU is about $T^{\text{fr},t} = \frac{20 \times 8}{1 \times 180 \times 10^3} < 0.1\text{ms}$. As a result, the $T^{\text{fr},t} \approx T^{\text{mid},t} \approx T^{\text{b},t}$. For eMBB, the packet size can be 100 times larger and the delay is not exceed the 0.6ms.

Therefore, in the following, we assume that the total delay is approximate to the processing delay ($T^{\text{tot}} \approx T^{\text{proc}}$).

1) *Processing Delay*: Assume the packet arrival of UEs follows a Poisson process with arrival rate $\lambda_{u(s,i)}$ for the i^{th} UE of the s^{th} service (or slice). Therefore, the mean arrival data rate of the s^{th} slice in the UPF layer is $\alpha_s^U = \sum_{u=1}^{U_s} \lambda_{u(s,i)}$. Assume the mean arrival data rate of the UPF layer for slice s (α_s^U) is approximately equal to the mean arrival data rate of the O-CU-UP layer (α_s^C) and the O-DU (α_s^D), i.e., $\alpha_s = \alpha_s^U \approx \alpha_s^C \approx \alpha_s^D$. This is because the amount of data transferred along the route (regardless of frame changes) is constant. In fact, according to Burke's theorem, the mean arrival data rate of the second and third layers, which are processed in the first layer, is still Poisson with rate α_s . It is assumed that there are load balancers in each layer for each service to divide equally the incoming traffic to VNFs. Suppose the baseband processing of each VNF is modeled by an M/M/1 processing queue. Each packet is processed by one of the VNFs of the corresponding slice. Therefore, the mean delay for the s^{th} slice in the O-DU, the O-CU, and the UPF is modeled as M/M/1 queue, and can be respectively [9], [21], [22],

$$T_s^{\text{DU}} = \frac{1}{\mu_s^d - \alpha_s / M_s^d}, \quad (15)$$

$$T_s^{\text{CU}} = \frac{1}{\mu_s^c - \alpha_s / M_s^c}, \quad (16)$$

$$T_s^{\text{UPF}} = \frac{1}{\mu_s^u - \alpha_s / M_s^u}, \quad (17)$$

where M_s^d , M_s^c and M_s^u represent the number of VNFs in O-DU, O-CU-UP and UPF, respectively. Moreover, $1/\mu_s^d$, $1/\mu_s^c$, and $1/\mu_s^u$ are the mean service time of the O-DU, O-CU, and the UPF layers, respectively. The arrival rate of each VNF in each layer for each slice s is α_s / M_s^i $i \in \{d, c, u\}$.

On the other hand, arrival data rate of wireless link for each UE i of service s is $\lambda_{u(s,i)}$, thus $\sum_{i=1}^{U_s} \lambda_{u(s,i)} = \alpha_s$. Moreover, the service time of transmission queue for UE i requesting service s has an exponential distribution with mean $1/R_{u(s,i)}$ and can be modeled as a M/M/1 queue [9], [21], [22]. Therefore, the mean delay of the transmission layer for UE i in slice s is

$$T_{u(s,i)}^{\text{RU}} = \frac{1}{R_{u(s,i)} - \lambda_{u(s,i)}}. \quad (18)$$

we assume $T_{u(s,i)}^{\text{tot}} \approx T_{u(s,i)}^{\text{proc}}$.

F. VNF Power

Assume the power consumption of each VNF in each logical node (O-DU, O-CU, and UPF) in the slice s , is represented by ϕ_s^d , ϕ_s^c , and ϕ_s^u , respectively. Then, the system's total cost of energy of all the slices can be represented as $\phi_{\text{tot}} = \sum_{s=1}^S \phi_s$,

A significant issue facing the industry is reducing energy consumption. Data centers are one of the most energy-consuming. As a result, restrictions are placed on data centers' energy, including VMs. So, one of our goals is

to limit the energy consumption of total VNFs that can be run as VM on data centers. So, by applying a custom policy on total power consumption, we can control data centers' power consumption ($\phi^{\text{tot}} \leq \phi^{\text{max}}$).

V. PROBLEM STATEMENT

Suppose the slice s (which is assigned to service s) has a priority factor δ_s (based on the priority of its hosting service) where $\sum_{s=1}^S \delta_s = 1$. The priority factor of each slice is obtained according to the service level agreement to promote a fairness in the system. This paper aims to maximize the sum-rate of all UEs subject to QoS constraints as follows.

$$\max_{\mathbf{P}, \mathbf{E}, \mathbf{M}, \mathbf{G}} \sum_{s=1}^S \sum_{i=1}^{U_s} \delta_s \bar{\mathcal{R}}_{u(s,i)} \quad (19a)$$

$$\text{subject to } P_r \leq P_r^{\text{max}} \quad \forall r, \quad (19b)$$

$$p_{r,u(s,i)}^k \geq 0 \quad \forall i, r, s, k, \quad (19c)$$

$$p_{r,u(s,i)}^k \leq P_s^{\text{max}} \quad \forall i, r, s, k, \quad (19d)$$

$$\bar{\mathcal{R}}_{u(s,i)} \geq \mathcal{R}_s^{\text{min}} \quad \forall s, \quad (19e)$$

$$C_r \leq C_r^{\text{max}} \quad \forall r, \quad (19f)$$

$$T_{u(s,i)}^{\text{tot}} \leq T_s^{\text{max}} \quad \forall i, s, \quad (19g)$$

$$\mu_s \geq \alpha_s / M_s \quad \forall s, \quad (19h)$$

$$\bar{\mathcal{R}}_{u(s,i)} \geq \lambda_{u(s,i)} \quad \forall i, s, \quad (19i)$$

$$0 \leq M_s \leq M_s^{\text{max}} \quad \forall s, \quad (19j)$$

$$\phi^{\text{tot}} \leq \phi^{\text{max}}, \quad (19k)$$

$$\sum_r g_{u(s,i)}^r = 1 \quad \forall s, i, \quad (19l)$$

$$\sum_{k=1}^{K_s} g_{u(s,i)}^r e_{r,u(s,i)}^k \geq 1 \quad \forall s, i, r, \quad (19m)$$

$$\sum_{s=1}^S \sum_{i=1}^{U_s} g_{u(s,i)}^r e_{r,u(s,i)}^k \leq 1 \quad \forall s, i, r, \quad (19n)$$

$$g_{u(s,i)}^r \in \{0, 1\} \quad \forall s, i, \quad (19o)$$

$$e_{r,u(s,i)}^k \in \{0, 1\} \quad \forall s, i. \quad (19p)$$

where $\bar{\mathcal{R}}_{u(s,i)}$ is derived by using $\bar{I}_{r,u(s,i)}^k$ instead of $I_{r,u(s,i)}^k$ in (8) and (5). In addition, $\mathbf{P} = [p_{r,u(s,i)}^k]$, $\forall s, i, r, k$, is the four-dimensional (4D) matrix of power for UEs, $\mathbf{E} = [e_{r,u(s,i)}^k]$, $\forall s, i, r, k$ indicates the binary 4D matrix for the PRB association. Moreover, $\mathbf{G} = [g_{u(s,i)}^r]$, $\forall s, i, r$ is a binary three dimensional (3D) matrix for the O-RU association. Furthermore, $\mathbf{M} = [M_s^d, M_s^c, M_s^u]$, $\forall s$ is a matrix containing the number of VNFs in each layer of slice. Notice that (19b), (19c) and (19d) limit the power of each O-RU and UE. Also, (19e) constrains the rate of each UE requesting each type of service, i.e., eMBB, mMTC, and URLLC, to be greater than a threshold. Meanwhile, (19f) and (19g) represent the limited fronthaul capacity and the limited end-to-end delay of the received signal, respectively. (19h) and (19i) are related to the stability of the M/M/1 queue, (19j) restrictes the number of VNFs in each slice due to the

limited resources, while (19l) and (19m) guarantee that the O-RU and PRB are associated with the UE, respectively. Also, (19n) ensures that each PRB can not be assigned to more than one UE associated with the same O-RU, (19k) indicates that the fixed cost of energy of VNFs in each slice does not exceed the threshold, while (19o) and (19p) constrain \mathbf{E} and \mathbf{G} to be binary matrices.

A. PRB Scheduling

In this section, we provide a brief study on the problem of PRB scheduling which can be completed in two steps to eliminate the inter-slice interference and guarantee the isolation of slices [34]. For this, firstly, we should assign the PRBs to the slices. Secondly, we assign PRBs of slices to UEs, find the optimal number of VNFs for each slice, allocate power of UEs, and assign O-RU to UEs, which uses the proposed Algorithm VI. Suppose, $\mathcal{R}_s^{\text{min}}$, and $\mathcal{R}_s^{\text{max}}$ are the minimum data rate and maximum data rate of each UE in slice s , respectively. Firstly, we need to find the average PRB number used by the UEs in each service. Since mMTC and URLLC require usually short packet transmissions, each UE in mMTC and URLLC requires 1 PRB. So if slice s serves mMTC or URLLC services, with U_s UEs, it requires $K_s = U_s \times 1$ PRBs. For eMBB, assume the average rate of each UE in slice s serving eMBB UEs is $\bar{R}_s = B \log_2(1 + \bar{\rho}_s)$, where, $\bar{\rho}_s$ is the average SINR of UEs in slice s . Therefore, the minimum number of PRBs that slice s with U_s UEs requires is $K_s^{\text{min}} = \lceil U_s \times \frac{\bar{R}_s}{\mathcal{R}_s^{\text{max}}} \rceil$. Moreover, the maximum number of PRB that the slice s with U_s UEs requires is $K_s^{\text{max}} = \lceil U_s \times \frac{\bar{R}_s}{\mathcal{R}_s^{\text{min}}} \rceil$. Also, $K_s = (K_s^{\text{min}} + K_s^{\text{max}})/2$ is the average number of required PRBs in slice s . Our goal is to obtain the number of PRBs assigned to each slice s (\bar{K}_s). The problem can be written as

$$\max_{\bar{\mathbf{K}}_s} \sum_{s=1}^S \delta_s K_s \ln(\bar{K}_s) \quad (20a)$$

$$\text{subject to } \sum_s \bar{K}_s \leq K \quad (20b)$$

$$K_s^{\text{min}} \leq \bar{K}_s \leq K_s^{\text{max}} \quad \forall s \in S_1, \quad (20c)$$

$$\bar{K}_s \leq K_s \quad \forall s \in S_2, S_3. \quad (20d)$$

We used logarithms to assign PRBs to all slices to make them equally fair, since proportional fairness is achieved by maximizing the log utility function [34]. Equation (20b) illustrates that the sum of PRBs of slices can not exceed the maximum number of PRBs (K). Equation (20c) restricts the number of PRBs of eMBB slices and (20d) limits the number of the PRBs of URLLC and mMTC slices. By relaxing \bar{K}_s , the objective function and constraints become convex and can be solved using the Lagrangian function.

B. Slice Management

In this subsection, we will look at the life cycle of network slicing on a practical level. The goal is to examine slice management, which includes creating, managing, and

deleting slices. Network slices generally have four life cycle stages [35];

- Preparation phase: the operator plans to create a network slice instance (NSI) by designing the its template, onboarding users, and preparing the environment. Also, the evaluation of requirements is performed in this step.
- Commissioning phase: the NSI is created, and the requirements are considered and allocated to the slice.
- Operation phase: the NSIs are activated, managed, monitored (e.g., KPIs), modified, and deactivated. As the slice enters the activated phase, it is ready to support services, and as the slice exits the de-activated phase, the slice is inactive, and communication services are stopped.
- Decommissioning phase: an NSI that is decommissioned no longer exists after this phase.

Since the requirements evaluation is considered in the preparation phase, we need an algorithm to estimate the UE traffic in the system at different times. Moreover, based on this estimation, we need to evaluate resources, including the optimal number of VNFs, PRB assignment of UEs for each slice, and the total power requirements. In this phase, we use our algorithm to calculate resources after estimating the system's traffic. As shown in Fig. 1, we have three different slices for eMBB, URLLC, and mMTC. The system must prepare VNFs for MAC/RLC functions in O-DU, PDCP/SDAP functions in O-CU, UPF, SMF, and AMF functionality layers for each slice. Moreover, O-RU, high PHY in O-DU, and O-CU-CP are shared between slices. Thus, we do not require evaluating and preparing for the share environments and platforms in the network slicing cycles. Moreover, the estimation of PRB and power is needed based on the proposed algorithm. After evaluating, assessing, and preparing the resources and environments for each slice, the commissioning phase is started. In this phase, the slices are created based on the previous phase estimation. These created slices are activated in the operation phase, and the actual resources are assigned based on the proposed algorithm. It is possible to modify the slice's resources even when the evaluation changes during the operation phase. If we need to remove a slice or any service not used in a zone, the unshared resources are released in the decommissioning phase.

VI. PROPOSED ALGORITHM

In this section, we first apply some simplifications to the system; Solving the problem (19) is complicated since this is non-convex mixed-integer non-linear problem (MINLP) with a binary variable and an integer variable. We applied some simplifications and use an iterative heuristic algorithm to solve the problem. We solve this problem in two levels, iteratively, until it converges [24].

At the first level, the main purpose is to assign appropriate PRBs and power to the UEs. Furthermore, sufficient activated VNFs are assigned to each slice. Hence, at this level, we would like to obtain the variables P , E , and M .

Despite the simplification of the problem (19), it is still NP-hard and challenging to solve. Therefore, we relax the variable E [8], [24] and reformulating the constraint (19g), to turn them into a jointly-convex problem; Afterward, we solve this problem using a conventional dual Lagrangian method. In the second level, finding the optimal O-RU association, G , is concerned with the fixed parameter of power, PRB allocation, and the number of activated VNFs. We repeat this procedure until the algorithm converges.

A. Sub-Problem 1

Suppose that G is fixed, we want to obtain P , E and M . Here, we first simplify and relax the parameters to convexify the problem. As we mentioned before, by replacing $p_{u(s,l)}^k$ and $p_{u(n,l)}^k$ in (3) with P_s^{\max} , an upper bound $\bar{I}_{r,u(s,i)}^k$ is obtained for $I_{r,u(s,i)}^k$, and also the lower bound $\bar{\rho}_{u(s,i)}^k$ is achieved for $\rho_{u(s,i)}^k$. Moreover, the lower bound $\bar{\mathcal{R}}_{u(s,i)}^k, \forall s, \forall i$ for $\mathcal{R}_{u(s,i)}^k$ is obtained by replacing $I_{r,u(s,i)}^k$ with $\bar{I}_{r,u(s,i)}^k$ in (8) and (5) and make these equations become concave functions.

Suppose $\bar{\rho}_{r,u(s,i)}^k = \frac{P_s^{\max} |\mathbf{h}_{r,u(s,i)}^H \mathbf{w}_{r,u(s,i)}^k|^2}{BN_0}$. we replace $\rho_{r,u(s,i)}^k$ with $\bar{\rho}_{r,u(s,i)}^k$ in (11), to convexify the (8) for the URLLC and mMTC services that have the short packet transmission. So, a lower bound for (8) is given that is a concave function.

$$\bar{\mathcal{R}}_{u(s_j,i)}^r = \sum_{k=1}^{K_{s_j}} B(\log_2(1 + \bar{\rho}_{u(s_j,i)}^k) - \hat{\zeta}_{u(s_j,i)}^k) e_{u(s_j,i)}^k, \quad (21a)$$

$$\bar{\mathcal{R}}_{u(s_j,i)} = \sum_{r=1}^R \bar{\mathcal{R}}_{u(s_j,i)}^r, \quad (21b)$$

$$\hat{\zeta}_{u(s_j,i)}^k = \log_2(e) Q^{-1}(\epsilon) \sqrt{\frac{\hat{\mathcal{C}}_{u(s_j,i)}^k}{N_{u(s_j,i)}^k}}, \quad (21c)$$

$$\hat{\mathcal{C}}_{u(s_j,i)}^k = 1 - \frac{1}{(1 + \hat{\rho}_{u(s_j,i)}^k)^2}. \quad (21d)$$

Without loss of generality, assume that UPF, O-CU and O-DU use the processors with the same processing capability. We notice that it makes the formulation simpler. However, losing this assumption does not change the formulation significantly and the problem can be solved in the same manner. Therefore, we have $\mu_s = \mu_s^u \approx \mu_s^c \approx \mu_s^d$. Moreover, as mentioned before, the mean arrival data rate of the UPF layer for a service s (α_s^U) is equal to the mean arrival data rate of the O-CU-UP layer (α_s^C) and O-DU (α_s^D). So $\alpha_s = \alpha_s^U \approx \alpha_s^C \approx \alpha_s^D$. Again, this assumption only simplifies the notations and losing it does not make the solution inefficient. These assumptions lead to having the same processing power for each layer $\phi_s^u = \phi_s^c = \phi_s^d$. As a result, we have $M_s = M_s^u = M_s^c = M_s^d$. Using the above assumption, we have $T_s^{\text{DU}} = T_s^{\text{CU}} = T_s^{\text{UPF}}$ and we have $T_s^{\text{proc}} = T_s^{\text{RU}} + T_s^{\text{DU}} + T_s^{\text{CU}} + T_s^{\text{UPF}}$. So, $T_s^{\text{proc}} = T_s^{\text{RU}} + 3 \times T_s^{\text{DU}}$.

The problem (19) is mixed-integer nonlinear programming with two integer variables, the PRB assignment, e ,

and the number of VNFs in slice s , M_s , and by relaxing the variables, the problem is also non-convex; therefore, this problem is NP-hard. Solving the problem is not trivial. To solve the problem by inspiring Stackelberg, we reformulate the equation in (19g) to reduce one of the variables (i.e., M_s) that can be solved after obtaining the rate of UEs. We notice that M_s is similar to the followers in Stackelberg Competition, and power and PRB assignment are identical to the leader. So, the new problem has two variables: power and PRB assignment. This new problem is convex by relaxing the binary variable, the PRB assignment, and estimating the lower bounds (21). The objective function and constraints of the problem are convex and can be solved by the Lagrangian function. After obtaining the power of UEs and PRB assignment, we can obtain the achievable rate of each UE so we can find the optimal number of VNFs in each slice (M_s).

In the following, we define a lemma to find the upper and lower bounds for the optimal number of VNFs based on the achievable rates. Afterward, we obtain the formula to attain the optimal number of VNFs.

Lemma 1. *The optimal number of VNFs in each slice s can be achieved by the $M_s = \max\{M_{u(s,i)} | i \in 1, 2, \dots, U_s\} \forall s$. where, $M_{u(s,i)} = \frac{\alpha_s(T_s^{\max} R_{u(s,i)} - T_s^{\max} \lambda_{u(s,i)} - 1)}{(T_s^{\max} \mu_s - 3)(R_{u(s,i)} - \lambda_{u(s,i)}) - \mu_s}$ for each UE i in slice s .*

Proof. In problem (19), the constraint (19g) can be reformulated as

$$T_s^{\max} \geq \frac{1}{R_{u(s,i)} - \lambda_{u(s,i)}} + \frac{3}{\mu_s - \alpha_s/M_s}, \quad (22a)$$

$$M_s \geq \frac{\alpha_s(T_s^{\max} R_{u(s,i)} - T_s^{\max} \lambda_{u(s,i)} - 1)}{(T_s^{\max} \mu_s - 3)(R_{u(s,i)} - \lambda_{u(s,i)}) - \mu_s}. \quad (22b)$$

Also from equations in (19k), (19h) and (19j), we have

$$\alpha_s/\mu_s \leq M_s \leq \min\{M^{\max}, \phi_{\max}/3\phi_s\}. \quad (23)$$

We denote $\mathfrak{M}_s = \min\{M^{\max}, \phi_{\max}/3\phi_s\}$. Thus, if we restrict constraint (19g) to equality, constraint (19g) is still valid. Also, we have the following inequality.

$$\frac{\alpha_s}{\mu_s} \leq \frac{\alpha_s(T_s^{\max} R_{u(s,i)} - T_s^{\max} \lambda_{u(s,i)} - 1)}{(T_s^{\max} \mu_s - 3)(R_{u(s,i)} - \lambda_{u(s,i)}) - \mu_s} \leq \mathfrak{M}_s. \quad (24)$$

In equation (24), $0 \leq \frac{\alpha_s(T_s^{\max} R_{u(s,i)} - T_s^{\max} \lambda_{u(s,i)} - 1)}{(T_s^{\max} \mu_s - 3)(R_{u(s,i)} - \lambda_{u(s,i)}) - \mu_s}$ is established due to the fact that the numerator and the denominator will both have the same sign. In the numerator, according to the (19i), $R_{u(s,i)} - \lambda_{u(s,i)} \geq 0$, and as we know that $\alpha_s \geq 0$, we have $\alpha_s(R_{u(s,i)} - \lambda_{u(s,i)}) \geq 0$. If we assume that the $(R_{u(s,i)} - \lambda_{u(s,i)})T_s^{\max} \geq 1$, the numerator will be positive. $(R_{u(s,i)} - \lambda_{u(s,i)})T_s^{\max} \geq 1$ since the order of T_s^{\max} is about milli second and the difference between achievable rate and packet rate can be more than $1/T_s^{\max}$. Therefore, to ensure that this constraint will be valid, we restrict constraint (19i) to $R_{u(s,i)} \geq \lambda_{u(s,i)} + 1/T_s^{\max}$. So the numerator will be positive. In the denominator, we can say that $(T_s^{\max} \mu_s)(R_{u(s,i)} - \lambda_{u(s,i)}) - \mu_s \geq 0$, since, $\mu_s \geq 0$

and $(R_{u(s,i)} - \lambda_{u(s,i)}) \geq 1/T_s^{\max}$ as mentioned above. The left side of the equation (24), leads to $R_{u(s,i)} \geq \lambda_{u(s,i)}$ that is the constraint (19i). For the right side, by reformulating the equation (24), we have a new constraint $\forall i, \forall s$ given by

$$\mathcal{R}_{u(s,i)} \geq \varpi_{u(s,i)}, \quad (25)$$

$$\varpi_{u(s,i)} = \lambda_{u(s,i)} + \frac{1}{T_s^{\max}} + \frac{3}{T_s^{\max} \mu_s - \alpha_s \frac{T_s^{\max}}{\mathfrak{M}_s} - 3}. \quad (26)$$

In addition, we denote $M_{u(s,i)} = \frac{\alpha_s(T_s^{\max} R_{u(s,i)} - T_s^{\max} \lambda_{u(s,i)} - 1)}{(T_s^{\max} \mu_s - 3)(R_{u(s,i)} - \lambda_{u(s,i)}) - \mu_s}$ for each UE i in slice s . So to obtain, the optimal number of activated VNF in each slice, we need to find the maximum of the $M_{u(s,i)}$ in each slice as $M_s = \max\{M_{u(s,i)} | i \in 1, 2, \dots, U_s\} \forall s$. \square

Despite simplifying the problem in (19), it is still non-convex and hard to solve. Therefore, the conventional approach to solve the problem of the PRB and the power allocation is to relax the variable \mathbf{E} into continuous value $e_{r,u(s,i)}^k \in [0, 1] \forall s, \forall i, \forall r, \forall k$ [8], [24]. Furthermore, the problem can be solved using the Lagrangian function and iterative algorithm.

In order to make (19) as a standard form of a convex optimization problem, it is required to change the variable of equations (13) to $P_r = \sigma_{q_r}^2 \times 2^{C_r}$ so the constraint (19f) is changed to $P_r \leq \sigma_{q_r}^2 \times 2^{C_{\max}}$. The combination of equations (19b) and (19f) leads to the following equation

$$P_r \leq \zeta_r = \min\{P_{\max}, \sigma_{q_r}^2 \times 2^{C_{\max}}\}. \quad (27)$$

Moreover, the combination of equations in (19e), (19i) and (25) leads to the following equation

$$\bar{\mathcal{R}}_{u(s,i)} \geq \eta_{u(s,i)} = \max\{\mathcal{R}_{u(s,i)}^{\min}, \lambda_{u(s,i)} + \frac{1}{T_s^{\max}}, \varpi_{u(s,i)}\}. \quad (28)$$

Assume \mathbf{v} , \mathbf{m} , \mathbf{h} , $\mathbf{\xi}$, $\mathbf{\chi}$, \mathbf{q} and $\mathbf{\kappa}$ are the matrix of Lagrangian multipliers that have non-zero positive elements. The Lagrangian function is written as

$$\begin{aligned} \mathcal{L}(P, E; \mathbf{v}, \mathbf{\chi}, \mathbf{h}, \mathbf{\xi}, \mathbf{\kappa}, \mathbf{m}) = & \sum_{s=1}^S \sum_{i=1}^{U_s} \delta_s \bar{\mathcal{R}}_{u(s,i)} \\ & + \sum_{s=1}^S \sum_{i=1}^{U_s} \mathbf{h}_{u(s,i)} (\bar{\mathcal{R}}_{u(s,i)} - \eta_{u(s,i)}) - \sum_{r=1}^R \mathbf{m}_r (P_r - \zeta_r) \\ & + \sum_{s=1}^S \sum_{i=1}^{U_s} \sum_{k=1}^K \sum_{r=1}^R \mathbf{\kappa}_{r,u(s,i)}^k p_{r,u(s,i)}^k \\ & + \sum_{s=1}^S \sum_{i=1}^{U_s} \sum_{k=1}^K \sum_{r=1}^R \mathbf{q}_{r,u(s,i)}^k (P_s^{\max} - p_{r,u(s,i)}^k) \\ & + \sum_{r=1}^R \sum_{s=1}^S \sum_{i=1}^{U_s} \mathbf{\chi}_{r,u(s,i)} \left(\sum_{k=1}^{K_s} e_{r,u(s,i)}^k - 1 \right) \\ & - \sum_{s=1}^S \sum_{i=1}^{U_s} \sum_{k=1}^K \sum_{r=1}^R \mathbf{v}_{r,u(s,i)}^k (e_{r,u(s,i)}^k - 1) \\ & + \sum_{s=1}^S \sum_{i=1}^{U_s} \sum_{k=1}^K \sum_{r=1}^R \mathbf{\xi}_{r,u(s,i)}^k e_{r,u(s,i)}^k. \end{aligned} \quad (29)$$

Lemma 2. The derivatives of the Lagrangian function (29) with respect to the \mathbf{P} and \mathbf{E} give the KKT conditions to obtain the optimal value of these two variables [8], [24].

Proof. Assume UE i in slice s , associated with O-RU r , is allocated to PRB k (i.e., $e_{r,u(s,i)}^k = 1$). Therefore, we have the following KKT condition

$$\frac{\partial \mathcal{L}}{\partial p_{r,u(s,i)}^k} = (\delta_s + \mathfrak{h}_{u(s,i)}) \mathfrak{B}_{r,u(s,i)}^k + (\mathfrak{s}_{r,u(s,i)}^k - \mathfrak{D}_{r,u(s,i)}^k) = 0, \quad (30)$$

where $\mathfrak{s}_{r,u(s,i)}^k = \kappa_{r,u(s,i)}^k - \mathfrak{q}_{r,u(s,i)}^k$ and other parameters are as follows:

$$\mathfrak{D}_{r,u(s,i)}^k = \mathfrak{m}_r |\mathbf{w}_{r,u(s,i)}^k|^2 g_{u(s,i)}^r e_{r,u(s,i)}^k, \quad (31a)$$

$$\mathfrak{B}_{r,u(s,i)}^k = \frac{B |\mathbf{h}_{r,u(s,i)}^H \mathbf{w}_{r,u(s,i)}^k|^2 g_{u(s,i)}^r e_{r,u(s,i)}^k}{\ln(2)} \mathfrak{S}_{r,u(s,i)}^k, \quad (31b)$$

$$\mathfrak{S}_{r,u(s,i)}^k = \frac{1}{|\mathbf{h}_{r,u(s,i)}^H \mathbf{w}_{r,u(s,i)}^k|^2 \mathfrak{k}_{r,u(s,i)}^k + BN_0 + I_{r,u(s,i)}^k}. \quad (31c)$$

Also, $\mathfrak{k}_{r,u(s,i)}^k = g_{u(s,i)}^r e_{r,u(s,i)}^k p_{r,u(s,i)}^k$. Thus, from equation (30), optimal power is obtained and power is allocated. We denote $\mathfrak{j}_{r,u(s,i)}^k = g_{u(s,i)}^r e_{r,u(s,i)}^k$. The optimal power is as follow.

$$p_{r,u(s,i)}^k = \left[\frac{(\delta_s + \mathfrak{h}_{u(s,i)}) B \mathfrak{j}_{r,u(s,i)}^k}{\ln 2 \times (-\mathfrak{s}_{r,u(s,i)}^k + \mathfrak{D}_{r,u(s,i)}^k)} - \frac{BN_0 + I_{r,u(s,i)}^k}{|\mathbf{h}_{r,u(s,i)}^H \mathbf{w}_{r,u(s,i)}^k|^2 \mathfrak{j}_{r,u(s,i)}^k} \right]^+. \quad (32)$$

Also $[a]^+ = \max(0, a)$. In addition, PRB assignment can be achieved from the derivatives of the Lagrangian function (29) with respect to the \mathbf{E} as follow.

$$\frac{\partial \mathcal{L}}{\partial e_{r,u(s,i)}^k} = \bar{\mathcal{R}}_{r,u(s,i)}^k (\delta_s + \mathfrak{h}_{u(s,i)}) - \mathfrak{m}_r |\mathbf{w}_{r,u(s,i)}^k|^2 p_{r,u(s,i)}^k g_{u(s,i)}^r + (\mathfrak{s}_{r,u(s,i)}^k - \mathfrak{v}_{r,u(s,i)}^k + \chi_{r,u(s,i)}) = 0. \quad (33)$$

So, the optimal \mathbf{E} is obtained using the KKT conditions, which require solving

$$e_{r,u(s,i)}^k \times (\mathfrak{F}_{r,u(s,i)}^k - \mathfrak{v}_{r,u(s,i)}^k - \mathfrak{m}_r |\mathbf{w}_{r,u(s,i)}^k|^2 p_{r,u(s,i)}^k g_{u(s,i)}^r) = 0, \quad (34)$$

where $\mathfrak{F}_{r,u(s,i)}^k = \bar{\mathcal{R}}_{r,u(s,i)}^k (\delta_s + \mathfrak{h}_{u(s,i)}) + (\mathfrak{s}_{r,u(s,i)}^k + \chi_{r,u(s,i)})$. Hence, from equation (33) and (34), PRB assignment is performed as follows

$$e_{r,u(s,i)}^k = \begin{cases} 1 & u(s,i) = \operatorname{argmax}_k \mathfrak{F}_{r,u(s,i)}^k \forall r, k \in K, s \in S, \\ 0 & \text{otherwise,} \end{cases} \quad (35)$$

where $\mathfrak{F}_{r,u(s,i)}^k = (\mathfrak{F}_{r,u(s,i)}^k - \mathfrak{v}_{r,u(s,i)}^k - \mathfrak{m}_r |\mathbf{w}_{r,u(s,i)}^k|^2 p_{r,u(s,i)}^k g_{u(s,i)}^r)$. \square

Thus, the user in slice s that has the most considerable

value of $\mathfrak{F}_{r,u(s,i)}^k$, should be allocated to PRB k . Since just one PRB can be allocated to a UE between those UEs (regardless of the services), that is associated to the same O-RU. The number of UEs are $\mathfrak{N} = \sum_{s=1}^S \sum_{i=1}^{U_s} 1$. Also, assume that the algorithm converges after T_{conv} times. The complexity order of this problem is about $O(T_{conv} \times \mathfrak{N} \times K)$.

B. Sub-Problem 2

After power allocation and PRB assignment, the remaining problem is to assign O-RU to each UE in each service.

Assume \mathbf{P} and \mathbf{E} are fixed, we want to find \mathbf{G} . Next, we introduce a greedy algorithm that assigns an O-RU to each UE.

Greedy Algorithm Assignment (GAA): The problem can be reformulated as follow

$$\max_{\mathbf{G}} \sum_{s=1}^S \sum_{i=1}^{U_s} \sum_{r=1}^R \delta_s g_{u(s,i)}^r \bar{\mathcal{R}}_{u(s,i)}^r \quad (36a)$$

$$\text{subject to} \sum_{s=1}^S \sum_{i=1}^{U_s} g_{u(s,i)}^r \psi_{r,u(s,i)} \leq \mathfrak{t}_r \quad \forall r, \quad (36b)$$

$$\sum_r g_{u(s,i)}^r = 1 \quad \forall s, i, \quad (36c)$$

$$g_{u(s,i)}^r \in \{0, 1\} \quad \forall s, i, \quad (36d)$$

where $\psi_{r,u(s,i)} = \sum_{k=1}^{K_s} |\mathbf{w}_{r,u(s,i)}^k|^2 p_{r,u(s,i)}^k e_{r,u(s,i)}^k$ and $\mathfrak{t}_r = \zeta_r - \sigma_r$ because of the equations (27) and (12). Since we obtained (28) in (VI-A), we can ignore this constraint in (36). The problem (36) is an NP-complete 0-1 multiple knapsack problem. We solve this problem using heuristic method (GAA method 1), which is a greedy algorithm [8], [36]. Firstly, we set all the variables to zero ($g_{u(s,i)}^r = 0, \forall s, i, r$). Then we define the parameter $\mathfrak{B}_{u(s,i)}^{rem}$. This parameter is used as a set of O-RUs that can be assigned to the UE i in slice s , which initially includes all the O-RUs ($\mathfrak{B}_{u(s,i)}^{rem} = \mathcal{R}, \forall s, i$). Also we introduce another parameter $\mathfrak{C}_r = \mathfrak{t}_r, \forall r$ which is the knapsack capacity of each O-RU. Next, we sort all the slices based on their priority. Afterward, based on the sorting of the UEs, we assign the O-RU that provides the highest achievable data rate for each UE on the condition that the value of the desired UE ($\psi_{r,u(s,i)}$) does not exceed the knapsack capacity of each O-RU (\mathfrak{C}_r). If it exceeds the capacity of the desired O-RU, we remove the specific O-RU from the set of O-RUs that can be assigned to that UE ($\mathfrak{B}_{u(s,i)}^{rem} = \mathfrak{B}_{u(s,i)}^{rem} \setminus \{r^*\}$). Then, the O-RU with the highest achievable data rate from the new set of O-RUs $\mathfrak{B}_{u(s,i)}^{rem}$ is selected. The complexity of sorting S slices based on their priority is $O(S \log(S))$. Depict $\mathfrak{N} = \sum_{s=1}^S \sum_{i=1}^{U_s} 1$ as the whole number of UEs in the system. The complexity order of this algorithm is about $O(S \log(S)) + O(R \times \mathfrak{N})$.

C. Iterative Proposed Algorithm

In Sections (VI-A) and (VI-B), the details of solving each sub-problem are depicted. Here, the iterative algorithm for the whole problem is demonstrated. Firstly, we fixed \mathbf{G}

Algorithm 1 Greedy Algorithm for Assignment of O-RU to UEs (GAA)

```

1: Set  $g_{u(s,i)}^r = 0$ ,  $\mathcal{C}_r = \mathbf{t}_r$ , and  $\mathfrak{B}_{u(s,i)}^{rem} = \mathcal{R} \ \forall s, \forall i, \forall r$ .
2: Sort slices according to their  $\delta_s$  in descending order
3: for  $s \leftarrow 1$  to  $S$  do
4:   for  $i \leftarrow 1$  to  $U_s$  do
5:      $RU = 0$ 
6:     for  $r \leftarrow 1$  to  $R$  do
7:       Acquire  $\mathfrak{G}_{u(s,i)}^r = \bar{\mathcal{R}}_{u(s,i)}^r$ 
8:     end for
9:     Obtain  $r^* = \operatorname{argmax}_{r \in \mathfrak{B}_{u(s,i)}^{rem}} \mathfrak{G}_{u(s,i)}^r$ 
10:    while  $RU == 0$  do
11:      if  $\mathcal{C}_{r^*} \geq \psi_{r^*,u(s,i)}$  then
12:        Set  $g_{u(s,i)}^{r^*} = 1$ 
13:        Set  $\mathcal{C}_{r^*} = \mathcal{C}_{r^*} - \psi_{r^*,u(s,i)}$ 
14:        Set  $RU = 1$ 
15:      else:  $\mathfrak{B}_{u(s,i)}^{rem} = \mathfrak{B}_{u(s,i)}^{rem} \setminus \{r^*\}$ 
16:      end if
17:    end while
18:  end for
19: end for

```

Algorithm 2 Iterative algorithm for the baseband resource allocation and VNF activation (IABV)

```

1: Set the maximum num. of iter.  $I_{\max}$ , convergence condition  $\epsilon > 0$ 
2: Assign Users to O-RU randomly (Initialize  $\mathbf{G}$ )
3: for  $i \leftarrow 1$  to  $I_{\max}$  do
4:   Acquire  $\mathbf{P}^{(i)}$ ,  $\mathbf{E}^{(i)}$  and  $\mathbf{M}^{(i)}$  using Lagrangian function and sub-gradient method based on (VI-A)
5:   Update  $\mathbf{G}^{(i)}$  based on algorithm GAAOU (1) in (VI-B)
6:   if the algorithm converged with the tolerance of  $\epsilon$  then
7:     Break
8:   else: Continue the algorithm
9:   end if
10: end for

```

to achieve \mathbf{P} and \mathbf{E} , using the Lagrangian method and the KKT conditions. Afterward, \mathbf{G} is updated using the GAA algorithm. This process is repeated until it converges. The whole algorithm (IABV method) is depicted as follows (Algorithm 2).

1) *Complexity Order:* The number of UEs are $\mathfrak{N} = \sum_{s=1}^S U_s$. Also, assume that the algorithm converges after T_{conv} times. As we mentioned before, the complexity order of the first sub-problem is about $O(T_{conv} \times \mathfrak{N} \times K)$ and the complexity order of the second sub-problem is about $O(S \log(S)) + O(R \times \mathfrak{N})$. So the complexity of the main problem (19) is $O(T_{conv} \times \mathfrak{N} \times K \times (S \log(S) + R\mathfrak{N}))$.

2) *Convergence Analysis:* Due to limited resources in power, the number of PRBs, and the number of activated VNFs and restrictions based on this limitation on power, energy, fronthaul capacity, etc., the objective function that is the summation of the achievable rates of UEs cannot exceed its optimal value and become infinite. Therefore, we can guarantee the convergence of the iterative algorithm if the objective function is the strictly ascending function concerning the number of iterations. Consequently, it converges to the optimum value. Consider the aggregate throughput as $\mathcal{T}(\mathbf{P}, \mathbf{E}, \mathbf{G}) = \sum_{s=1}^S \sum_{i=1}^{U_s} \delta_s \bar{\mathcal{R}}_{u(s,i)}$. In the

first step of the iteration i of the algorithm 2 (IABV), we have $\mathcal{T}(\mathbf{P}^i, \mathbf{E}^i, \mathbf{G}^{i-1})$. In this step, optimal power and PRB allocation are obtained for the fixed O-RU association, so we have $\mathcal{T}(\mathbf{P}^i, \mathbf{E}^i, \mathbf{G}^{i-1}) \geq \mathcal{T}(\mathbf{P}^{i-1}, \mathbf{E}^{i-1}, \mathbf{G}^{i-1})$. In the second step of the iteration i , the optimal O-RU association is achieved to maximize the aggregate throughput. So we have this inequality $\mathcal{T}(\mathbf{P}^i, \mathbf{E}^i, \mathbf{G}^i) \geq \mathcal{T}(\mathbf{P}^i, \mathbf{E}^i, \mathbf{G}^{i-1})$. As a result, we have $\mathcal{T}(\mathbf{P}^i, \mathbf{E}^i, \mathbf{G}^i) \geq \mathcal{T}(\mathbf{P}^{i-1}, \mathbf{E}^{i-1}, \mathbf{G}^{i-1})$. Hence, in each step of the iteration, the aggregate throughput increased. Note that $\mathcal{T}^*(\mathbf{P}^*, \mathbf{E}^*, \mathbf{G}^*)$ is the achieved aggregate throughput for all the feasible resource allocation solutions of $\{\mathbf{P}, \mathbf{E}, \mathbf{G}\}$. So, $\mathcal{T}^*(\mathbf{P}^*, \mathbf{E}^*, \mathbf{G}^*) \geq \mathcal{T}(\mathbf{P}^i, \mathbf{E}^i, \mathbf{G}^i)$ and thus in each iteration, the aggregate throughput can not be larger than the optimal solution. So the the aggregate throughput is ascending function concerning the number of iterations and it will converge to the sub-optimal solution. In addition, if we assume that the interference is set to be zero $I_{r,u(s,i)}^k = 0$, and we suppose that each UE has the maximum power $p_{r,u(s,i)}^k = P_s^{max}$, and we consider that all PRB is assigned to all UE $e_{r,u(s,i)}^k = 1 \ \forall s, \forall i$ and each UE is assigned to the nearest O-RU with the best channel quality. So, the solution of this allocation, is the upper bound for the aggregate throughput. Thus, we can guarantee the convergence of our iterative algorithm since the objective function \mathcal{T} is the ascending function concerning the number of iterations and it has the upper bound.

Algorithm 3 Fast Algorithm (FA) to Check the Convergence

```

1: Set count = 0
2: Set  $p_{r,u(s,i)}^k = 0$ ,  $e_{r,u(s,i)}^k = 0$  and  $g_{u(s,i)}^r = 0 \ \forall r, k, s, i$ 
3: for  $s \leftarrow 1$  to  $S$  do
4:   for  $i \leftarrow 1$  to  $U_s$  do
5:     count = count + 1
6:      $r^* = \operatorname{argmin}_r d_{r,u(s,i)} \ \forall r$ 
7:      $g_{u(s,i)}^{r^*} = 1$ 
8:     temp = mod(count, K)
9:     if temp=0 then
10:       $e_{r^*,u(s,i)}^K = 1$ 
11:      Set  $p_{r^*,u(s,i)}^K = \min\{P_s^{max}, P_r^{max}/\mathfrak{N}\}$ 
12:    else:  $e_{r^*,u(s,i)}^{temp} = 1 \ \forall r$ 
13:      Set  $p_{r^*,u(s,i)}^{temp} = \min\{P_s^{max}, P_r^{max}/\mathfrak{N}\}$ 
14:    end if
15:   end for
16: end for

```

VII. NUMERICAL RESULTS AND THE FEASIBLE REGION

In this section, firstly, we describe the initial points and the comparison algorithms. Then, we talk about the feasible region of our system model. Afterward, we illustrate the numerical results.

A. The initial Points and The Comparison Algorithms

In this part, numerical results for the main problem are depicted to evaluate the performance of the algorithms using

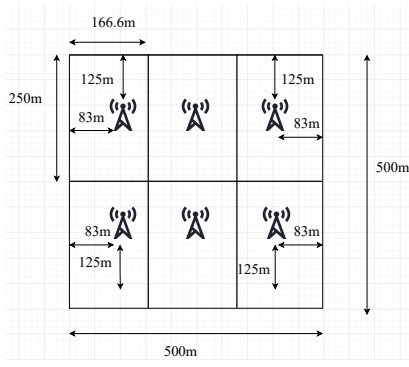


Fig. 2: O-RU placement in a cell

the Monte-Carlo method. We consider three network slices for eMBB, URLLC, and mMTC services. Assume we have six 4-antenna O-RU (MISO) located in a place with a diameter of 500 meters as shown in Fig. 2. In addition, we consider the users placed randomly in this area.

Here, the channel vector from the O-RU r to the UE i in service s is set as $\mathbf{h}_{r,u(s,i)}^k = d_{r,u(s,i)}^{-\mathcal{L}} \Omega_{r,u(s,i)}^k$, where $d_{r,u(s,i)}^{-\mathcal{L}}$ is the distance between the O-RU r and UE i in service s and $\mathcal{L} = 3.8$ is the path-loss exponent. Also, $\Omega_{r,u(s,i)}^k$ is the random variable that is generated by the Rayleigh distribution and it is the Rayleigh fading channel between the UE and O-RU. We consider 25 PRBs in the network. The packet size for mMTC is equal to 20 bytes, and for URLLC is equal to 32 bytes [37]. The maximum number of VNF for each slice is 25 and the mean arrival data rate for the eMBB service is $\lambda = 3\text{Mbps}$ and for the mMTC service and the URLLC service is $\lambda = 0.2\text{Mbps}$. Also, the quantization noise is assumed to be 10^{-13} . Moreover, we set $\eta_{u(s,i)} = \eta_{u(s,i)}/200$, $\mathbf{m}_r = \zeta_r/10$ and $\mathbf{q}_{r,u(s,i)}^k = P_s^{\max}/100$. The other parameters of these simulations are depicted in Table II [37]–[40].

TABLE II: Simulation Parameters

Parameter	Value
noise power	-174 dBm
bandwidth	180 KHz
maximum transmit power of each O-RU	40 dBm
maximum delay for eMBB	4 msec
maximum delay for URLLC	1 msec
maximum delay for mMTC	5 msec
maximum fronthaul capacity	46 bps
minimum data rate for eMBB	20 bps
minimum data rate for URLLC and mMTC	2 bps
maximum received power for mMTC	20 dBm
maximum received power for eMBB and URLLC	33 dBm

Finding a feasible initial value is almost tricky. We use a fast method discussed in VII-B to overcome this challenge. Two different methods are used to compare the performance of the proposed method (IABV) and show the optimality of our approach. The first one is a baseline scheme, which uses random PRB allocation. Therefore, the allocation of PRB to each UE is random when we have low interference, but in figures with high interference, we randomly assign just one RB to each UE. Also, the association of O-RU is carried out based on distance. It means that each UE is assigned

to the nearest O-RU. The optimal power is obtained using the CVX of Matlab, which uses the successive convex approximation (SCA) method since the problem is convex. After achieving power and other parameters, the achievable rate will be obtained, and the optimal number of VNF is achieved from Lemma (1).

For the second one, we use the idea of the fixed BBU capacity and dynamic resource allocation (FBDR) algorithm proposed in [8] and named it the dynamic resource allocation scheme (DR scheme). We have services with different QoS in this work, similar to tenants with different QoS introduced in [8]. Therefore, we can use the DR scheme similar to the FBDR method adapted to our conditions for comparison. Instead of BBU in C-RAN, we have O-DU and O-CU in O-RAN. Since we do not talk about O-DU and O-CU capacity, we use the dynamic resource allocation scheme (DR scheme) algorithm and do not consider BBU capacity. In the DR scheme, PRB and power are dynamically allocated. The number of VNFs is obtained from the simulation. The UEs are associated with the O-RU based on the quality of their channels and the channel distance instead of using the greedy algorithm 1 (GAA algorithm) for O-RU assignment. The figures in [8] show that dynamic BBU capacity and dynamic resource allocation (DBDR) perform better than FBDR for the same priority area. The numerical results section also indicates that our proposed algorithm performs better than the DR scheme.

B. Feasible Region

Applying the correct initial point to make the system feasible and converge is a significant step in our work. To solve this problem, we investigated the non-converging and converging simulation for models with fixed initial parameters and random channel gains of UEs. We experimentally found that in non-converged simulations, there are UEs at the edge of the boundaries or far away from the O-RU and have a weak channel gain. One solution is to eliminate UEs who undermine system convergence. For a large number of UEs with a fixed number of PRBs, the probability of having an infeasible solution increases due to a large number of UE interference. Another solution is to remove the simulations in the Monte-Carlo that do not converge using the fast algorithm (FA) to check the convergence before the proposed algorithm (IABV). Therefore, if more than half of the iterations have a feasible solution for the initial condition, the simulation can be displayed as a feasible model. If the conditions in (28), (27), (19d) and (19c) are met in the fast algorithm (FA), the given algorithm will converge. Assume, the number of UEs is $\mathfrak{N} = \sum_{s=1}^S \sum_{i=1}^{U_s} 1$, the number of PRBs is K , and the distance between the r^{th} O-RU to the UE i in slice s is $d_{r,u(s,i)}$. The FA algorithm is represented in Algorithm 3. The complexity order of this algorithm is $O(R \times \mathfrak{N})$ which is remarkably lower than the complexity order of the IABV method. In the FA algorithm, the O-RU association is based on the distance of the UE to the O-RU. Each UE is associated with the nearest O-RU. Also, the power of each

UE is set to be the minimum of the maximum power of each UE and the maximum power of each O-RU divided by the total number of UEs ($\min\{P_s^{\max}, P_r^{\max}/\mathcal{N}\}$). Moreover, the allocation of PRBs to UEs is based on dividing the number of UEs by the total number of PRBs.

C. Numerical Results

In Fig. 3, the aggregate throughput is demonstrated versus the different number of UEs in each service for these three methods. Suppose we have one service instance for each type of service, so we have three various services in this figure. Also, we have between 6 to 48 UEs in the system. Here, we did not consider the priority. The figure presented that the proposed method, IABV, is 18.6% higher throughput than the baseline scheme. As the number of UEs increases in each service, the aggregated throughput initially increases. Still, due to the interference and the power constraint, it will be saturated from 12 UEs in each service.

Fig. 4 depicts the number of activated VNFs for five different mean service times of one URLLC service vs. the mean arrival time for 12 UEs. This figure presents that as the mean arrival rate increases, the number of activated VNF increases. Moreover, the number of activated VNFs decreases when the mean service rate increases.

In Fig. 5, the aggregate throughput is depicted vs. the maximum power of UE for three different instances of eMBB service using proposed method (IABV), DR scheme and the baseline scheme. Here, we suppose that we have 12 UEs in each service. We assume that these three services require 5bits/sec/Hz, 10bits/sec/Hz, and 15bits/sec/Hz. As you can see in the figure, increasing the maximum power increases the aggregate throughput. Moreover, the proposed method (IABV), gives higher aggregate rates in compared to the DR scheme and the baseline scheme.

Fig. 6 illustrates the mean total delay of a UE in a URLLC service regarding the mean arrival rate of the UE and the number of UEs in the service for the proposed method (IABV). It is shown that the delay is an ascending function of the mean arrival rate (when the mean service time is fixed) and the number of UEs in the service. Moreover, we can see that the mean delay of a URLLC service does not reach the maximum threshold of the delay.

Fig. 7 is the same as Fig. 6 that presented the mean total delay of a UE in a URLLC service regarding the mean arrival rate of the UE for 20 UEs using three different methods. As you can see, the proposed method (IABV) outperforms the other scenarios.

Fig. 8, represents the aggregate throughput concerning the number of UEs in each service and the maximum power for three different mMTC service instances. Assume each UE in each mMTC service instance requires 0.1 bits/sec/Hz data rate and is not sensitive to the end-to-end delay. There is no restriction on fronthaul link capacity and the number of VNFs. The figure depicts that by increasing the number of UEs in each instance of the service, or by increasing the maximum power of each UE in each instance of mMTC service, the aggregate throughput increases.

Assume we have two types of eMBB service instances. In Fig. 9, the aggregate throughput (by considering the priority factor δ_s) is depicted for two eMBB service instances. Here we consider 4 UEs in each service. The Fig. 9 presented that by increasing the priority factor for one service instance, more resources are allocated to this service instance, and the aggregate throughput of this service is increased and vice versa. Also, we can realize from this figure that the aggregate throughput has the most significant value at the same priority.

In Fig. 10, the aggregate throughput is shown according to the number of iteration (outer loop) of the proposed algorithm (IABV) for different numbers of UEs for one service. In this figure, the convergence of the IABV method is illustrated. The minimum data rate for each UE is assumed to be 2 Mbps. After four iterations, IABV converges to the fixed value.

In Fig. 11, the mean total delay of URLLC service is indicated according to the number of iterations of the proposed algorithm (IABV) for different numbers of UEs for one URLLC service. This figure shows that the algorithm converged to the fixed value after four iterations.

In Fig. 12, the aggregate throughput is shown according to the number of UEs for two different methods, namely the proposed algorithm (IABV) and the optimal method for URLLC service for the low interference. The minimum data rate is 5bits/sec/Hz for each UE and the maximum delay is 0.1ms. Also the mean arrival rate is set to be 0.2Mbps and the mean service rate is 0.5Mbps. The optimal approach is obtained from the two-step joint exhaustive search and using CVX. In each iteration in the first step, the PRB allocation and O-RU association are obtained from brute force, and in the second step, we use CVX to get optimal power. Our solution is close to the optimal value in a small number of UEs.

In Fig. 13, the aggregate throughput is depicted vs. the maximum interference for two different maximum power thresholds of O-RU. Here we assume that with the increase of every ten dBm of interference power, it is assumed that ten users have been added to the system. In -105 dBm, we have 5 UEs, and at the end, we have 55 UEs in the system. Since the amount of interference in the system is entered as a fixed value, the allocation of PRBs is not considered. The higher maximum power threshold leads to a greater aggregate throughput. The aggregate throughput first increases with the number of UEs and at the same time the amount of the maximum interference, then it becomes almost fixed and finally decreases so much. When the aggregate throughput decreases, the maximum interference is so high that it takes the system out of feasibility.

In Fig. 14, the aggregate throughput is shown versus the number of UEs for an eMBB service with low interference for the IABV and FA methods in the feasible region. The minimum data rate for each UE is 1Mb/s/Hz. The maximum power for each O-RU is 34dBm, and the maximum power for each UE is 30dBm. We assume that the system is not sensitive to fronthaul capacity and end-to-end delay and has enough VNF resources. By increasing the number of UEs,

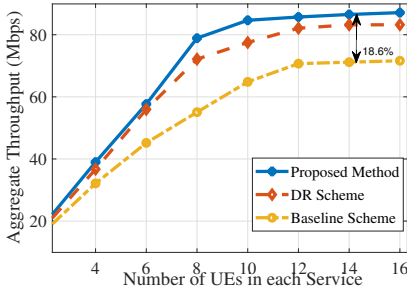


Fig. 3 Aggr. throughput vs. number of UEs in each service

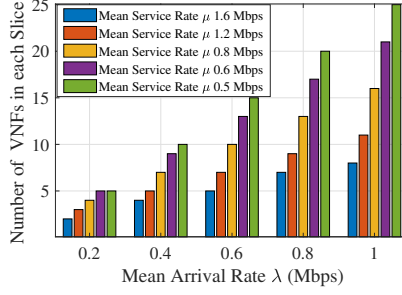


Fig. 4 Number of VNFs in each service vs. arrival rate

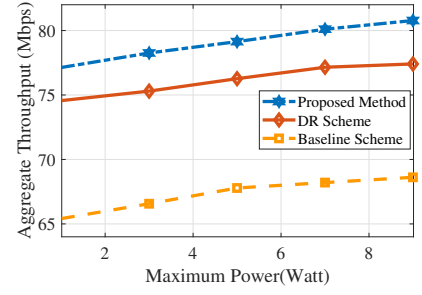


Fig. 5 Aggr. throughput for eMBB vs. maximum power

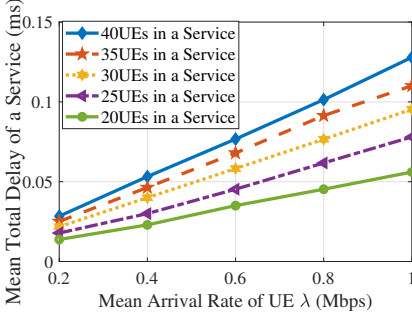


Fig. 6 Total Delay of a URLLC vs. the arrival rate of a UE

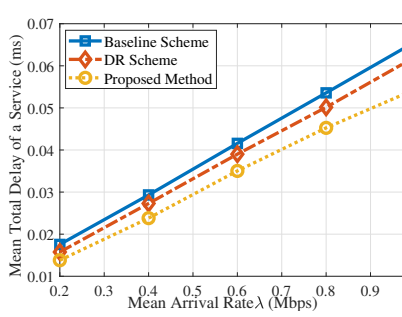


Fig. 7 Total Delay of a URLLC vs. arrival rate of a UE

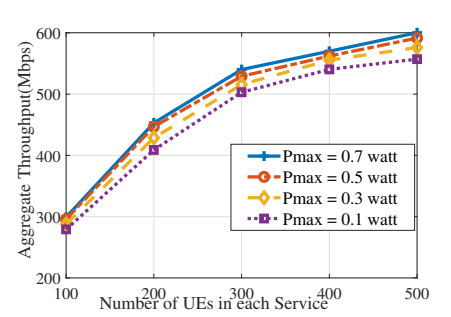


Fig. 8 Aggr. throughput vs. the number of mMTC UEs

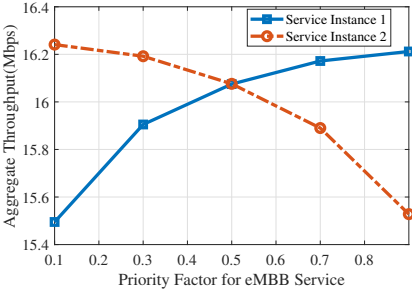


Fig. 9 Throughput of eMBB services vs. priority of the first one

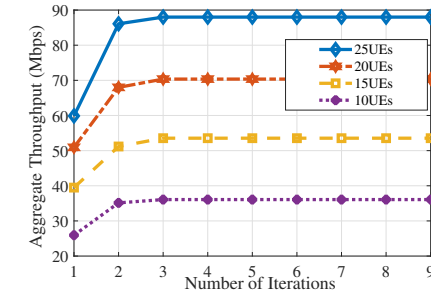


Fig. 10 Aggregate throughput vs. number of iterations

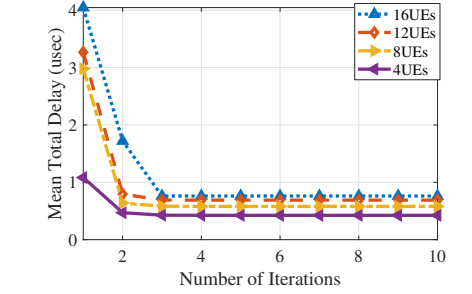


Fig. 11 Mean Delay time vs. number of iterations

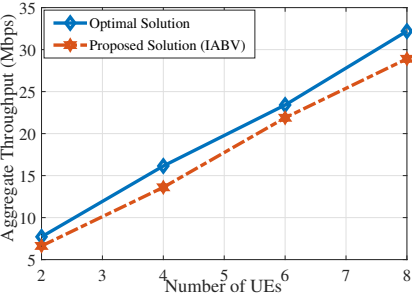


Fig. 12 Aggregate throughput vs. number of UEs

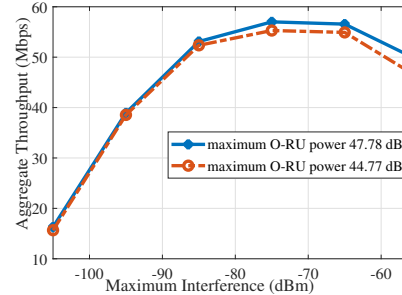


Fig. 13 Aggregate throughput vs. maximum interference

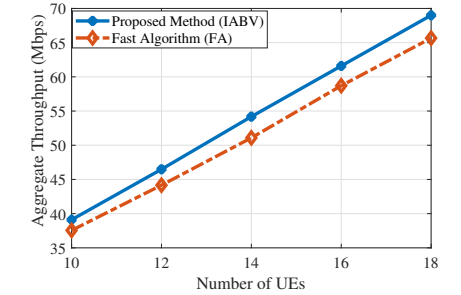


Fig. 14 Aggregate throughput vs. number of UEs

TABLE III: Execution Time vs. Number of UEs

Number of UEs	5	10	15	20	25
Execution Time (usec)	12.156	19.156	29.140	44.573	67.912

the aggregate throughput raises. And we can see that the IABV method is better than the FA method.

Table III shows the execution time versus the number of UEs for one service. We run our simulation on the system with configures (RAM = 8 GB, CPU = Core i5, SSD Hard Disk). As the number of UEs in the system increases, the execution time increases polynomially.

VIII. CONCLUSION

In this paper, we modeled the downlink of the O-RAN system using isolated network slicing for different 5G services, i.e., eMBB, mMTC, and URLLC. We aim to obtain the number of activated VNFs in each service, RU

association, power, and PRB allocation to maximize the aggregate throughput. The limited fronthaul capacity and the mean end-to-end delay for each service are considered. The problem is mixed-integer non-linear programming that is solved by a two-step iterative algorithm. In the first step, we reformulated the problem to achieve the number of activated VNFs as a function of data rate. Then, we obtained PRB association and power allocation using the Lagrangian method. In the second step, the O-RU association is carried out. The performance of our proposed method (i.e., IABV) is compared with the baseline scheme and DR scheme in [8]. In addition, the feasible region is discussed, and the FA algorithm is introduced to check the feasibility of the initial values. Also, we assume distinct scenarios for each service, i.e., eMBB, mMTC, and URLLC, based on their requirement QoS. Simulation results show that the proposed method (i.e., IABV) achieves 18.6% higher data rate than the baseline scheme. Moreover, simulation results illustrate more minor delays for the proposed method (IABV) than DR scheme and the baseline scheme.

REFERENCES

- [1] X. Shen *et al.*, "AI-assisted network-slicing based next-generation wireless networks," *IEEE Open Journal of Vehicular Technology*, vol. 1, pp. 45–66, Jan. 2020.
- [2] M. Setayesh *et al.*, "Joint PRB and power allocation for slicing eMBB and URLLC services in 5G C-RAN," in *IEEE Global Communications Conference (GLOBECOM)*, Taipei, Taiwan, Dec. 2020, pp. 1–6.
- [3] P. Popovski *et al.*, "5G wireless network slicing for eMBB, URLLC, and mMTC: A communication-theoretic view," *IEEE Access*, vol. 6, pp. 55 765–55 779, Sep. 2018.
- [4] A. Dogra *et al.*, "A survey on beyond 5G network with the advent of 6G: Architecture and emerging technologies," *IEEE Access*, vol. 9, pp. 67 512–67 547, Oct. 2020.
- [5] M. Alsenwi *et al.*, "Intelligent resource slicing for eMBB and URLLC coexistence in 5G and beyond: A deep reinforcement learning based approach," *IEEE Transactions on Wireless Communications*, vol. 20, no. 7, pp. 4585 – 4600, Feb. 2021.
- [6] S. Riolo *et al.*, "Modeling and analysis of tagged preamble transmissions in random access procedure for mMTC scenarios," *IEEE Transactions on Wireless Communications*, vol. 20, no. 7, pp. 4296–4312, 2021.
- [7] R. Mijumbi *et al.*, "Network function virtualization: State-of-the-art and research challenges," *IEEE Communications surveys & tutorials*, vol. 18, no. 1, pp. 236–262, Sep. 2015.
- [8] Y. Lee *et al.*, "Dynamic network slicing for multitenant heterogeneous cloud radio access networks," *IEEE Transactions on Wireless Communications*, vol. 17, no. 4, pp. 2146–2161, Apr. 2018.
- [9] J. Tang *et al.*, "System cost minimization in cloud RAN with limited fronthaul capacity," *IEEE Transactions on Wireless Communications*, vol. 16, no. 5, pp. 3371–3384, May 2017.
- [10] "O-RAN architecture description," O-RAN Alliance, Tech. Rep., 2020.
- [11] L. Gavrilovska *et al.*, "From Cloud RAN to Open RAN," *Wirel. Pers. Commun.*, vol. 113, no. 3, pp. 1523–1539, Mar. 2020.
- [12] N. Kazemifard and V. Shah-Mansouri, "Minimum delay function placement and resource allocation for Open RAN (O-RAN) 5G networks," *Computer Networks*, vol. 188, p. 107809, Apr. 2021.
- [13] L. Feng *et al.*, "Dynamic resource allocation with RAN slicing and scheduling for uRLLC and eMBB hybrid services," *IEEE Access*, vol. 8, pp. 34 538–34 551, Feb. 2020.
- [14] Y. L. Lee *et al.*, "A new network slicing framework for multi-tenant heterogeneous cloud radio access networks," in *International Conference on Advances in Electrical, Electronic and Systems Engineering (ICAEEES)*. Putrajaya, Malaysia: IEEE, Nov. 2016, pp. 414–420.
- [15] H. Xiang *et al.*, "A realization of fog-RAN slicing via deep reinforcement learning," *IEEE Transactions on Wireless Communications*, vol. 19, no. 4, pp. 2515–2527, Jan. 2020.
- [16] S. E. Elayoubi *et al.*, "5G RAN slicing for verticals: Enablers and challenges," *IEEE Communications Magazine*, vol. 57, no. 1, pp. 28–34, Jan. 2019.
- [17] S. D'Orto *et al.*, "Toward operator-to-waveform 5G radio access network slicing," *IEEE Communications Magazine*, vol. 58, no. 4, pp. 18–23, Apr. 2020.
- [18] P. Yang *et al.*, "How should I orchestrate resources of my slices for bursty URLLC service provision?" *IEEE Transactions on Communications*, vol. 69, no. 2, pp. 1134–1146, Nov. 2020.
- [19] F. Saggese *et al.*, "Power minimization of downlink spectrum slicing for eMBB and URLLC users," *arXiv preprint arXiv:2106.08847*, Jun. 2021.
- [20] K. Guo *et al.*, "Exploiting hybrid clustering and computation provisioning for green C-RAN," *IEEE Journal on Selected Areas in Communications*, vol. 34, no. 12, pp. 4063–4076, Nov. 2016.
- [21] P. Luong *et al.*, "Joint virtual computing and radio resource allocation in limited fronthaul green C-RANs," *IEEE Transactions on Wireless Communications*, vol. 17, no. 4, pp. 2602–2617, Feb. 2018.
- [22] —, "A novel energy-efficient resource allocation approach in limited fronthaul virtualized C-RANs," in *IEEE 87th Vehicular Technology Conference (VTC Spring)*, Jun. 2018, pp. 1–6.
- [23] S. Ali *et al.*, "Energy-efficient resource allocation and RRH association in multi-tier 5G H-CRANs," *Transactions on Emerging Telecommunications Technologies*, vol. 30, no. 1, p. e3521, Jan. 2019.
- [24] —, "Joint RRH-association, sub-channel assignment and power allocation in multi-tier 5G C-RANs," *IEEE Access*, vol. 6, pp. 34 393–34 402, Jun. 2018.
- [25] N. Amani *et al.*, "Power-Efficient resource allocation in massive MIMO aided Cloud RANs," *arXiv preprint arXiv:1908.07568*, Aug. 2019.
- [26] B. Han, L. Liu *et al.*, "Research on resource migration based on novel RRH-BBU mapping in cloud radio access network for HSR scenarios," *IEEE Access*, vol. 7, pp. 108 542–108 550, Aug. 2019.
- [27] S. Niknam *et al.*, "Intelligent O-RAN for beyond 5G and 6G wireless networks," *arXiv preprint arXiv:2005.08374*, May 2020.
- [28] C. B. Both *et al.*, "System intelligence for UAV-Based mission critical with challenging 5G/b5G connectivity," *arXiv preprint arXiv:2102.02318*, Feb. 2021.
- [29] O. W. G. 2, "AI/ML workflow description and requirements," O-RAN Alliance, Tech. Rep., Mar. 2020.
- [30] B. S. Lin, "Toward an AI-Enabled O-RAN-based and SDN/NFV-driven 5G and IoT network era," *Network and Communication Technologies*, vol. 6, no. 1, pp. 6–15, Jun. 2021.
- [31] J. Mei *et al.*, "Intelligent radio access network slicing for service provisioning in 6G: A hierarchical deep reinforcement learning approach," *IEEE Transactions on Communications*, vol. 69, no. 9, pp. 6063–6078, 2021.
- [32] S. Huang *et al.*, "User selection for multiuser MIMO downlink with zero-forcing beamforming," *IEEE Transactions on Vehicular Technology*, vol. 62, no. 7, pp. 3084–3097, 2013.
- [33] J. Cavazos. (2020) 5G testing: What is O-RAN? – part 2. [Online]. Available: <https://blogs.keysight.com/blogs/inds.entry.html>
- [34] D. Marabissi and R. Fantacci, "Highly flexible RAN slicing approach to manage isolation, priority, efficiency," *IEEE Access*, vol. 7, pp. 97 130–97 142, 2019.
- [35] ETSI-TS-128-530-V15.0.0, "5G management and orchestration; concepts, use cases and requirements (3GPP TS 28.530 version 15.0.0 release 15)," 2018-10.
- [36] Y. Akçay *et al.*, "Greedy algorithm for the general multidimensional knapsack problem," *Annals of Operations Research*, vol. 150, no. 1, pp. 17–29, Dec. 2007.
- [37] ETSI-TR-138-913-V14.3.0, "5G; study on scenarios and requirements for next generation access technologies (3GPP TR 38.913 version 14.3.0 release 14)," 2017-10.
- [38] 3GPP-TS-36.104-V13.3.0, "Evolved universal terrestrial radio access (E-UTRA); base station (BS) radio transmission and reception (release= 13)," 2016-03.
- [39] 3GPP-TR-36.931-V13.0.0, "Evolved universal terrestrial radio access (E-UTRA); radio frequency (RF) requirements for LTE pico Node B (release 13)," 2016-01.
- [40] E. Mohyeldin, "Minimum technical performance requirements for IMT radio interface(s)," 2020.