

Resource Allocation in an Open RAN System using Network Slicing

Mojdeh Karbalaee Motalleb

School of ECE, College of Engineering, University of Tehran, Iran

Email: {mojdeh.karbalaee}@ut.ac.ir,

Abstract—Taking advantage of both virtual RAN (v-RAN) and Cloud RAN (C-RAN), Open RAN (O-RAN) is introduced as the next generation of RAN systems which leads to increase flexibility, Openness, and reduce operational costs and allow them to add new capabilities to the network more quickly. O-RAN separate RAN into three different units, namely Radio Unit (O-RU), Distributed Unit (O-DU), and Central Unit (O-CU). In this paper, we study the problem of baseband resource allocation and virtual network function (VNF) activation in O-RAN architecture based on their service priority for a different types of 5G services includes enhanced mobile broadband (eMBB), ultra-reliable low latency communications (URLLC) and massive Machine Type Communications (mMTC) services. According to the concept of network slicing, the isolation of different types of services in O-DU, O-CU, and user plane function (UPF) is performed. Limited fronthaul capacity and the restriction of end-to-end delay are considered in the problem. The optimization of baseband resources includes O-RU assignment, physical resource block (PRB), and power allocation. The main problem is mixed-integer non-linear programming that is tremendously difficult. To solve the challenging problem, we broke it down into two different steps that the iterative algorithm solves. In the first step, we reformulated and simplified the problem to find the power allocation, PRB assignment, and the number of activated VNFs. In the second step, the O-RU association is carried out. The proposed method is confirmed by the simulation results in a way that the simulations illustrate a higher achievable data rate than a baseline scheme that only optimizes one of the baseband resources.

Index Terms—Open Radio Access Network (O-RAN), Virtual Network Function (VNF)

I. INTRODUCTION

One of the goals of the fifth generation of wireless system is to achieve the desired QoS (such as rate, delay, power, ...) for different type of services. Network slicing is the best solution for this aim. A network slice is an end-to-end logical network which offers services with special needs. Multiple isolated network slices run, manage, and work independently on the same infrastructure. There are several implementations of network slicing, including network core slicing, RAN slicing, and slicing of both sections. Different type of services includes enhanced mobile broadband (eMBB) and ultra-reliable low latency communications (URLLC) and massive Machine Type Communications (mMTC) services are introduced in 5th generation of mobile network. Each type of service requires special slice of network based on its QoS [1]–[3].

Recently, RAN virtualization attracts significant attention from both industry and academia since it has remarkable benefits which leads to increase the flexibility and reduce operator costs such as CAPEX and OPEX and also allow

them to add new capabilities to the network more quickly. In addition to RAN virtualization, openness and RAN intelligence are two other fundamental points that encourage Open Radio Access Network (O-RAN) Alliance to establish O-RAN as a next generation of RAN systems. The idea of O-RAN comes from the integration of virtual RAN (vRAN) and cloud RAN (CRAN) and it takes the advantage of both. CRAN, divides RAN into two parts radio remote head (RRH) and base band unit (BBU). More than one distributed RRHs can be connected to a centralized BBU which is named BBU-pool [6]. Unlike previous generation of RAN that divide RAN into two parts, O-RAN separate RAN into three different units, namely Radio Unit (O-RU), Distributed Unit (O-DU) and Central Unit (O-CU). O-RU is a logical node contains RF and lower PHY. Moreover, the O-DU expresses another logical node that includes higher PHY, MAC and RLC. In addition, the O-CU depicts the logical node contains two parts, which are O-CU user plane (O-CU-UP) and O-CU control plane (O-CU-CP). O-CU-UP hosts PDCP-UP and SDAP and O-CU-CP hosts PDCP-CP and RRC. O-DU and O-CU are connected to each other via an open and well-defined interface F_1 . Moreover, O-DU is connected to radio unit (O-RU) with an open fronthaul interface. The architecture of O-RAN contains other principal logical nodes called Orchestration and Automation, RAN Intelligent Controller (RIC)- Near Real Time and O-Cloud. One of the necessities of new generation of wireless networks is its intelligence. Based on the requirement of smart wireless network, O-RAN offers machine learning techniques. The two logical nodes RIC-Non Real Time (which is placed in Orchestration and Automation node) and RIC- Near Real Time implement the algorithms for network intelligence [7]–[13].

To improve system performance in the fifth generation of telecommunications, the separation of network software and hardware elements has been done and introduced as network function virtualization (NFV), and virtual network functions (VNF) are system function blocks. The key idea of the implementation of NFV is to decouple software from physical hardware, dynamic scaling and the deployment of flexible network function. A usual NFV offer is to execute VNFs on a virtual machines or containers in a cloud system [4], [5]. As a result, some O-RAN components that includes user plane function (UPF), O-CU, O-DU and RAN Intelligent Controller (RIC)-near real time, are virtualized and implemented as a VNF that can be run on virtual machines (VMs) or containers.

A. Related Works

Network slicing is increasingly receiving research attention. Many researcher, studied the problem of resource allocation in network slicing for multitenant cellular networks [14]–[16]. In [15], dynamic network slicing in multitenant heterogeneous CRAN (H-CRAN) is considered. The process of allocating network resources to users is discussed. The network slicing scheme includes a higher level, which manages user acceptance control, user communication which includes radio unit association (RRH association to maximize user rates and allocate base band resource capacity), the allocation of BBU capacity and a lower level, which is the allocation of power and physical resource blocks (PRB) among users. In article [17], network slicing in the radio section is considered for fog or F-RAN structure, in which two network slices are set for hotspots and vehicle scenarios with related infrastructure. In [18], [19] the implementation of RAN level slicing is discussed in mobile network operator (MNO). Also the problem of resource allocation is considered. Moreover, the challenges facing RAN slicing have also been explored, one of which involves designing and managing multiple slices in the shared infrastructure in an efficient manner while guaranteeing the agreed service level agreement (SLA) for each.

Multiplexing of eMBB and URLLC services on the same RAN and sharing the resources for these services is a challenging problem that many researchers pay attention to this topic. In [2], [20] the problem of resource allocation in the coexistence of these two services (URLLC and eMBB) is considered based on their QoS. In [21], the problem of resource allocation for joint eMBB and URLLC is formulated and solved by deep reinforcement learning. In [22] the problem of power minimization for these two type services (URLLC and eMBB) is presented for non-orthogonal multiple access (NOMA) and orthogonal multiple access (OMA). In [23], the authors proposed to allocate RAN resources for the network slicing system in the coexistence of eMBB and URLLC services. The system guarantee the latency, the service rate and the maintenance of the reliability.

Virtualization technique for RAN and core is one of the interested topics. In [24], [25], the authors solve the problem of obtaining beamforming and VMs activation in C-RAN system with limited fronthaul capacity. The goal of this paper is to minimize the energy cost, with the system delay, fronthaul capacity and rate constraint. To guarantee delay for services, transmission and processing delay is modeled based on M/M/1 queuing theory. In [26], [27], the problem of joint virtual computing resource allocation with beamforming is formulated; Also the association of RRH to UE is considered and solved using novel methods.

In this paper, as depicted in Figure 1, the downlink of the ORAN system is studied.

II. SYSTEM MODEL AND PROBLEM FORMULATION

In this section, first, we present the system model. Then, we obtain achievable data rates, power of O-RU and the

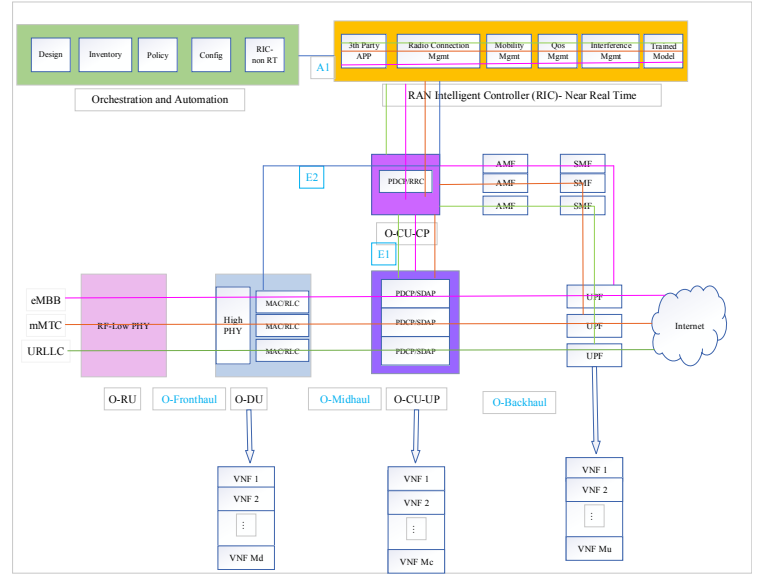


Fig. 1: Network sliced ORAN system

fronthaul capacity for the downlink (DL) of the ORAN system. Afterward, we discuss about the mean delay and the power of VNFs. Finally, the main problem is expressed.

A. System Model

Suppose we have three service types includes mMTC, eMBB and URLLC which support different applications.

Assume we have S_1 , S_2 and S_3 different applications for the first, second and third service type, respectively ($S = S_1 + S_2 + S_3$). So, we have S preallocated slices serving these S services; 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). So each service request s served by its corresponding slice.

Each Service $s_j \in \{1, 2, \dots, S_j\}$ consists of U_s request from the single-antenna UEs which require certain QoS to be able to use the requested program ($j \in \{1, 2, 3\}$ indicate service type). There are different application request which fall into one of these service categories. Each application request requires specific QoS. Based on the request for the application and QoS, UE may be admitted and allocated to the resources. Each slice $s_j \in \{1, 2, \dots, S_j\}$, $j \in \{1, 2\}$ consists of preallocated virtual resource blocks that are mapped to the Physical Resource Blocks (PRBs), 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.

All K PRBs can be assigned to the all UE in each service. Also, each VNF instance is running on the virtual machine (VM) that are using resources from the data centers. Each VM, requires enough resources of CPU, memory, storage and network bandwidth.

In addition, there are R multi-antenna O-RU that are shared between slices. Each O-RU $r \in \{1, 2, \dots, R\}$ has J antenna for transmitting and receiving data. Also $\mathcal{R} = \{r | r \in 1, 2, \dots, R\}$ depicts the set of O-RUs. Moreover, all O-RUs, have access to the all PRBs.

B. The Achievable Rate

The SNR of i^{th} UE served at slice s on PRB k is obtained from

$$\rho_{r,u(s,i)}^k = \frac{|p_{r,u(s,i)}^k \mathbf{h}_{r,u(s,i)}^H \mathbf{w}_{r,u(s,i)}^k|^2}{BN_0 + I_{r,u(s,i)}^k}, \quad (1)$$

where $p_{r,u(s,i)}^k$ represents the transmission power from O-RU r to i^{th} UE served at slice s on PRB k . $\mathbf{h}_{r,u(s,i)}^k \in \mathbb{C}^J$ is the vector of channel gain of a wireless link from r^{th} O-RU to the i^{th} UE in s^{th} slice. In addition, $\mathbf{w}_{r,u(s,i)}^k \in \mathbb{C}^J$ depicts the transmit beamforming vector from r^{th} O-RU to the i^{th} UE in s^{th} slice that is the zero forcing beamforming vector to minimize the interference which is indicated as below

$$\mathbf{w}_{r,u(s,i)}^k = \mathbf{h}_{r,u(s,i)}^k (\mathbf{h}_{r,u(s,i)}^H \mathbf{h}_{r,u(s,i)}^k)^{-1} \quad (2)$$

Moreover, $g_{u(s,i)}^r \in \{0, 1\}$ is the binary variable that illustrates whether O-RU r served the i^{th} UE that is allocated to s^{th} slice or not. Also, BN_0 denotes the power of Gaussian additive noise, and $I_{r,u(s,i)}^k$ is the power of interfering signals represented as follow

$$\begin{aligned} I_{r,u(s,i)}^k &= \underbrace{\sum_{\substack{l=1 \\ l \neq i}}^{U_s} \gamma_1 p_{u(s,l)}^k \sum_{\substack{r'=1 \\ r' \neq r}}^R |\mathbf{h}_{r',u(s,i)}^H \mathbf{w}_{r',u(s,l)}^k g_{u(s,l)}^{r'}|^2}_{(\text{intra-slice interference})} \\ &+ \underbrace{\sum_{\substack{n=1 \\ n \neq s}}^S \sum_{l=1}^{U_s} \gamma_2 p_{u(n,l)}^k \sum_{\substack{r'=1 \\ r' \neq r}}^R |\mathbf{h}_{r',u(s,i)}^H \mathbf{w}_{r',u(n,l)}^k g_{u(n,l)}^{r'}|^2}_{(\text{inter-slice interference})} \\ &+ \underbrace{\sum_{j=1}^R \sigma_{q_{r_j}}^2 |\mathbf{h}_{r,u(s,i)}|^2}_{(\text{Quantization Noise Interference})} \end{aligned} \quad (3)$$

where $\gamma_1 = e_{u(s,i)}^k e_{u(s,l)}^k$ and $\gamma_2 = e_{u(s,i)}^k e_{u(n,l)}^k$. $e_{u(s,i)}^k$ is the binary variable to show whether the k^{th} PRB is allocated to the UE i in slice s , assigned to r^{th} o-RU.

To obtain SNR as formulated in (1), let $y_{u(s,i)}$ be the received signal user i in s^{th} service

$$y_{u(s,i)} = \sum_{r=1}^R \sum_{k=1}^{K_s} \mathbf{h}_{r,u(s,i)}^H g_{u(s,i)}^r e_{r,u(s,i)}^k \mathbf{y}_{r,u(s,i)}^k + z_{u(s,i)}, \quad (4)$$

where $\mathbf{y}_{r,u(s,i)}^k = \mathbf{w}_{r,u(s,i)}^k p_{r,u(s,i)}^{k \frac{1}{2}} x_{u(s,i)} + \mathbf{q}_r$ and $x_{u(s,i)}$ depicts the transmitted symbol vector of UE i in s^{th} set of service, $z_{u(s,i)}$ is the additive Gaussian noise $z_{u(s,i)} \sim \mathcal{N}(0, N_0)$ and N_0 is the noise power. In addition, $\mathbf{q}_r \in \mathbb{C}^J$ indicates the quantization noise, which is made from signal compression in O-DU.

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 $\mathcal{R}_{u(s_1,i)}$ that is formulated as below.

$$\begin{aligned} \mathcal{R}_{r,u(s_1,i)}^k &= B \log_2(1 + \rho_{r,u(s_1,i)}^k), \\ \mathcal{R}_{u(s_1,i)}^r &= \sum_{k=1}^K B \log_2(1 + \rho_{r,u(s_1,i)}^k e_{r,u(s_1,i)}^k), \\ \mathcal{R}_{u(s_1,i)} &= \sum_{r=1}^R \mathcal{R}_{u(s_1,i)}^r g_{u(s_1,i)}^r \end{aligned} \quad (5)$$

where B is the bandwidth of system. $\mathcal{R}_{u(s_1,i)}^r$ is the achievable rate of each RU r to UE i in slice s_1 . Since the blocklength in URLLC and mMTC is finite, the achievable data rate for the i^{th} UE request in the s_j^{th} ($j \in \{2, 3\}$) application of service type 2 (URLLC) and 3 (mMTC) is not achieved from Shannon Capacity formula. So, for the short packet transmission the achievable data rate is approximated from follow

$$\begin{aligned} \mathcal{R}_{r,u(s_j,i)}^k &= B \log_2(1 + \rho_{r,u(s_j,i)}^k - \zeta_{u(s_j,i)}^k) e_{u(s_j,i)}^k, \\ \mathcal{R}_{u(s_j,i)}^r &= \sum_{k=1}^K B (\log_2(1 + \rho_{u(s_j,i)}^k) - \zeta_{u(s_j,i)}^k) e_{u(s_j,i)}^k, \\ \mathcal{R}_{u(s_j,i)} &= \sum_{r=1}^R \mathcal{R}_{u(s_j,i)}^r g_{u(s_j,i)}^r \end{aligned} \quad (6)$$

Where $j \in \{1, 2\}$. Also we have

$$\zeta_{u(s_j,i)}^k = \log_2(e) Q^{-1}(\epsilon) \sqrt{\frac{\mathfrak{C}_{u(s_j,i)}^k}{N_{u(s_j,i)}^k}} \quad (7)$$

Where, ϵ is the transmission probability, Q^{-1} is the inverse of Q- function (Gaussian), $\mathfrak{C}_{u(s_j,i)}^k = 1 - \frac{1}{(1 + \rho_{u(s_j,i)}^k)^2}$ depicts the channel dispersion of UE i at slice s_j , experiencing PRB k and $N_{u(s_j,i)}^k$ represents the blocklength of it. $\mathcal{R}_{u(s_j,i)}^{e,r}$ is the achievable rate of each O-RU r to UE i in slice s_j .

If we replace $p_{u(s,l)}^k$ and $p_{u(n,l)}^k$ in (3) by P_{max} , an upper bound $\bar{I}_{r,u(s,i)}^k$ is obtained for $I_{r,u(s,i)}^k$. Therefore, $\bar{\mathcal{R}}_{u(s,i)} \forall s, \forall i$ is derived by using $\bar{I}_{r,u(s,i)}^k$ instead of $I_{r,u(s,i)}^k$ in (6) and (5).

C. Power of O-RU and Fronthaul Capacity

Let P_r denote the power of transmitted signal from the r^{th} O-RU to UEs served by it. From (4), we have,

$$P_r = \sum_{s=1}^S \sum_{k=1}^{K_s} \sum_{i=1}^{U_s} |\mathbf{w}_{r,u(s,i)}^k|^2 p_{r,u(s,i)}^k g_{u(s,i)}^r e_{r,u(s,i)}^k + \sigma_{q_r}^2. \quad (8)$$

Since we have 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

$$C_r = \log \left(1 + \frac{\sum_{s=1}^S \sum_{k=1}^{K_s} \sum_{i=1}^{U_s} |\mathbf{w}_{r,u(s,i)}^k|^2 \alpha_{r,u(s,i)}^k}{\sigma_{q_r}^2} \right), \quad (9)$$

Where, $\alpha_{r,u(s,i)}^k = p_{r,u(s,i)}^k g_{u(s,i)}^r e_{r,u(s,i)}^k$ and $\sigma_{q_r}^2$ is the power of quantization noise.

D. Mean Delay

In this part, the end to end mean delay for a service is obtained. Suppose the mean total delay is depicted as T_{tot} .

$$\begin{aligned} T_{tot} &= T_{process} + T_{transmission} + T_{propagation} \\ T_{process} &= T_{RU} + T_{DU} + T_{CU} + T_{UPF} \\ T_{transmission} &= T_{front} + T_{mid} + T_{back} + T_{trans2net} \\ T_{propagation} &= T_{front} + T_{mid} + T_{back} + T_{trans2net} \end{aligned} \quad (10)$$

Total delay is sum of processing delay, transmission delay and propagation delay. The propagation delay is the time takes for a signal to reach to its destination. So it has a constant value based on the length of fiber link ($T = L/c$, where L is the length of link and c is the speed of signal). Also, the transmission delay is the amount of time required to push all the packets into the fiber link. ($T = \frac{\alpha}{R}$ Where, R is the rate of transmission in each link and α is the mean arrival data rate of the each link which is constant in this model.) Here we assume the value of propagation delay and transmission is negligible compared to the rest.

$$T_{tot} \approx T_{process} \quad (11)$$

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} slice. Therefore, the mean arrival data rate of the s^{th} slice in the UPF layer is $\alpha_s^1 = \sum_{u=1}^{U_s} a_{u(s,i)} \lambda_{u(s,i)}$, where $a_{u(s,i)}$ is a binary variable which indicates whether the i^{th} UE requested s^{th} service is admitted or not.

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 O-DU (α_s^D). so $\alpha_s = \alpha_s^U \approx \alpha_s^C \approx \alpha_s^D$. since, by using Burkes Theorem, the mean arrival data rate of the second and third layer 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 the incoming traffic to VNFs equally. Suppose the baseband processing of each VNF is depicted as M/M/1 processing queue. Each packet is processed by one of the VNFs of a slice. So, the mean delay for the s^{th} slice in the first and the second layer, modeled as M/M/1 queue, is formulated as follow, respectively

$$\begin{aligned} T_{DU}^s &= \frac{1}{\mu_s^d - \alpha_s/M_s^d}, \\ T_{CU}^s &= \frac{1}{\mu_s^c - \alpha_s/M_s^c}, \\ T_{UPF}^s &= \frac{1}{\mu_s^u - \alpha_s/M_s^u} \end{aligned} \quad (12)$$

Where M_s^d , M_s^c and M_s^u are the variables that depict the sum 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. Besides, α_s is the arrival rate which is divided by load balancer before arriving to the VNFs. The arrival rate of each VNF in each layer for each slice s is α_s/M_s^i $i \in \{d, c, u\}$.

In addition, $T_{RU}^{u(s,i)}$ is the mean transmission delay of i^{th} UE in s^{th} service on the wireless link. The arrival data rate of wireless link for each UE i in service s is $\lambda_{u(s,i)}$. As a result we have, $\sum_{i=1}^{U_s} \lambda_{u(s,i)} = \alpha_s$. 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.

Therefore, the mean delay of the transmission layer for UE i in slice s is

$$T_{RU}^{u(s,i)} = \frac{1}{R_{u(s,i)} - \lambda_{u(s,i)}}; \quad (13)$$

So the mean processing delay for each UE i in slice s is

$$T_{process}^{u(s,i)} = T_{RU}^{u(s,i)} + T_{DU}^s + T_{CU}^s + T_{UPF}^s \quad (14)$$

Hence, $T_{tot}^{u(s,i)} \approx T_{process}^{u(s,i)}$

E. VNF Power

Assume the power consumption of baseband processing at each DC d that is connected to VNFs of a slice s is depicted as ϕ_s . So the total power of the system for all active DCs that are connected to slices can be represented as

$$\phi_{tot} = \sum_{s=1}^S \phi_s.$$

Where, ϕ_s is obtained from below

$$\phi_s = M_s^u \phi_s^u + M_s^c \phi_s^c + M_s^d \phi_s^d \quad (15)$$

Moreover, ϕ_s^u , ϕ_s^c and ϕ_s^d are the static cost of energy in UPF, O-CU and O-DU, respectively.

F. Problem Statement

Suppose each slice s has priority factor δ_s where $\sum_{s=1}^S \delta_s = 1$. The optimization problem is formulated as follow. The aim of this paper is to maximize the sum rate of all UEs with the presence of constraints which is written as follow,

$$\max_{P, E, M, G} \sum_{s=1}^S \sum_{i=1}^{U_s} \delta_s \bar{R}_{u(s,i)} \quad (16a)$$

$$\text{subject to } P_r \leq P_{max} \quad \forall r \quad (16b)$$

$$p_{r,u(s,i)}^k \geq 0 \quad \forall i, \forall r, \forall s, \forall k, \quad (16c)$$

$$\bar{R}_{u(s,j)} \geq \mathcal{R}_{min}^j \quad \forall s, j \in \{1, 2, 3\}, \quad (16d)$$

$$C^r \leq C_{max}^r \quad \forall r, \quad (16e)$$

$$T_{tot}^{u(s,i)} \leq T_{max}^s \quad \forall i, \forall s, \quad (16f)$$

$$\mu_s \geq \alpha_s/M_s \quad \forall s, \quad (16g)$$

$$\bar{R}_{u(s,i)} \geq \lambda_{u(s,i)} \quad \forall i, \forall s, \quad (16h)$$

$$0 \leq M_s \leq M^{max} \quad \forall s, \quad (16i)$$

$$\sum_r g_{u(s,i)}^r = 1 \quad \forall s, \forall i, \quad (16j)$$

$$\sum_{k=1}^{K_s} g_{u(s,i)}^r e_{r,u(s,i)}^k \geq 1 \quad \forall s, \forall i, \forall r \quad (16k)$$

$$\sum_{s=1}^S \sum_{i=1}^{U_s} g_{u(s,i)}^r e_{r,u(s,i)}^k \leq 1 \quad \forall s, \forall i, \forall r \quad (16l)$$

$$\phi_{tot} \leq \phi_{max}, \quad (16m)$$

$$g_{u(s,i)}^r \in \{0, 1\} \quad \forall s, \forall i, \quad (16n)$$

$$e_{r,u(s,i)}^k \in \{0, 1\} \quad \forall s, \forall i, \quad (16o)$$

where $\mathbf{P} = [p_{r,u(s,i)}^k] \quad \forall s, \forall i, \forall r, \forall k$, is the matrix of power for UEs, $\mathbf{E} = [e_{r,u(s,i)}^k] \quad \forall s, \forall i, \forall r, \forall k$ indicate the binary variable for PRB association. Moreover, $\mathbf{G} = [g_{u(s,i)}^r] \quad \forall s, \forall i, \forall r$ is a binary variable for O-RU association. Furthermore, $\mathbf{M} = [M_s^d, M_s^c, M_s^u] \quad \forall s$ is the matrix that shown the number of VNFs in each layer of slice. (16b), and (16c), indicate that the power of each RU do not exceed the maximum power, and the power of each UE is a positive integer value, respectively. Also (16d) shows that the rate of each UE requesting eMBB, URLLC and mMTC is more than a threshold, respectively. (16e) and (16f) expressed the limited capacity of the fronthaul link, and the limited delay of receiving signal, respectively. (16g) and (16h) denoted the stability of the M/M/1 queue model. (16i) restricted the number of VNF in each slice due to the limited resources. (16j) and (16k) guarantee that O-RU and PRB is associated to the UE, respectively. Also, (16l) ensure that each PRB can not be assigned to more than one UE associated to the same O-RU. In addition, (16m) indicate that the static cost of energy of VNFs in each slice do not exceed from the threshold. Moreover, (16n) and (16o) depict that \mathbf{E} and \mathbf{G} are matrix of binary variables.

III. PROPOSED ALGORITHM SCHEME

In this section, we first apply some simplifications to the system; Solving problem (16) is complicated due to the fact that this problem is a non-convex problem and it is a mixed integer non-linear problem (MINLP) with a binary variable and an integer variable. In the following, we apply the simplifications to reformulate MINLP parts and use iterative heuristic algorithm to solve the reformulated problem. We solve this problem in two level iteratively until it converges; In the first level, parameters ($\mathbf{P}, \mathbf{E}, \mathbf{M}$) are obtained by relaxing and reformulating parameters and turn it to convex problem; Afterward we solve it by dual optimization problem. In the second level, finding optimal O-RU association (\mathbf{G}) is concerned with the fixed parameter of power, PRB allocation and number of VNFs. We repeat this procedure until the algorithm converges.

A. Sub-Problem 1

Suppose that \mathbf{G} is fixed, we want to obtain \mathbf{P}, \mathbf{E} and \mathbf{M} . Here, we first simplify and relax the parameters to convexify the problem.

As we mentioned before, by replacing $p_{u(s,i)}^k$ and $p_{u(n,i)}^k$ in (3) by P_{max} , an upper bound $\bar{I}_{r,u(s,i)}^k$ for $I_{r,u(s,i)}^k$, the lower bound $\bar{\rho}_{u(s,i)}^k$ for $\rho_{u(s,i)}^k$ and the lower bound $\bar{\mathcal{R}}_{u(s,i)} \quad \forall s, \forall i$ for $\mathcal{R}_{u(s,i)}$ is obtained by replacing with $I_{r,u(s,i)}^k \quad \bar{I}_{r,u(s,i)}^k$ in (6) and (5) and make them concave.

Suppose $\hat{\rho}_{r,u(s,i)}^k = \frac{|P_{max} \mathbf{h}_{r,u(s,i)}^H \mathbf{w}_{r,u(s,i)}^k g_{u(s,i)}^r|}{BN_0}$. To convexify (6) (for the short packet transmission), we replace

$\rho_{r,u(s,i)}^k$ with $\hat{\rho}_{r,u(s,i)}^k$ in (7). So, a lower bound for (6) is given that is a concave function.

$$\begin{aligned} \bar{\mathcal{R}}_{u(s,i)}^r &= \sum_{k=1}^{K_{s_j}} B(\log_2(1 + \bar{\rho}_{u(s,i)}^k) - \hat{\zeta}_{u(s,i)}^k) e_{u(s,i)}^k \\ \bar{\mathcal{R}}_{u(s,i)} &= \sum_{r=1}^R \bar{\mathcal{R}}_{u(s,i)}^r \\ \hat{\zeta}_{u(s,i)}^k &= \log_2(e) Q^{-1}(\epsilon) \sqrt{\frac{\hat{\mathcal{C}}_{u(s,i)}^k}{N_{u(s,i)}^k}} \\ \hat{\mathcal{C}}_{u(s,i)}^k &= 1 - \frac{1}{(1 + \hat{\rho}_{u(s,i)}^k)^2} \end{aligned} \quad (17)$$

Consider UPF, O-CU and O-DU have the same processor (for simplification), so 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 approximately equal to the mean arrival data rate of the O-CU-UP layer (α_s^C) and O-DU (α_s^D). so $\alpha_s = \alpha_s^U \approx \alpha_s^C \approx \alpha_s^D$. So the given assumption leads to have same energy for each layer $\phi_s^u = \phi_s^c = \phi_s^d$. As a result of these assumption, for simplicity, we can assume that $M_s = M_s^u = M_s^c = M_s^d$. Using the above assumption, we have $T_{DU}^s = T_{CU}^s = T_{UPF}^s$

$$\begin{aligned} T_{process}^s &= T_{RU}^s + T_{DU}^s + T_{CU}^s + T_{UPF}^s \\ T_{process}^s &= T_{RU}^s + 3 \times T_{DU}^s. \end{aligned} \quad (18)$$

Lemma 1. In problem (16), the constraint (16f) can be reformulated as below $\forall i, \forall s$

$$\begin{aligned} T_{max}^s &\geq \frac{1}{R_{u(s,i)} - \lambda_{u(s,i)}} + \frac{3}{\mu_s - \alpha_s/M_s} \\ M_s &\geq \frac{\alpha_s(T_{max}^s R_{u(s,i)} - T_{max}^s \lambda_{u(s,i)} - 1)}{(T_{max}^s \mu_s - 3)(R_{u(s,i)} - \lambda_{u(s,i)}) - \mu_s} \end{aligned} \quad (19)$$

Also from equation (16m), (16g) and (16i) we have

$$0 \leq M_s \leq \min\{M^{max}, \alpha_s/\mu_s, \phi_{max}/3\phi_s\} \quad (20)$$

We denote $\mathfrak{M}_s = \min\{M^{max}, \alpha_s/\mu_s, \phi_{max}/3\phi_s\}$. Thus, if we restrict (16f) to equality we have

$$0 \leq \frac{\alpha_s(T_{max}^s R_{u(s,i)} - T_{max}^s \lambda_{u(s,i)} - 1)}{(T_{max}^s \mu_s - 3)(R_{u(s,i)} - \lambda_{u(s,i)}) - \mu_s} \leq \mathfrak{M}_s \quad (21)$$

In (21), $0 \leq \frac{\alpha_s(T_{max}^s R_{u(s,i)} - T_{max}^s \lambda_{u(s,i)} - 1)}{(T_{max}^s \mu_s - 3)(R_{u(s,i)} - \lambda_{u(s,i)}) - \mu_s}$ is established due to the fact that the numerator and the denominator will both have same sign. Using (16h), in numerator, $\alpha_s \geq 0$, $R_{u(s,i)} - \lambda_{u(s,i)} \geq 0$ and to simplify the problem, assume $(R_{u(s,i)} - \lambda_{u(s,i)})T_{max}^s \geq 1$ since the order of T_{max}^s is about milli second and the difference between achievable rate and packet rate can be more than $1/T_{max}^s$. Therefore, we restrict constraint (16h) to $R_{u(s,i)} \geq \lambda_{u(s,i)} + 1/T_{max}^s$. So the numerator is positive. In denominator, it can be said approximately that $(T_{max}^s \mu_s)(R_{u(s,i)} - \lambda_{u(s,i)}) - \mu_s \geq 0$, since, $(R_{u(s,i)} - \lambda_{u(s,i)}) \geq 1/T_{max}^s$ as mentioned above. Therefore, we just need to have constraint below

$$\frac{\alpha_s(T_{max}^s R_{u(s,i)} - T_{max}^s \lambda_{u(s,i)} - 1)}{(T_{max}^s \mu_s - 3)(R_{u(s,i)} - \lambda_{u(s,i)}) - \mu_s} \leq \mathfrak{M}_s \quad (22)$$

So by reformulating the equation (22), we have a new constraint $\forall i, \forall s$ as below,

$$\begin{aligned}\mathcal{R}_{u(s,i)} &\geq \frac{\mathfrak{M}_s((T_{max}^s \mu_s - 3)\lambda_{u(s,i)} + \mu_s) - \alpha_s(T_{max}^s \lambda_{u(s,i)} + 1)}{\mathfrak{M}_s(T_{max}^s \mu_s - 3) - \alpha_s T_{max}^s}, \\ \varpi_{u(s,i)} &= \frac{\mathfrak{M}_s((T_{max}^s \mu_s - 3)\lambda_{u(s,i)} + \mu_s) - \alpha_s(T_{max}^s \lambda_{u(s,i)} + 1)}{\mathfrak{M}_s(T_{max}^s \mu_s - 3) - \alpha_s T_{max}^s}, \\ \mathcal{R}_{u(s,i)} &\geq \varpi_{u(s,i)}.\end{aligned}\quad (23)$$

In addition, we denote $\mathfrak{M}_{u(s,i)} = \frac{\alpha_s(T_{max}^s R_{u(s,i)} - T_{max}^s \lambda_{u(s,i)} - 1)}{(T_{max}^s \mu_s - 3)(R_{u(s,i)} - \lambda_{u(s,i)}) - \mu_s}$. So we have,

$$M_s = \max\{\mathfrak{M}_{u(s,i)} | i \in 1, 2, \dots, U_s\} \quad \forall s. \quad (24)$$

Despite simplifying the problem (16), it is still non-convex and hard to be solved. So the simplest approach is to relax \mathbf{E} into continuous value $e_{r,u(s,i)}^k \in [0, 1] \forall s, \forall i, \forall r, \forall k$. Furthermore, the problem can be solved using the Lagrangian function and iterative algorithm.

In order to make (16) as a standard form of a convex optimization problem, it is required to change the variable of equations (9) to $P_r = \sigma_{q_r}^2 \times 2^{C_r}$ so the constraint (16e) is changed to $P_r \leq \sigma_{q_r}^2 \times 2^{C_{max}^r}$. The combination of equations (16d) and (16e) leads to the following equation

$$\begin{aligned}\zeta_r &= \min\{P_{max}, \sigma_{q_r}^2 \times 2^{C_{max}^r}\}, \\ P_r &\leq \zeta_r.\end{aligned}\quad (25)$$

Moreover, the combination of equations (16d), (16h) and (23) leads to the following equation

$$\begin{aligned}\eta_{u(s,i)} &= \max\{\mathcal{R}_{u(s,i)}^{max}, \lambda_{u(s,i)} + 1/T_{max}^s, \varpi_{u(s,i)}\}, \\ \bar{\mathcal{R}}_{u(s,i)} &\geq \eta_{u(s,i)}.\end{aligned}\quad (26)$$

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

The Lagrangian function is written as follow

$$\mathcal{L}(P, \mathbf{E}; \mathbf{v}, \mathbf{\chi}, \mathbf{h}, \mathbf{\xi}, \mathbf{\kappa}, \mathbf{m}) = \sum_{s=1}^S \sum_{i=1}^{U_s} \delta_s \bar{\mathcal{R}}_{u(s,i)} \quad (27a)$$

$$+ \sum_{s=1}^S \sum_{i=1}^{U_s} \mathfrak{h}_{u(s,i)} (\bar{\mathcal{R}}_{u(s,i)} - \eta_{u(s,i)}) \quad (27b)$$

$$- \sum_{r=1}^R \mathbf{m}_r (P_r - \zeta_r) \quad (27c)$$

$$+ \sum_{s=1}^S \sum_{i=1}^{U_s} \sum_{k=1}^K \sum_{r=1}^R \kappa_{r,u(s,i)}^k p_{r,u(s,i)}^k \quad (27d)$$

$$+ \sum_{r=1}^R \sum_{s=1}^S \sum_{i=1}^{U_s} \chi_{r,u(s,i)} \left(\sum_{k=1}^{K_s} e_{r,u(s,i)}^k - 1 \right) \quad (27e)$$

$$- \sum_{s=1}^S \sum_{i=1}^{U_s} \sum_{k=1}^K \sum_{r=1}^R \mathbf{v}_{r,u(s,i)}^k (e_{r,u(s,i)}^k - 1) \quad (27f)$$

$$+ \sum_{s=1}^S \sum_{i=1}^{U_s} \sum_{k=1}^K \sum_{r=1}^R \xi_{r,u(s,i)}^k e_{r,u(s,i)}^k. \quad (27g)$$

Lemma 2. By taking derivatives of (27) (the lagrangian function), with respect to the \mathbf{P} and the \mathbf{E} , these two variables are obtained. Assume, $e_{r,u(s,i)}^k = 1$

$$\frac{\partial \mathcal{L}}{\partial p_{r,u(s,i)}^k} = (\delta_s + \mathfrak{h}_{u(s,i)}) \mathfrak{B}_{r,u(s,i)}^k + (\kappa_{r,u(s,i)}^k - \mathbf{m}_r \mathfrak{D}_{r,u(s,i)}^k) = 0 \quad (28)$$

Where

$$\begin{aligned}\mathfrak{D}_{r,u(s,i)}^k &= |\mathbf{w}_{r,u(s,i)}^k|^2 g_{u(s,i)}^r e_{r,u(s,i)}^k, \\ \mathfrak{B}_{r,u(s,i)}^k &= \frac{B |\mathbf{h}_{r,u(s,i)}^{Hk} \mathbf{w}_{r,u(s,i)}^k|^2 g_{u(s,i)}^r e_{r,u(s,i)}^k}{\ln(2)} \mathfrak{S}_{r,u(s,i)}^k, \\ \mathfrak{S}_{r,u(s,i)}^k &= \frac{1}{|\mathbf{h}_{r,u(s,i)}^{Hk} \mathbf{w}_{r,u(s,i)}^k|^2 \mathfrak{I}_{r,u(s,i)}^k + BN_0 + I_{r,u(s,i)}^k}.\end{aligned}\quad (29)$$

where $\mathfrak{I}_{r,u(s,i)}^k = g_{u(s,i)}^r e_{r,u(s,i)}^k p_{r,u(s,i)}^k$. Thus, from equation (28), optimal power is obtained and power is allocated. We denote $\mathfrak{J}_{r,u(s,i)}^k = g_{u(s,i)}^r e_{r,u(s,i)}^k$.

$$p_{r,u(s,i)}^k = \left[\frac{(\delta_s + \mathfrak{h}_{u(s,i)}) \mathfrak{B}_{r,u(s,i)}^k}{\kappa_{r,u(s,i)}^k - \mathbf{m}_r \mathfrak{D}_{r,u(s,i)}^k} - \frac{BN_0 + I_{r,u(s,i)}^k}{|\mathbf{h}_{r,u(s,i)}^{Hk} \mathbf{w}_{r,u(s,i)}^k|^2 \mathfrak{J}_{r,u(s,i)}^k} \right]^+.$$

Also $[a]^+ = \max(0, a)$. In addition, PRB assignment is obtained as follow

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial e_{r,u(s,i)}^k} &= \bar{\mathcal{R}}_{r,u(s,i)}^k (\delta_s + \mathfrak{h}_{u(s,i)}) \\ &\quad - \mathbf{m}_r |\mathbf{w}_{r,u(s,i)}^k|^2 p_{r,u(s,i)}^k g_{u(s,i)}^r \\ &\quad + (\xi_{r,u(s,i)}^k - \mathbf{v}_{r,u(s,i)}^k + \chi_{r,u(s,i)}) = 0.\end{aligned}\quad (31)$$

Using KKT conditions, we have

$$e_{r,u(s,i)}^k \times (\mathfrak{F}_{r,u(s,i)}^k - \mathbf{v}_{r,u(s,i)}^k - \mathbf{m}_r |\mathbf{w}_{r,u(s,i)}^k|^2 p_{r,u(s,i)}^k g_{u(s,i)}^r) = 0. \quad (32)$$

Where $\mathfrak{F}_{r,u(s,i)}^k = \bar{\mathcal{R}}_{r,u(s,i)}^k (\delta_s + \mathfrak{h}_{u(s,i)}) + (\xi_{r,u(s,i)}^k + \chi_{r,u(s,i)})$. Hence, from equation (31) and (32), PRB assignment is performed as follow.

$$e_{r,u(s,i)}^k = \begin{cases} 1 & u(s,i) = \arg\max \mathfrak{F}_{r,u(s,i)}^k \forall s, \forall r, \forall k \\ 0 & \text{otherwise} \end{cases} \quad (33)$$

Thus the user in each slice s that have the largest value of $\mathfrak{F}_{r,u(s,i)}^k$, should be allocated to the PRB k ; Due to the fact that just one PRB can be allocated to a UE between those UEs (regardless to the services) that are associated to the same O-RU.

B. Sub-Problem 2

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

Greedy Algorithm for Non-Comp O-RU Assignment: The problem can be reformulated as follow

$$\max_{\mathbf{G}} \sum_{s=1}^S \sum_{i=1}^{U_s} \sum_{r=1}^R \delta_s g_{u(s,i)}^r \bar{\mathcal{R}}_{u(s,i)}^r \quad (34a)$$

$$\text{subject to } \sum_{s=1}^S \sum_{i=1}^{U_s} g_{u(s,i)}^r \psi_{r,u(s,i)} \leq \mathbf{t}_r \quad \forall r \quad (34b)$$

$$\sum_r g_{u(s,i)}^r = 1 \quad \forall s, \forall i, \quad (34c)$$

$$g_{u(s,i)}^r \in \{0, 1\} \quad \forall s, \forall i, \quad (34d)$$

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 $\mathbf{t}_r = \zeta_r - \sigma_r$. Since we obtained (26) in (III-A), we can ignore this constraint in (34). The problem (34) is an NP-complete 0-1 multiple knapsack problem. We solve this problem using GAAOU which is a greedy algorithm (1) as follow [15], [28]. Firstly, we set all variables $g_{u(s,i)}^r = 0, \forall s, \forall i, \forall r$. Then we define $\mathfrak{B}_{u(s,i)}^{rem} = \mathcal{R} \forall s, \forall i$ and $\mathfrak{C}_r = \mathbf{t}_r, \forall r$ as a set of all O-RUs and value of each O-RU, respectively. Next, we sort all slices based on their priority. Afterward, we assign the O-RU that provides the highest achievable rate for each UE (we start from the UEs on the slices with highest priority) on the condition that it does not exceed the value of each O-RU (that is a function of maximum power and capacity of O-RU). If it exceeds the value of O-RU, then O-RU with the next highest achievable rate is selected. The complexity of sorting S slices based on their priority is $O(\text{Slog}(S))$. Depict $\mathfrak{N} = \sum_{s=1}^S \sum_{i=1}^{U_s} 1$. The complexity of this algorithm is about $O(\text{Slog}(S)) + O(R \times \mathfrak{N})$.

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

```

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

```

C. Iterative Proposed Algorithm

In (III-A) and (III-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} , then \mathbf{P} and \mathbf{E} is achieved using Lagrangian method. Afterward, \mathbf{G} is updated using GAAOU algorithm. This process is repeated until it converges. The whole algorithm is depicted as follow (Algorithm (2)).

Algorithm 2 Iterative Algorithm for Power Allocation, PRB, VNF and O-RU Association (IAPPVO)

```

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

```

REFERENCES

- [1] X. Shen, J. Gao, W. Wu, K. Lyu, M. Li, W. Zhuang, X. Li, and J. Rao, "Ai-assisted network-slicing based next-generation wireless networks," *IEEE Open Journal of Vehicular Technology*, vol. 1, pp. 45–66, 2020.
- [2] M. Setayesh, S. Bahrani, and V. W. Wong, "Joint prb and power allocation for slicing embb and urllc services in 5g c-ran," in *GLOBECOM 2020-2020 IEEE Global Communications Conference*. IEEE, 2020, pp. 1–6.
- [3] P. Popovski, K. F. Trillingsgaard, O. Simeone, and G. Durisi, "5g wireless network slicing for embb, urllc, and mmcc: A communication-theoretic view," *Ieee Access*, vol. 6, pp. 55 765–55 779, 2018.
- [4] R. Mijumbi, J. Serrat, J.-L. Gorricho, N. Bouten, F. De Turck, and R. Boutaba, "Network function virtualization: State-of-the-art and research challenges," *IEEE Communications surveys & tutorials*, vol. 18, no. 1, pp. 236–262, 2015.
- [5] Z. Luo and C. Wu, "An online algorithm for vnf service chain scaling in datacenters," *IEEE/ACM Transactions on Networking*, vol. 28, no. 3, pp. 1061–1073, 2020.
- [6] B. Han, L. Liu, J. Zhang, C. Tao, C. Qiu, T. Zhou, Z. Li, and Z. Piao, "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, 2019.
- [7] L. Gavrilovska, V. Rakovic, and D. Denkovski, "From cloud ran to open ran," *Wirel. Pers. Commun.*, vol. 113, no. 3, pp. 1523–1539, 2020.
- [8] S. Niknam, A. Roy, H. S. Dhillon, S. Singh, R. Banerji, J. H. Reed, N. Saxena, and S. Yoon, "Intelligent o-ran for beyond 5g and 6g wireless networks," *arXiv preprint arXiv:2005.08374*, 2020.
- [9] 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, 2021.
- [10] C. B. Both, J. Borges, L. Gonçalves, C. Nahum, C. Macedo, A. Klautau, and K. Cardoso, "System intelligence for uav-based mission critical with challenging 5g/b5g connectivity," *arXiv preprint arXiv:2102.02318*, 2021.
- [11] "O-ran architecture description," O-RAN Alliance, Tech. Rep., 2020.
- [12] O.-R. W. G. 2, "Ai/ml workflow description and requirements," O-RAN Alliance, Tech. Rep., 2020.
- [13] B.-S. Lin, "Toward an ai-enabled o-ran-based and sdn/nfv-driven 5g& iot network era," *Network and Communication Technologies*, vol. 6, no. 1, pp. 6–15, 2021.

- [14] L. Feng, Y. Zi, W. Li, F. Zhou, P. Yu, and M. Kadoch, "Dynamic resource allocation with ran slicing and scheduling for urllc and embb hybrid services," *IEEE Access*, vol. 8, pp. 34 538–34 551, 2020.
- [15] Y. L. Lee, J. Loo, T. C. Chuah, and L.-C. Wang, "Dynamic network slicing for multitenant heterogeneous cloud radio access networks," *IEEE Transactions on Wireless Communications*, vol. 17, no. 4, pp. 2146–2161, 2018.
- [16] Y. L. Lee, J. Loo, and T. C. Chuah, "A new network slicing framework for multi-tenant heterogeneous cloud radio access networks," in *2016 International Conference on Advances in Electrical, Electronic and Systems Engineering (ICAEEES)*. IEEE, 2016, pp. 414–420.
- [17] H. Xiang, S. Yan, and M. Peng, "A realization of fog-ran slicing via deep reinforcement learning," *IEEE Transactions on Wireless Communications*, vol. 19, no. 4, pp. 2515–2527, 2020.
- [18] S. E. Elayoubi, S. B. Jemaa, Z. Altman, and A. Galindo-Serrano, "5g ran slicing for verticals: Enablers and challenges," *IEEE Communications Magazine*, vol. 57, no. 1, pp. 28–34, 2019.
- [19] S. D'Oro, F. Restuccia, and T. Melodia, "Toward operator-to-waveform 5g radio access network slicing," *IEEE Communications Magazine*, vol. 58, no. 4, pp. 18–23, 2020.
- [20] P. Yang, X. Xi, T. Q. Quek, J. Chen, X. Cao, and D. Wu, "How should i orchestrate resources of my slices for bursty urllc service provision?" *IEEE Transactions on Communications*, vol. 69, no. 2, pp. 1134–1146, 2020.
- [21] M. Alsenwi, N. H. Tran, M. Bennis, S. R. Pandey, A. K. Bairagi, and C. S. Hong, "Intelligent resource slicing for embb and urllc coexistence in 5g and beyond: A deep reinforcement learning based approach," *IEEE Transactions on Wireless Communications*, 2021.
- [22] F. Saggese, M. Moretti, and P. Popovski, "Power minimization of downlink spectrum slicing for embb and urllc users," *arXiv preprint arXiv:2106.08847*, 2021.
- [23] P. Korrai, E. Lagunas, S. K. Sharma, S. Chatzinotas, A. Bandi, and B. Ottersten, "A ran resource slicing mechanism for multiplexing of embb and urllc services in ofdma based 5g wireless networks," *IEEE Access*, vol. 8, pp. 45 674–45 688, 2020.
- [24] J. Tang, W. P. Tay, T. Q. Quek, and B. Liang, "System cost minimization in cloud ran with limited fronthaul capacity," *IEEE Transactions on Wireless Communications*, vol. 16, no. 5, pp. 3371–3384, 2017.
- [25] K. Guo, M. Sheng, J. Tang, T. Q. Quek, and Z. Qiu, "Exploiting hybrid clustering and computation provisioning for green c-ran," *IEEE Journal on Selected Areas in Communications*, vol. 34, no. 12, pp. 4063–4076, 2016.
- [26] P. Luong, F. Gagnon, C. Despins, and L.-N. Tran, "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, 2018.
- [27] P. Luong, C. Despins, F. Gagnon, and L.-N. Tran, "A novel energy-efficient resource allocation approach in limited fronthaul virtualized c-rans," in *2018 IEEE 87th Vehicular Technology Conference (VTC Spring)*. IEEE, 2018, pp. 1–6.
- [28] Y. Akçay, H. Li, and S. H. Xu, "Greedy algorithm for the general multidimensional knapsack problem," *Annals of Operations Research*, vol. 150, no. 1, pp. 17–29, 2007.