

Network Slicing and Resource Allocation in an Open RAN System

Mojdeh Karbalaee Motalleb

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

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

Abstract—

Index Terms—

I. INTRODUCTION

II. SYSTEM MODEL AND PROBLEM FORMULATION

In this section, first, we present the system model. Then, we obtain achievable data rates and delays for the downlink (DL) of the ORAN system. Afterward, we discuss about assignment of physical data center resources. Finally, the main problem is expressed.

A. System Model

Suppose we have two service types includes eMBB and URLLC. Assume we have V_1 and V_2 different applications for the first and second service type, respectively ($V = V_1 + V_2$). Assume we have S preallocated slices serving V services; There are S_1 slices for the first service type (eMBB) and S_2 slices for the second service type (URLLC) ($S = S_1 + S_2$). Each Service $v_j \in \{1, 2, \dots, V_j\}$ consists of U_{v_j} request from the single-antenna UEs which require certain QoS to be able to use the requested program ($j \in \{1, 2\}$ 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 slice. Each slice $s_j \in \{1, 2, \dots, S_j\}$, $j \in \{1, 2\}$ consists of K_{s_j} , $j \in \{1, 2\}$ preallocated virtual resource blocks that are mapped to 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.

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 RU that are shared between slices. Each RU $r \in \{1, 2, \dots, R\}$ has J antenna for transmitting and receiving data. Moreover, all RUs, have access to PRBs.

B. The Achievable Rate

The SNR of i^{th} UE in v^{th} service experienced at slice s on PRB k which is obtained from $\rho_{u(v,i)}^{k,s}$ is the SNR

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

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

$$\mathbf{w}_{r,u(v,i)}^{k,s} = \mathbf{h}_{r,u(v,i)}^{k,s} (\mathbf{h}_{r,u(v,i)}^{H k,s1} \mathbf{h}_{r,u(v,i)}^{k,s})^{-1} \quad (2)$$

Moreover, $g_{u(v,i)}^r \in \{0, 1\}$ is a binary variable that illustrates whether RU r is mapped to the i^{th} UE in v^{th} service or not. Also, BN_0 denotes the power of Gaussian additive noise, and $I_{r,u(v,i)}^{k,s1}$ is the power of interfering signals represented as follow

$$I_{r,u(v,i)}^{k,s} = \underbrace{\sum_{\substack{l=1 \\ l \neq i}}^{U_v} \gamma_1 p_{u(v,i)}^{k,s} \sum_{\substack{r'=1 \\ r' \neq r}}^R |\mathbf{h}_{r',u(v,i)}^{H k,s} \mathbf{w}_{r',u(v,i)}^{k,s1} g_{u(v,i)}^{r'}|^2}_{\text{(intra-service interference)}} + \underbrace{\sum_{\substack{y=1 \\ y \neq v}}^V \sum_{n=1}^S \sum_{l=1}^{U_y} \gamma_2 p_{u(y,l)}^{k,n} \sum_{\substack{r'=1 \\ r' \neq r}}^R |\mathbf{h}_{r',u(v,i)}^{H k,s} \mathbf{w}_{r',u(y,l)}^{k,n} g_{u(y,l)}^{r'}|^2}_{\text{(inter-service interference)}} \quad (3)$$

where $\gamma_1 = e_{r,u(v,i)}^{k,s} e_{r',u(v,i)}^{k,s} a_{u(v,i)} a_{u(v,l)} b_{v,s}$ and $\gamma_2 = e_{r,u(v,i)}^{k,s} e_{r',u(y,l)}^{k,n} a_{u(v,i)} a_{u(y,l)} b_{v,s} b_{y,n}$. Where $a_{u(v,i)} \in \{0, 1\}$ is a binary variable to depict user admission. $b_{v,s}$ is a binary variable that illustrates whether slice s is allocated to service v or not. $e_{r,u(v,i)}^{k,s}$ is the binary variable to show whether the k^{th} PRB is allocated to the UE i in service v , assigned to r^{th} o-RU using slice s or not.

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

$$\mathcal{R}_{u(v,i)}^e = \sum_{s=1}^{S_1} \sum_{k=1}^{K_{s_1}} B \log_2(1 + \rho_{u(v,i)}^{k,s1}) a_{u(v,i)} b_{v,s1} e_{r,u(v,i)}^{k,s1}, \quad (4)$$

where B is the bandwidth of system.

Since the blocklength in URLLC is finite, the achievable data rate for the i^{th} UE request in the v^{th} application of service type 2 (URLLC) is not achieved from Shannon

Capacity formula. So, for the short packet transmission the achievable data rate is approximated from follow

$$\mathcal{R}_{u(v_2,i)}^u = \sum_{s_2=1}^{S_2} \sum_{k=1}^{K_{s_2}} B(\log_2(1 + \rho_{u(v_2,i)}^{k,s_2}) - \zeta_{u(v_2,i)}^{k,s_2}) \beta_{u(v_2,i)}^{k,s_2} \quad (5)$$

Where $\beta_{u(v_2,i)}^{k,s_2} = a_{u(v_2,i)} b_{v_2,s_1} e_{u(v_2,i)}^{k,s_2}$ and $\zeta_{u(v_2,i)}^{k,s_2} = \log_2(e) Q^{-1}(\epsilon) \sqrt{\frac{C_{u(v_2,i)}^{k,s_2}}{N_{u(v_2,i)}^{k,s_2}}}$ Where, ϵ is the transmission probability, Q^{-1} is the inverse of Q- function (Gaussian), $C_{u(v_2,i)}^{k,s_2} = 1 - \frac{1}{(1 + \rho_{u(v_2,i)}^{k,s_2})}$ depicts the channel dispersion of UE i requesting service v_2 , experiencing PRB k at slice s_2 and $N_{u(v_2,i)}^{k,s_2}$ represents the blocklength of it.

C. 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} .

$$\begin{aligned} T_{tot} &= T_{process} + T_{transmission} + T_{propagation} \\ T_{process} &= T_{RU} + T_{DU} + T_{CU} + T_{UPF} \\ T_{transmission} &= T_{front} + T_{mid} + T_{back} + T_{trans2net} \\ T_{propagation} &= T_{front} + T_{mid} + T_{back} + T_{trans2net} \end{aligned} \quad (6)$$

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). Here we assume the value of propagation delay is negligible compared to the rest.

1) *Processing Delay*: Assume the packet arrival of UEs follows a Poisson process with arrival rate $\lambda_{u(v,i)}$ for the i^{th} UE of the v^{th} service. Therefore, the mean arrival data rate of the v^{th} service in the UPF layer is $\alpha_v^1 = \sum_{u=1}^{U_v} a_{u(v,i)} \lambda_{u(v,i)}$, where $a_{u(v,i)}$ is a binary variable which indicates whether the i^{th} UE requested v^{th} service is admitted or not.

Assume the mean arrival data rate of the UPF layer (α_v^U) is approximately equal to the mean arrival data rate of the O-CU-UP layer (α_v^C) and O-DU (α_v^D). so $\alpha_v = \alpha_v^U \approx \alpha_v^C \approx \alpha_v^D$. since, by using Burkes Theorem, the mean arrival data rate of the second layer which is processed in the first layer is still Poisson with rate α_v . 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 v^{th} service in the first and the second layer, modeled as M/M/1 queue, is formulated as follow, respectively

$$\begin{aligned} T_{DU}^v &= \frac{1}{\mu_d - \alpha_v / M_{v,d}}, \\ T_{CU}^v &= \frac{1}{\mu_c - \alpha_v / M_{v,c}}, \\ T_{UPF}^v &= \frac{1}{\mu_u - \alpha_v / M_{v,u}} \end{aligned} \quad (7)$$

Where $M_{v,d} = \sum_{s=1}^S M_s^d b_{v,s}$, $M_{v,c} = \sum_{s=1}^S M_s^c b_{v,s}$ and $M_{v,u} = \sum_{s=1}^S M_s^u b_{v,s}$ are the sum of VNFs in O-DU, O-CU-UP and UPF, respectively. Since, each service may use more than one slice, the data packets of each service is divided between VNFs of slices allocated to service. Moreover, $1/\mu_d$, $1/\mu_c$ and $1/\mu_u$ are the mean service time of the O-DU, O-CU and the UPF layers respectively. Besides, α_v 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 service v is $\alpha_v / M_{v,i}$ $i \in \{d, c, u\}$.

In addition, T_{RU}^v is the mean transmission delay of v^{th} service on the wireless link. The arrival data rate of wireless link is equal to the arrival data rate of load balancers for each service. Moreover, it is assumed that the service time of transmission queue for each slice s has an exponential distribution with mean $1/(R_{totv})$ and can be modeled as a M/M/1 queue. Therefore, the mean delay of the transmission layer is

$$T_{RU}^v = \frac{1}{R_{totv} - \alpha_v}; \quad (8)$$

where, $R_{totv} = \sum_{u=1}^{U_v} a_{u(v,i)} R_{u(v,i)}$ is the total achievable rate of each service. So the mean processing delay for each UE in service v is

$$T_{process}^v = T_{RU}^v + T_{DU}^v + T_{CU}^v + T_{UPF}^v \quad (9)$$

2) *Transmission Delay*: The transmission delay is the amount of time required to push all the packets into the fiber link. Here, we have transmission delay in fronthaul, midhaul, backhaul and the link to transmit data to internet.

$$\begin{aligned} T_{front} &= \frac{\alpha_v^f}{R_f} \\ T_{mid} &= \frac{\alpha_v^m}{R_m} \\ T_{back} &= \frac{\alpha_v^b}{R_b} \\ T_{trans2net} &= \frac{\alpha_v^t}{R_t} \end{aligned} \quad (10)$$

Where, R_f , R_m , R_b and R_t are the rate of transmission in fronthaul, midhaul, backhaul and the link to transmit data to internet, respectively. Furthermore, the mean arrival data rate of the each link (α_v^i , $i \in \{f, m, b, t\}$) is approximately equal to others ($\alpha_v \approx \alpha_v^i$, $i \in \{f, m, b, t\}$).

D. Physical Data Center Resource

Each VNF requires physical resources that contain memory, storage, CPU and Network Bandwidth. Let the required resources for VNF f in slice s is represented by a tuple as

$$\bar{\Omega}_s^f = \{\Omega_{M,s}^f, \Omega_{S,s}^f, \Omega_{C,s}^f, \Omega_{N,s}^f\}, \quad (11)$$

where $\bar{\Omega}_s^f \in \mathbb{C}^4$ and $\Omega_{M,s}^f, \Omega_{S,s}^f, \Omega_{C,s}^f, \Omega_{N,s}^f$ indicate the amount of required memory, storage, CPU and Network Bandwidth, respectively. Moreover, the total amount of

required memory, storage, CPU and Network Bandwidth of all VNFs of a slice is defined as

$$\bar{\Omega}_{\mathfrak{z},s}^{tot} = \sum_{f=1}^{F_s} \bar{\Omega}_{\mathfrak{z},s}^f \quad \mathfrak{z} \in \{M, S, C, N\}. \quad (12)$$

Where, $F_s = M_s^d + M_s^c + M_s^u$. Also, there are D_c data centers (DC), serving the VNFs. Each DC contains several servers that supply VNF requirements. The amount of memory, storage, CPU and and Network Bandwidth is denoted by τ_{M_j} , τ_{S_j} , τ_{C_j} and τ_{N_j} for the j^{th} DC, respectively

$$\tau_j = \{\tau_{M_j}, \tau_{S_j}, \tau_{C_j}, \tau_{N_j}\},$$

In this system model, the assignment of physical DC resources to VNFs is considered. Let $y_{s,d}$ be a binary variable indicating whether the d^{th} DC is allocated the resources to the VNFs of s^{th} slice or not.

E. Problem Statement

In this system, the goal is to minimize the cost of the system. Power of each O-RU is obtained as below

$$P_r = \sum_{v=1}^V \sum_{i=1}^{U_v} \sum_{s=1}^S \sum_{k=1}^{K_s} p_{r,u(v,i)}^{k,s} a_{u(v,i)} b_{v,s} e_{r,u(v,i)}^{k,s} g_{u(v,i)}^r. \quad (13)$$

The total power cost of O-RUs for transmitting data to UE is depicted as follow

$$P_{tot} = \sum_{r=1}^R P_r \quad (14)$$

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

$$\phi_{tot} = \sum_{s=1}^S \phi_s + \sum_{d=1}^{D_c} z_d \psi_d.$$

Where, z_d is shown that whether the d^{th} DC is turned on or not and ψ_d is a static cost when a DC is active.

$$z_d = \begin{cases} 1 & \sum_{s=1}^S y_{s,d} \geq 1 \\ 0 & \text{otherwise} \end{cases} \quad (15)$$

In addition, $\phi_{s,d}$ is obtained from below

$$\phi_s = M_s^u \phi_s^u + M_s^c \phi_s^c + M_s^d \phi_s^d \quad (16)$$

So the optimization problem is formulated as follow The aim of this paper is minimize the total power of all RUs and the total power consumption of baseband processing at all DCs simultaneously, with the presence of constraints which is written as follow,

$$\min_{P, A, E, M, G, Y} P_{tot} + \phi_{tot} \quad (17a)$$

$$\text{subject to } P_r \leq P_{max} \quad \forall r \quad (17b)$$

$$p_{r,u(v,i)}^{k,s} \geq 0 \quad \forall v, \forall i, \forall r, \forall s, \forall k, \quad (17c)$$

$$\mathcal{R}_{u(v_1,k)}^e \geq \mathcal{R}_{min}^{v_1,e} \quad \forall v_1, \quad (17d)$$

$$\mathcal{R}_{u(v_2,k)}^u \geq \mathcal{R}_{min}^{v_2,u} \quad \forall v_2, \quad (17e)$$

$$T_{tot}^v \leq T_{tot}^{max,v} \quad \forall v, \quad (17f)$$

$$a_{u(v,i)} \leq a_{u(v,i)} \sum_r g_{u(v,i)}^r \leq 1 \quad \forall v, \forall i \quad (17g)$$

$$\sum_{s=1}^S \sum_{v=1}^V \sum_{i=1}^{U_v} a_{u(v,i)} g_{u(v,i)}^r e_{r,u(v,i)}^{k,s} \leq 1 \quad \forall r, \quad (17h)$$

$$\sum_{s=1}^S b_{v,s} \geq 1 \quad \forall v, \quad (17i)$$

$$\sum_{d=1}^{D_c} \sum_{v=1}^V y_{s,d} a_{v,s} \geq 1 \times \sum_{v=1}^V a_{v,s} \quad \forall s, \quad (17j)$$

$$\sum_{s=1}^S y_{s,d} \bar{\Omega}_{\mathfrak{z},s}^{tot} \leq \tau_{\mathfrak{z},d} \quad \forall d, \forall \mathfrak{z} \in \mathcal{E}; \quad (17k)$$