

Joint Power Allocation and Network Slicing in an End-to-End O-RAN System

1st Mojdeh Karbalaee Motalleb

*Electrical and Computer Engineering
Tehran University
Tehran, Iran
mojdeh.karbalaee@ut.ac.ir*

2nd Vahid Shah-Mansouri

*Electrical and Computer Engineering
Tehran University
Tehran, Iran
vmansouri@ut.ac.ir*

3rd Salar Nouri Naghadeh

*Electrical and Computer Engineering
Tehran University
Tehran, Iran
salar.nouri@ut.ac.ir*

Abstract—Many major telecommunication companies confirmed the unification of the xRAN Group with the C-RAN Alliance to establish a more flexible and openness radio access network which is the Open-RAN (O-RAN) for the fifth generation of wireless technology.

To increase energy efficiency and optimize the allocation of resources, Network Slicing (NS) is considered as the best method for the fifth generation (5G) in order to virtualize the common physical network into several logical end-to-end networks. Every slice consists of a part of core network resources, network functions, and radio access network resources as a functional end-to-end network.

In this paper, we elaborate joint NS in RAN and Core of O-RAN system to investigate the power of each User Equipment (UE), map slices to services and also map physical Data Centers (DC) to slices to jointly maximize energy efficiency and minimize consumption power of RUs and the cost of physical resources in a downlink channel. The problem is formulated as a mixed-integer optimization problem that can be decomposed into two independent sub-problems due to the fact that sub-problems are independent. Heuristic algorithms are proposed to each of sub-problems in order to map slices to services, optimize power consumption and map slices to physical resources to minimize the cost of total DCs simultaneously.

Index Terms—O-RAN, Network Slicing, Energy Efficiency, Data Center (DC)

I. INTRODUCTION

O-RAN, as the integration and expansion of C-RAN and xRAN, is expected to be a key technology in 5G Networks to extensively enhance the RAN performance. The core idea of C-RAN is to split the radio remote head (RU) from baseband unit (BBU). Several BBUs operating on a cloud server will create a BBU-Pool, providing unified baseband signal processing with powerful computing capabilities [1]. To communicate between BBU-Pool and RUs, the fronthaul fiber link interface, is assumed [2]–[4]. On the other hand, xRAN technology has three fundamental features. The control plane is decoupled from user plane. Besides, a modular eNB software stack is built to operate on common-off-the-shelf (COTS) hardware. Moreover, open north-bound and south-bound interfaces are introduced [5].

O-RAN virtualizes the elements of the radio access networks, separate them and define appropriate open interfaces for connection of these elements. Moreover, O-RAN uses machine learning techniques to develop smarter RAN layers in

its architecture. In an innovative O-RAN system, the programmable RAN software is decoupled from hardware, which can be run on any specific processing platforms [6].

Openness is the specific feature of O-RAN. Open interface is one of the most crucial properties for O-RAN to enable mobile network operators (MNOs) to define their own services. The concept of software defined network (SDN), which is the separation of control plane (CP), from user plane (UP), is deployed in an intelligent O-RAN architecture. Moreover, this separation promotes RRM to use non-realtime (RT) and near-realtime (NRT) RAN Intelligent Controller (RIC). In the O-RAN architecture, distributed unit (DU) is a logical node hosting RLC/MAC/High-PHY layer. Moreover, central unit (CU) is a logical node hosting RRC, SDAP and PDCP. Also, radio unit (RU), is a logical node hosting LOW-PHY layer and RF processing [7]. O-RAN introduces interfaces such as open fronthaul interface that maps DU and RRU, an E2 interface, and an A1 interface between orchestration/NMS layer containing the non-real-time RIC (RIC non-RT) function and the eNB/gNB containing the near-real-time RIC (RIC near-RT) function.

To evolve servicing in 5G networks, separation of elements of software and hardware of network is employed and is referred to as network functions virtualization (NFV). Virtual network function (VNF) are functional blocks of the system in this case. 5G networks are expected to host several services with different requirements simultaneously. Network slicing is considered as a solution for such demand. A network slice is a logical end-to-end network which provide service with specific requirements. Multiple network slices run on the same physical network infrastructure which are managed and operate independently [8]. Slicing can be defined for the RAN, for the core or both of these networks [9]. In core-RAN Slicing, each slice of RAN is attached to a slice of core.

Based on their type of request, UEs are classified into various services where UEs in the same service has similar service requests. In addition, each service is mapped to one or more slices based on the resource of the slices. Each VNF of the ORAN is mapped to one or some virtual machines (VMs) in the data center. For simplicity of analysis and without loss of generality, we assume a VNF is mapped to one VM and requires specific physical resources including storage, memory,

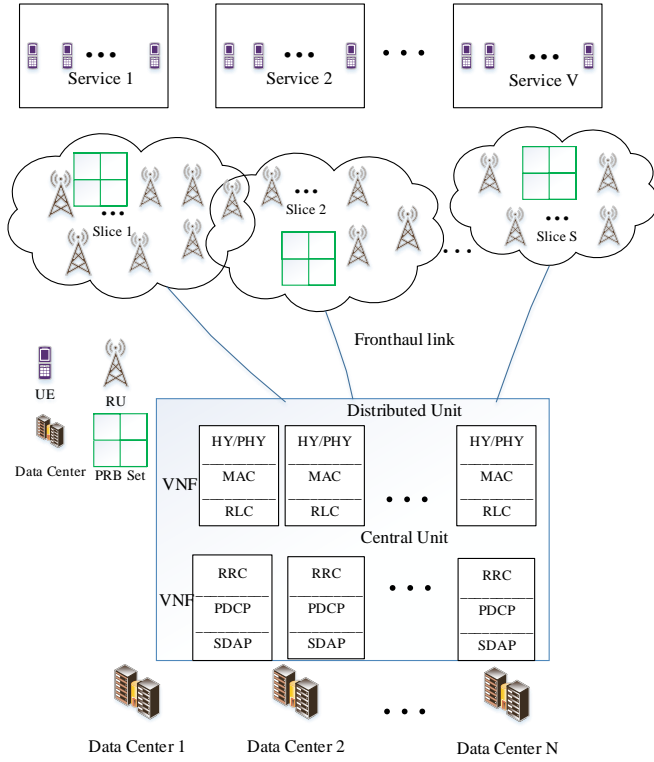


Fig. 1: Network sliced O-RAN system

and processing [2], [10], [11].

In [12], the evolution of RAN in the 5G network is described. The architecture of C-RAN and X-RAN is expressed and development of the O-RAN from C-RAN and X-RAN is explained. Also requirements of the network in terms of capacity and latency is taken account. In [13], the O-RAN architecture, is studied. The eight work groups include focus areas of O-RAN alliance which contains use-cases, overall architecture, open interfaces, and the virtualization and modularization of hardware and software is discussed. In [4], power allocation in C-RAN architecture, with limited fronthaul capacity in clustered network is studied. In [8], dynamic network slicing is considered in Heterogeneous CRAN (H-CRAN) to maximize the weighted sum-rate. Also, two level network slicing (upper and lower level) is done. The admission of UEs is maximized using the heuristic algorithms. In addition, the allocation of BBU capacity to UEs is considered.

In [2], [10], [11], [14] minimization of cost obtained by power consumption, also cloud processing and limited power for each RU is considered. Furthermore, the processing delay of each VM and wireless transmission delay taken account. In [2], the processing and transmission delay is modeled as a M/M/1 queue. In [14], M/D/1 queue is modeled for transmission delay.

In this paper, as depicted in Figure 1, the downlink of the O-RAN system is studied. The RAN is divided to three layers including RU, DU, and CU with open interfaces. DU and CU are running on general purpose data centers. UEs are divided to different services according to their service requirements.

RAN resources are decoupled to slices to provide requirements of services. Optimal power allocation and joint mapping of slices to services are applied. In addition, mapping slices to physical resources is taken to account. The contribution of the paper is expressed as follow:

- In this paper, joint network slicing and resource allocation is considered in an end to end O-RAN system.
- UEs are categorized into different class of services based on their service requirements. RAN which contains RUs, PRBs, and VNFs, in the two layer of processing, are classified into different group of slices.
- The problem is decomposed into two independent sub-problems.
- Heuristic algorithms are applied for the sub-problems to efficiently obtain the solution.

The rest of the paper is organized as follows. In Section II, system model which contains obtaining achievable rates, processing and transmission delay, and physical data center resources is expressed. Then, problem statement is explained and decomposed into two independent sub-problems. In Section III, heuristic algorithms for the sub-problems is presented. In Section IV, numerical results are provided for the problem.

II. SYSTEM MODEL AND PROBLEM FORMULATION

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

A. System Model

Suppose there are S slices Serving V services. Each Service $v \in \{1, 2, \dots, V\}$, consists of U_v single-antenna user equipments (UEs) that require certain service. Each slice $s \in \{1, 2, \dots, S\}$ consists of R_s RUs and K_s Physical Resource Blocks (PRBs). All RUs in a slice that is mapped to a service, transmit signals cooperatively to all the UEs in specific service [4], [15]. Each RU $r \in \{1, 2, \dots, R\}$ is mapped to DU via an optical fiber link with limited fronthaul capacity. Suppose we have two processing layers in the DU/CU of O-RAN system which contain VNFs. The lower layer consist of high-PHY, MAC, and RLC, and the upper layer consist of RRC, PDCP and SDAP. Assume we have M_1 VNFs in the first layer and M_2 VNFs in the second layer for processing data. Each VNF in both layers connects to one or more slices. So in s^{th} slice, there are $M_{s,1}$ VNFs in the first layer and $M_{s,2}$ VNFs in the second layer. All VNFs in first and second layer has the computational capacity that is equal to μ_1 and μ_2 , respectively. Also, RUs and PRBs can serve more than one slices.

B. The Achievable Rate

The achievable data rate for the i^{th} UE in the v^{th} service can be written as

$$\mathcal{R}_{u(v,i)} = B \log_2(1 + \rho_{u(v,i)}), \quad (1)$$

where B is the bandwidth of system and $\rho_{u(v,i)}$ is the SNR of i^{th} UE in v^{th} service which is obtained from

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

where $p_{u(v,i)}$ represents the transmission power allocated by RUs to i^{th} UE in v^{th} service, and $\mathbf{h}_{\mathcal{R}_s, u(v,i)} \in \mathbb{C}^{R_s}$ is the vector of channel gain of a wireless link from RUs in the s^{th} slice to the i^{th} UE in v^{th} service. In addition, $\mathbf{w}_{\mathcal{R}_s, u(v,i)} \in \mathbb{C}^{R_s}$ depicts the transmit beamforming vector from RUs in the s^{th} slice to the i^{th} UE in v^{th} service. Moreover, BN_0 denotes the power of Gaussian additive noise, and $I_{u(v,i)}$ is the power of interfering signals. Moreover, $a_{v,s} \in \{0, 1\}$ is a binary variable that illustrates whether slice s is mapped to service v or not. If $a_{v,s} = 1$ then, v^{th} service is mapped to s^{th} slice; otherwise, it is not mapped.

To obtain SNR as formulated in (2), let $\mathbf{y}_{U_v} \in \mathbb{C}^{U_v}$ be the received signal's vector of all users in v^{th} service

$$\mathbf{y}_{U_v} = \sum_{s=1}^S \sum_{k=1}^{K_s} \mathbf{H}_{\mathcal{R}_s, \mathcal{U}_v}^H \boldsymbol{\eta}_{\mathcal{R}_s} \zeta_{U_v, k, s} a_{v,s} + \mathbf{z}_{U_v}, \quad (3)$$

where, $\boldsymbol{\eta}_{\mathcal{R}_s} = \mathbf{W}_{\mathcal{R}_s, \mathcal{U}_v} \mathbf{P}_{U_v}^{\frac{1}{2}} \mathbf{x}_{U_v} + \mathbf{q}_{\mathcal{R}_s}$. Where $\mathbf{x}_{U_v} = [x_{u(v,1)}, \dots, x_{u(v, U_v)}]^T \in \mathbb{C}^{R_s}$ depicts the transmitted symbol vector of UEs in v^{th} set of service, \mathbf{z}_{U_v} is the additive Gaussian noise $\mathbf{z}_{U_v} \sim \mathcal{N}(0, N_0 \mathbf{I}_{U_v})$ and N_0 is the noise power. In addition, $\mathbf{q}_{\mathcal{R}_s} \in \mathbb{C}^{R_s}$ indicates the quantization noise which, is made from signal compression in DU. Besides, $\mathbf{P}_{U_v} = \text{diag}(p_{u(v,1)}, \dots, p_{u(v, U_v)})$.

Furthermore, $\zeta_{k,s}^{U_v} \triangleq \{\zeta_{k,s}^{u(v,1)}, \zeta_{k,s}^{u(v,2)}, \dots, \zeta_{k,s}^{u(v, N_{U_v})}\}$, $\zeta_{k,s}^{u(v,i)} \in \{0, 1\}$ is a binary parameter, which demonstrates whether i^{th} UE in v^{th} service can transmit its signals through k^{th} PRB and also this PRB belongs to s^{th} slice or not. $\mathbf{H}_{\mathcal{R}_s, \mathcal{U}_v} = [\mathbf{h}_{\mathcal{R}_s, u(v,1)}, \dots, \mathbf{h}_{\mathcal{R}_s, u(v, U_v)}]^T \in \mathbb{C}^{R_s \times U_v}$ shows the channel matrix between RU set \mathcal{R}_s to UE set \mathcal{U}_v , besides. What's more, it is assumed we have perfect channel state information (CSI).

Moreover, $\mathbf{W}_{\mathcal{R}_s, \mathcal{U}_v} = [\mathbf{w}_{\mathcal{R}_s, u(v,1)}, \dots, \mathbf{w}_{\mathcal{R}_s, u(v, U_v)}] \in \mathbb{C}^{R_s \times U_v}$ is the zero forcing beamforming vector to minimize the interference which is indicated as below

$$\mathbf{W}_{\mathcal{R}_s, \mathcal{U}_v} = \mathbf{H}_{\mathcal{R}_s, \mathcal{U}_v} (\mathbf{H}_{\mathcal{R}_s, \mathcal{U}_v}^H \mathbf{H}_{\mathcal{R}_s, \mathcal{U}_v})^{-1}. \quad (4)$$

Hence, the interference power of i^{th} UE in v^{th} service can be represented as follow

$$\begin{aligned} I_{u(v,i)} = & \underbrace{\sum_{s=1}^S \sum_{n=1}^S \sum_{\substack{l=1 \\ l \neq i}}^{U_v} \gamma_1 p_{u(v,l)} a_{v,s} \zeta_{u(v,i), n, s} \zeta_{u(v,l), n, s}}_{\text{(intra-service interference)}} \\ & + \underbrace{\sum_{\substack{y=1 \\ l \neq v}}^V \sum_{s=1}^S \sum_{n=1}^S \sum_{l=1}^{U_y} \gamma_2 p_{u(y,l)} a_{y,s} \zeta_{u(v,i), n, s} \zeta_{u(y,l), n, s}}_{\text{(inter-service interference)}} \\ & + \underbrace{\sum_{s=1}^S \sum_{j=1}^{R_s} \sigma_{q_{r(s,j)}}^2 |\mathbf{h}_{r(s,j), u(v,i)}|^2 a_{v,s}}_{\text{(quantization noise interference)}}. \end{aligned} \quad (5)$$

where, $\gamma_1 = |\mathbf{h}_{\mathcal{R}_s, u(v,i)}^H \mathbf{w}_{\mathcal{R}_s, u(v,l)}|^2$ and $\gamma_2 = |\mathbf{h}_{\mathcal{R}_s, u(v,i)}^H \mathbf{w}_{\mathcal{R}_s, u(y,l)}|^2$. Moreover, $\sigma_{q_{r(s,j)}}$ is the variance of quantization noise of j^{th} RU in s^{th} slice. As it is clear, Interference signal for each UE is coming from UEs using the same PRB. If we replace $p_{u(v,l)}$ and $p_{u(y,l)}$ by P_{max} , an upper bound $\bar{I}_{u(v,i)}$ is obtained for $I_{u(v,i)}$. Therefore, $\bar{\mathcal{R}}_{u(v,i)} \forall v, \forall i$ is derived by using $\bar{I}_{u(v,i)}$ instead of $I_{u(v,i)}$ in (1) and (2).

let $\bar{p}_{r(s,j)}$ denote the power of transmitted signal from j^{th} RU in s^{th} slice. from (3) we have,

$$\bar{p}_{r(s,j)} = \sum_{v=1}^V \mathbf{w}_{r(s,j), \mathcal{U}_v} \mathbf{P}_{U_v}^{\frac{1}{2}} \mathbf{P}_{U_v}^{H \frac{1}{2}} \mathbf{w}_{r(s,j), \mathcal{U}_v}^H a_{v,s} + \sigma_{q_{r(s,j)}}^2. \quad (6)$$

Nevertheless, the rate of users on the fronthaul link between DU and the j^{th} RU in s^{th} slice is formulated as [3], [16]

$$C_{R(s,j)} = \log \left(1 + \sum_{v=1}^V \frac{\mathbf{w}_{r(s,j), \mathcal{D}_s} \mathbf{P}_{U_v}^{\frac{1}{2}} \mathbf{P}_{U_v}^{H \frac{1}{2}} \mathbf{w}_{r(s,j), \mathcal{U}_v}^H a_{v,s}}{\sigma_{q_{r(s,j)}}^2} \right), \quad (7)$$

where, $a_{v,s}$ is a binary variable denotes whether the slice s is mapped to service v or not.

C. Mean Delay

Let the packet arrival of UEs have 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 UEs map to s^{th} slice in the first layer of the O-RAN which contains SDAP, PDCP, and RLC is $\alpha_{s1} = \sum_{v=1}^V \sum_{u=2}^{U_v} a_{v,s} \lambda_{u(v,i)}$, where, $a_{v,s}$ is a binary variable which indicates whether the v^{th} service is mapped to the s^{th} slice or not. Furthermore, the mean arrival data rate of the second layer which contains MAC and HY/PHY is approximately equal to the mean arrival data rate of the first layer $\alpha_s = \alpha_{s1} \approx \alpha_{s2}$ since, by using Burke's Theorem, the mean arrival data rate of the second layer which is processed in the first layer is still Poisson with rate α_s . It is assumed there are dispatchers in each layer for each slice to divide the incoming traffic to VNFs equally. The dispatcher processes the incoming traffic and divides it equally to processed by

VNFs [2], [10], [11]. Suppose the baseband processing of each VNF is depicted as a M/M/1 processing queue. Each packet is processed by one of the VNFs of slices. So the mean delay of the s^{th} slice in the first and the second layer, which is modeled as M/M/1 queue from queuing theory, is formulated as follow ,respectively

$$\begin{aligned} d_{s_1} &= \frac{1}{\mu_1 - \alpha_s/M_{s,1}} \\ d_{s_2} &= \frac{1}{\mu_2 - \alpha_s/M_{s,2}}. \end{aligned} \quad (8)$$

where, $1/\mu_1$ and $1/\mu_2$ are the mean service time of the first and the second layer respectively, due to the fact that the service time of each data packet has an exponential distribution for $\mu_i > \alpha_s/M_{s,i}$ $i \in \{1, 2\}$. Besides, α_s is the data arrival rate which is divided by dispatcher before arriving to VNFs. So the data arrival rate of each VNF in each layer of slice s is $\alpha_s/M_{s,i}$ $i \in \{1, 2\}$. In addition, $d_{s_{tr}}$ is the transmission delay for s^{th} slice as a result of wireless transmission . The arrival data rate of wireless transmission is equal to the arrival data rate of dispatchers (used as a traffic divider) for each slice [2]. Moreover, it is assumed that the service time of transmission queue for each slice s has an exponential distribution with mean $1/(R_{tot_s})$ and can be modeled as a M/M/1 queue [2], [10], [11], [14]. Therefore, the mean delay of the transmission layer is

$$d_{s_{tr}} = \frac{1}{R_{tot_s} - \alpha_s}; \quad (9)$$

where, $R_{tot_s} = \sum_{v=1}^V \sum_{u=2}^{U_v} a_{v,s} R_{u(v,i)}$ is the total achievable rate of each slice that is mapped to specific service. Mean delay of each slice is

$$D_s = d_{s_1} + d_{s_2} + d_{s_{tr}} \forall s. \quad (10)$$

D. Physical Data Center Resource

Each VNF requires physical resources that contain memory, storage and CPU. Let the required resources for VNF f in slice s is represented by a three-dimensional vector as follow

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

where, $\bar{\Omega}_s^f \in \mathbb{C}^3$ and $\Omega_{M,s}^f, \Omega_{S,s}^f, \Omega_{C,s}^f$ indicate the amount of required memory, storage, and CPU, respectively. Moreover the total amount of required memory, storage and CPU of all VNFs of a slice is a three-dimensional vector which is defined as

$$\bar{\Omega}_{3,s}^{tot} = \sum_{f=1}^{M_{s_1}+M_{s_2}} \bar{\Omega}_{3,s}^f \quad 3 \in \{M, S, C\}; \quad (12)$$

Also, there are D_c data centers (DC), serving VNFs. Each DC contains several servers that supply VNF's requirements. The amount of memory, storage and CPU is denoted by τ_{M_j}, τ_{S_j} and τ_{C_j} for j^{th} DC, respectively.

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

In this system model, the placement of physical DC resources to VNFs is considered. $y_{s,d}$ is a binary variable indicates whether d^{th} DC is connected to VNFs of s^{th} slice or not.

E. Problem Statement

One of the most important parameters to measure the optimality of the system is energy efficiency which is represented as a sum-rate to sum-power

$$\eta(\mathbf{P}, \mathbf{A}) := \frac{\sum_{v=1}^V \sum_{k=1}^{U_v} \mathcal{R}_{u(v,k)}}{\sum_{s=1}^S \sum_{i=1}^{R_s} \bar{p}_{r(s,i)}} = \frac{\mathfrak{R}_{tot}(\mathbf{P}, \mathbf{A})}{P_r^{tot}(\mathbf{P}, \mathbf{A})}, \quad (13)$$

where, $P_r^{tot}(\mathbf{P}, \mathbf{A}) = \sum_{s=1}^S \sum_{i=1}^{R_s} \bar{p}_{r(s,i)}$ is the total power consumption of all RUs in all slices. Also, $\mathfrak{R}_{tot}(\mathbf{P}, \mathbf{A}) = \sum_{v=1}^V \sum_{k=1}^{U_v} \mathcal{R}_{u(v,k)}$ is the total rates of all UEs applied for all types of services. 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

$$\phi_{tot} = \sum_{s=1}^S \sum_{d=1}^{D_c} y_{s,d} \phi_{s,d}.$$

Also, a cost function for the placement of VNFs into DCs is defined as

$$\psi_{tot} = \phi_{tot} - \nu \sum_{d=1}^{D_c} \sum_{v=1}^V y_{s,d} a_{v,s} \quad (14)$$

where, ν is a design variable to value between the first term of (14) which is the total power consumption of physical resources and the second term that is shown the amount of admitted slices to have physical resources. Our goal is to maximize sum-rate and minimize sum-power (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,

$$\max_{\mathbf{P}, \mathbf{A}, \mathbf{Y}} \quad \eta(\mathbf{P}, \mathbf{A}) + \frac{1}{\psi_{tot}(\mathbf{Y})} \quad (15a)$$

$$\text{subject to} \quad \bar{p}_{r(s,i)} \leq P_{max} \quad \forall s, \forall i, \quad (15b)$$

$$p_{u(v,k)} \geq 0 \quad \forall v, \forall k, \quad (15c)$$

$$\mathcal{R}_{u(v,k)} \geq \mathcal{R}_{u(v,k)}^{min} \quad \forall v, \forall k, \quad (15d)$$

$$C_{r(s,i)} \leq C_{r(s,i)}^{max} \quad \forall s, \forall i, \quad (15e)$$

$$D_s \leq D_s^{max} \quad \forall s, \quad (15f)$$

$$\sum_{s=1}^S a_{v,s} \geq 1 \quad \forall s, \quad (15g)$$

$$\sum_{d=1}^{D_c} \sum_{v=1}^V y_{s,d} a_{v,s} \geq 1 \times \sum_{v=1}^V a_{v,s} \quad \forall s, \quad (15h)$$

$$\bar{\Omega}_{3,s}^{tot} = \sum_{f=1}^{F_s} \bar{\Omega}_{3,s}^f \leq \sum_{d=1}^{D_c} y_{s,d} \tau_{3d} \quad \forall s, \forall 3 \in \mathcal{E}; \quad (15i)$$

where, $\mathbf{P} = [p_{u(v,k)}] \forall v, \forall k$, is the matrix of power for UEs, $\mathbf{A} = [a_{v,s}] \forall v, \forall s$ denotes the binary variable for connecting

slices to services and $\mathbf{Y} = [y_{s,d}] \forall s, \forall d$ is a binary variable shown whether the physical DC is connected to a VNFs of a slice or not. (15b) and (15c) 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 (15d) shows that the rate of each UE is more than a threshold. (15e) and (15f) expressed the limited capacity of the fronthaul link, and the limited delay of receiving signal, respectively. Furthermore, (15g) ensures that each service is connected to at least one slice. Also, (15h), guarantees that each slice (VNFs in two layers of slices) has been placed to one or more physical resources of DCs. Moreover, in (15i) $\mathcal{E} = \{M, S, C\}$ and the constraint supports that we have enough physical resources for VNFs of each slice.

The optimization problem in (15), can be decomposed into two independent optimization problems A and B since the variables can be obtained independently and respectively. In problem A, the variables \mathbf{P} and \mathbf{A} is obtained. The problem A is as follow

$$\max_{\mathbf{P}, \mathbf{A}} \quad \eta(\mathbf{P}, \mathbf{A}) \quad (16a)$$

$$\text{subject to} \quad \bar{p}_{r(s,i)} \leq P_{max} \quad \forall s, \forall i, \quad (16b)$$

$$p_{u(v,k)} \geq 0 \quad \forall v, \forall k, \quad (16c)$$

$$\mathcal{R}_{u(v,k)} \geq \mathcal{R}_{u(v,k)}^{min} \quad \forall v, \forall k, \quad (16d)$$

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

$$D_s \leq D_s^{max} \quad \forall s, \quad (16f)$$

$$\sum_{s=1}^S a_{v,s} \geq 1 \quad \forall s. \quad (16g)$$

In problem B, \mathbf{Y} is obtained. The problem B is

$$\min_{\mathbf{Y}} \quad \psi_{tot}(\mathbf{Y}) \quad (17a)$$

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

$$\bar{\Omega}_{3,s}^{tot} = \sum_{f=1}^{F_s} \bar{\Omega}_{3,s}^f \leq \sum_{d=1}^{D_c} y_{s,d} \tau_{3d} \quad \forall s, \forall \mathbf{j} \in \mathcal{E}; \quad (17c)$$

III. HUERISTIC METHOD

In this subsection, the proposed method is applied to solve the optimization problem. We want to solve (16). Since the problem is non-convex and NP-Hard, iterative algorithm is applied. To solve the problem and obtain \mathbf{A} and \mathbf{P} we divide problem (16) into two different part that can be solved iteratively. Firstly, in the first part of sub-problem A, we need to obtain \mathbf{A} by fixing $\mathbf{P} = P_{max}$ in the problem (16). Also we fixed $\eta = 0$. Afterward, by achieving \mathbf{A} , in the second part, we just need to acquire \mathbf{P} . This part of the problem can be approximated and converted to a convex problem, so the problem can be solved by convex methods. After solving \mathbf{P} , η is updated. Then in the next iteration, with new \mathbf{P} and η

the two parts of the problems is solved until the algorithm converged.

A. First Part of Sub-Problem A

Two different methods are applied to acquire \mathbf{A} . The details of the heuristic algorithm are represented in algorithm 1.

Algorithm 1 Mapping Slice to Service

- 1: Sort services according to the number of UEs in it and their requirements in the descending order.
 - 2: Sort slices according to the weighted linear combination of number of PRBs, RUs and VNFs in two layers and the Capacity of their resources in the descending order.
 - 3: **for** $i \leftarrow 1$ to S **do**
 - 4: **for** $j \leftarrow 1$ to V **do**
 - 5: Set $a_{i,j} = 1$
 - 6: Obtain Parameters of Systems (power and rate of UEs, rate of fronthaul links, power of RUs)
 - 7: **if** conditions (15b), (15c), (15d) and (15e) is not applied **then**
 - 8: Set $a_{i,j} = 0$;
 - 9: **else**
 - 10: break from inner loop;
 - 11: **end if**
 - 12: **end for**
 - 13: **end for**
-

B. Second Part of Sub-Problem A

In this part, by assuming that \mathbf{A} is fixed, the power of UEs in each service is achieved.

Theorem 1. η^* which is the optimum energy efficiency can be achieved if

$$\begin{aligned} \max_{\mathbf{P}} (\mathfrak{R}_{tot}(\mathbf{P}) - \eta^* P_{r_{tot}}(\mathbf{P})) = \\ \mathfrak{R}_{tot}(\mathbf{P}^*) - \eta^* P_{r_{tot}}(\mathbf{P}^*) = 0. \end{aligned} \quad (18)$$

Proof. See [17, Appendix A] □

The second sub-problem can be solved using the Lagrangian function and iterative algorithm. Since, Interference is a function of the power of UEs, to make it tractable, we assume an upper bound $\bar{I}_{u(v,i)}$ for interference (the worst-case). In order to make (16) as a standard form of a convex optimization problem, it is required to change the variable of equations (16e) and (16f) ($p_{r(s,i)} = \sigma_{q_{r(s,j)}}^2 \times 2^{C_{r(s,i)}}$ and $1/(D_s - d_{s1} + d_{s2}) + \alpha_s$ respectively). Assume λ, μ, ξ , and

κ are the matrix of Lagrangian multipliers that have non-zero positive elements. The Lagrangian function is written as follow

$$\begin{aligned} \mathcal{L}(\mathbf{P}; \lambda, \mu, \xi, \kappa) = & \sum_{v=1}^V \sum_{k=1}^{U_v} \bar{\mathcal{R}}_{u(v,k)} - \eta \sum_{v=1}^V \sum_{i=1}^{R_s} \bar{p}_{r(s,i)} \\ & + \sum_{s=1}^S \sum_{k=1}^{U_v} \lambda_{u(v,k)} (\bar{\mathcal{R}}_{d(s,k)} - \mathcal{R}_{u(v,k)}^{max}) \\ & - \sum_{s=1}^S \sum_{i=1}^{R_s} \mu_{r(s,i)} (\bar{p}_{r(s,i)} - P_{max}) \\ & - \sum_{s=1}^S \sum_{i=1}^{R_s} \xi_{r(s,i)} (\bar{p}_{r(s,i)} - \sigma_{q_{r(s,i)}}^2 2^{C_{r(s,i)}^{max}}). \\ & + \sum_{v=1}^V \sum_{k=1}^{U_v} \kappa_{u(v,k)} \sum_{s=1}^S (R_{u(v,k)} - \mathfrak{D}_s) a_{v,s}. \end{aligned} \quad (19)$$

where, $\mathfrak{D}_s = \frac{1}{D_s^{max} - d_{s1} - d_{s2}} + \alpha_s$. Optimal power is obtained from (19)

$$p_{u(v,i)}^* = \left[\frac{\eta_{u(v,i)} \mathfrak{w}_{u(v,i)} - \mathfrak{r}_{u(v,i)} \mathfrak{z}_{u(v,i)}}{\mathfrak{r}_{u(v,i)} \mathfrak{w}_{u(v,i)}} \right]^+ \quad (20)$$

where, $\eta_{u(v,i)} = (\lambda_{u(v,i)} + \kappa_{u(v,k)} + 1) \frac{B}{Ln_2}$ and $\mathfrak{w}_{u(v,i)} = \sum_{s=1}^S |\mathbf{h}_{R_s, u(v,i)}^H \mathbf{w}_{R_s, u(v,i)}|^2 a_{v,s}$. Also $\mathfrak{z}_{u(v,i)} = BN_0 + \bar{I}_{u(v,i)}$ and $\mathfrak{r}_{u(v,i)} = \sum_{s=1}^S \sum_{i=1}^{R_s} (\mu_{r(s,i)} + \xi_{r(s,i)} + \eta) \|w_{r(s,i), u(v,i)}\|^2$. By using sub-gradient method, the optimal power \mathbf{P} is obtained [15].

C. Solving two part of Sub-problem A iteratively

In (III-A) and (III-B), the details of solving each part of the sub-problem are depicted. Firstly, we obtain \mathbf{A} by fixing $\mathbf{P} = \mathbf{P}_{max}$ in the problem (16) and using algorithm (1). Also we fixed $\eta = 0$. Afterward, by achieving \mathbf{A} , in the second part, we acquire \mathbf{P} using the sub-gradient method. After solving \mathbf{P} , η is updated. Then in the next iteration, with new \mathbf{P} and η the two parts of the problems is solved until the algorithm converged. Here, the algorithm of solving sub-problem A is shown in algorithm (2)

D. Sub-Problem B

In this subsection, we would like to solve (17), which is the placement of virtual resources to physical resources in order to minimize the cost function ψ_{tot} . To achieve optimum \mathbf{Y} heuristic algorithm is applied. The details of the heuristic algorithm are written in algorithm (3). In this algorithm, firstly, we sort slices and DCs according to their sum-weighted of their requirements (line 1 and line 2 of algorithm 3). We define a weighted parameter for $\Omega_{3,s}^{tot}$ and τ_j^3 as follow

$$\begin{aligned} \hat{\Omega}_s^{tot} &= w_M \bar{\Omega}_{M,s}^{tot} + w_S \bar{\Omega}_{S,s}^{tot} + w_C \bar{\Omega}_{C,s}^{tot} \\ \hat{\tau}_j &= w_M \tau_j^M + w_S \tau_j^S + w_C \tau_j^C, \end{aligned} \quad (21)$$

where, $\mathbf{w} = \{w_M, w_S, w_C\}$ is the weight of memory, storage and CPU. Secondly, we start mapping from the most needed

Algorithm 2 Joint Network Slicing and Power Allocation

- 1: Set the maximum number of iterations I_{max} , convergence condition ϵ_η and the initial value $\eta^{(1)} = 0$
- 2: Set $\mathbf{P} = \mathbf{P}_{max}$
- 3: **for** counter $\leftarrow 1$ to I_{max} **do**
- 4: Achieve \mathbf{A} by applying Algorithm (1)
- 5: Obtain \mathbf{P} by using sub-gradient method which is mentioned in (III-B).
- 6: **if** $\Re_{tot}(\mathbf{P}^{(i)}, \mathbf{A}^{(i)}) - \eta^{(i)} P_{r_{tot}}(\mathbf{P}^{(i)}, \mathbf{A}^{(i)}) < \epsilon_\eta$ **then**
- 7: Set $\mathbf{P}^* = \mathbf{P}^{(i)}$, $\mathbf{A}^* = \mathbf{A}^{(i)}$ and $\eta^* = \eta^{(i)}$;
- 8: **break**;
- 9: **else**
- 10: $i = i + 1$, Setting $\mathbf{P} = \mathbf{P}^{(i)}$;
- 11: **end if**
- 12: **end for**

TABLE I: Simulation Parameter

Parameter	Value
Noise power	-174dBm
Bandwidth	120 KHZ
Maximum transmit Power of each RU	40dBm
Maximum transmit Power of each UE	30dBm
Minimum delay	300usec
Maximum fronthaul capacity	200bits/sec/Hz
Minimum data rate	10 bits/sec/Hz

slices to the DC with the most physical resources (from line 4 to line 11 of algorithm 3). After mapping DCs to slices, if some slices are not admitted, we start admitting remaining slices to more than one DCs (from line 12 to line 25 of algorithm 3). At the end, if DC with the lowest physical resources is free and can be served instead of DC with the highest physical resource, the slices remapped to new DC with the lowest physical resource since it has the lowest power consumption (line 26 of algorithm 3). If a slice does not admit to a specific DC, it remains for the next placement. In next placement the remaining slices, map to more than one DC according to their requirements.

IV. NUMERICAL RESULTS

In this section, simulation and numerical results for the main problem are depicted. In Fig. 2, energy efficiency is depicted for two different numbers of services with the different mean number of UEs in each service (the parameters for simulation listed in table I). The optimal method is obtained by using MOSEK toolbox to obtain \mathbf{A} and CVX toolbox to obtain \mathbf{P} and iteratively updates \mathbf{A} , η and \mathbf{P} . The optimal method is 0.1 bits/J/Hz better than the proposed method for $V = 6$ and $E[U_v] = 10$, and also, 0.09 bits/J/Hz better than the proposed method for $V = 3$ and, $E[U_v] = 10$. As it is shown, the Energy Efficiency is increased as a mean number of UEs rises. In Fig. 3, the ratio of admitted slices is demonstrated for two different numbers of DCs with the different number of slices (the parameters for simulation listed in table II and also we

Algorithm 3 Plection of Physical resources into Virtual resources

```

1: Sort Slices according to  $\hat{\Omega}_s^{tot}, \forall s$  in descending order.
2: Sort DCs according to  $\hat{\tau}_j, \forall j$  in descending order.
3:  $\mathbf{Y} = \mathbf{0}$ 
4: for  $d \leftarrow 1$  to  $D_c$  do
5:   for  $s \leftarrow 1$  to  $S$  do
6:     if  $\sum_{d=1}^{D_c} y_{s,d} == 0$  and  $\bar{\Omega}_{\mathbf{z}(s)}^{tot} \leq \tau_{\mathbf{z}(s)} \forall \mathbf{z}, \forall s$  then
7:       Set  $y_{s,d} = 1$ ;
8:        $\tau_j^{\mathbf{z}} \leftarrow \tau_j^{\mathbf{z}} - \bar{\Omega}_{\mathbf{z}(s)}^{tot} \quad \mathbf{z} \in \{M, S, C\}$ 
9:     end if
10:   end for
11: end for
12:  $ind_{rem} = \{s | (\sum_{d=1}^{D_c} y_{s,d} == 0)\}$ 
13: Sort remaining amount of DCs same as before in descending order.
14: Sort remaining slices same as before in descending order.
15: for  $r \leftarrow 1$  to  $S_{rem}$  do
16:   for  $n \leftarrow 1$  to  $D_c$  do
17:     Set  $y_{s,d} = 1$ ;
18:      $\bar{\Omega}_{\mathbf{z}(s)}^{tot} \leftarrow \bar{\Omega}_{\mathbf{z}(s)}^{tot} - \tau_j^{\mathbf{z}}$ 
19:     if  $\bar{\Omega}_{\mathbf{z}(s)}^{tot} == 0$  then
20:       Set  $y_{s,d} = 1$ ;
21:        $\tau_j^{\mathbf{z}} \leftarrow \tau_j^{\mathbf{z}} - \bar{\Omega}_{\mathbf{z}(s)}^{tot} \quad \mathbf{z} \in \{M, S, C\}$ 
22:       break inner loop
23:     end if
24:   end for
25: end for
26: Remapping DCs must be done to prevent wasting Energy

```

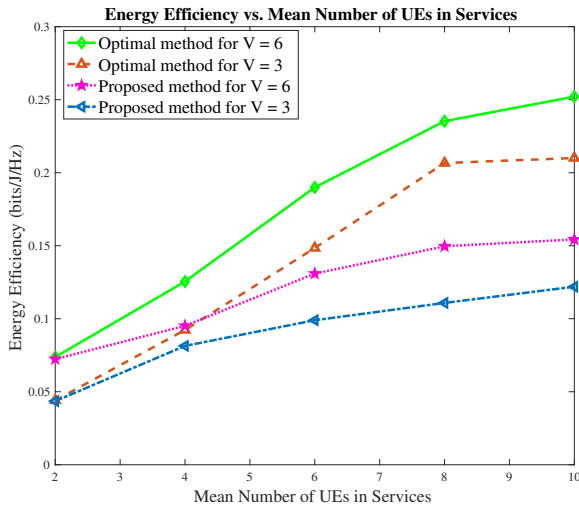


Fig. 2: Energy Efficiency vs. Mean Number of UEs in each Service

TABLE II: Simulation Parameter

Parameter	Value
Mean of CPU for DCs	320GHz
Mean of Memory for DCs	1T
Mean of Storage for DCs	100T
Mean of CPU for Slices	32GHz
Mean of Memory for Slices	100G
Mean of Storage for Slices	10T

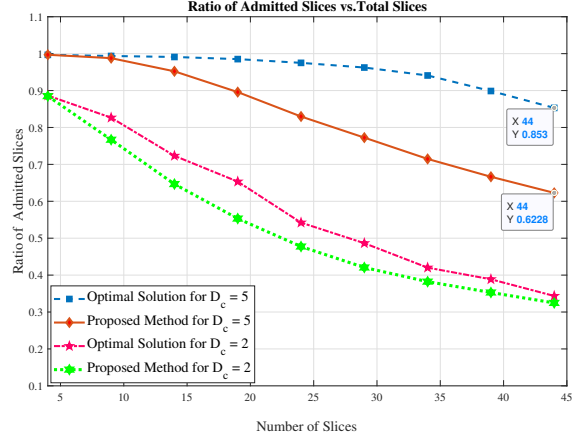


Fig. 3: Ratio of Admitted Slices connected to just one DC vs. Total slices

set $w_C = 320$, $w_S = 100$, $w_M = 1$). In this simulation, it is assumed that just one DC can serve each slice, and it is not admitted by more than one DC. The proposed method is based on Algorithm 3, and the optimal method is done by the MOSEK toolbox. When we have two DCs, the proposed method and optimal method have approximately the same ratio of admitted slices. But by increasing the number of DCs to five, the performance of the proposed method reduced. Using five DCs, the difference between the proposed method and the optimal method in the worst case (44 slices) is about 23 percentage. In Fig. 4, the normalized resource consumption is depicted due to the number of slices (the parameters for simulation listed in table II and also we set $w_C = 320$, $w_S = 100$, $w_M = 1$). In this simulation, it is assumed that the number of DCs is entirely enough to cover all slices. The optimality of the placement of DCs to slices is measured. It is shown that how much resources of active DCs are not used. For ten slices, the difference between the optimal solution and the proposed solution is about 15 percent.

V. CONCLUSION

In this paper, joint network slicing and power allocation is considered in an end to end O-RAN system. It is assumed that UEs are classified based on their service requirements. Also, there is a number of slices served the services. Each slice, is consist of PRBs, RUs, and VNFs that run on VMs. The limited fronthaul capacity is considered for the fiber links between RUs and DU The target is to maximize the sum-rate

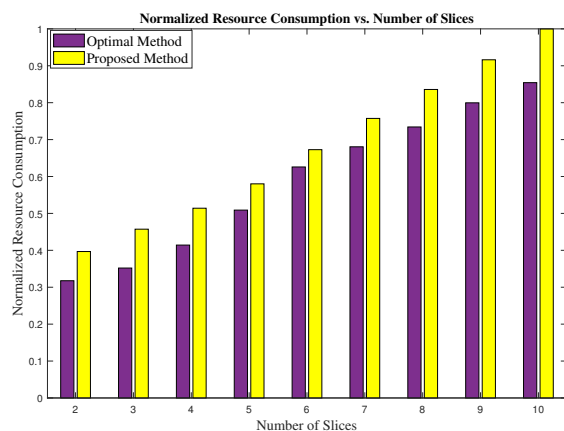


Fig. 4: Normalized Resource Consumption vs. Number of Slices

and minimize the power consumption and energy cost of data centers simultaneously. The problem is decomposed into two sub-problems. Each sub-problems are solved separately by a heuristic algorithm. For the sub-problem A, energy efficiency vs. mean number of UEs in each service is depicted. The energy efficiency exceeded by increasing the mean number of UEs in each service. For the sub-problem B, two figures are shown. In the first one, the ratio of admitted slices that connect to just one DC, for a different number of slices is denoted. In the second figure, the normalized resource consumption of DCs is depicted. In each figure, the heuristic algorithm is compared with the optimal method and the difference between them is discussed.

REFERENCES

- [1] J. Wu, Z. Zhang, Y. Hong, and Y. Wen, "Cloud radio access network (c-ran): a primer," *IEEE Network*, vol. 29, no. 1, pp. 35–41, 2015.
- [2] 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.
- [3] O. Simeone, J. Kang, J. Kang, and S. Shamai, "Cloud radio access networks: Uplink channel estimation and downlink precoding," *arXiv preprint arXiv:1608.07358*, 2016.
- [4] M. K. Motalleb, A. Kabiri, and M. J. Emadi, "Optimal power allocation for distributed mimo c-ran system with limited fronthaul capacity," in *2017 Iranian Conference on Electrical Engineering (ICEE)*. IEEE, 2017, pp. 1978–1982.
- [5] (2018) xran forum merges with c-ran alliance to form oran alliance. [Online]. Available: <https://www.businesswire.com/news/home/20180227005673/en/xRAN-Forum-Merges-C-RAN-Alliance-Form-ORAN>
- [6] B. Fletcher. (2019) Vodafone initiates first open ran trials in the u.k., challenging traditional vendors. [Online]. Available: <https://www.fiercewireless.com/tech/vodafone-initiates-first-open-ran-trials-uk-challenging-traditional-vendors>
- [7] I. Scales. (2018) The open ran (oran) alliance formed to lever open 5g for 'other' technologies? and much more... [Online]. Available: <https://www.o-ran.org/resources>
- [8] 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.
- [9] X. Zhou, R. Li, T. Chen, and H. Zhang, "Network slicing as a service: enabling enterprises' own software-defined cellular networks," *IEEE Communications Magazine*, vol. 54, no. 7, pp. 146–153, 2016.

- [10] 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.
- [11] 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.
- [12] P. Sehier, P. Chanclou, N. Benzaoui, D. Chen, K. Kettunen, M. Lemke, Y. Pointurier, and P. Dom, "Transport evolution for the ran of the future," *IEEE/OSA Journal of Optical Communications and Networking*, vol. 11, no. 4, pp. B97–B108, 2019.
- [13] S. A. T. Kawahara and A. U. R. Matsukawa, "O-ran alliance standardization trends," 2019.
- [14] 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.
- [15] P.-R. Li, T.-S. Chang, and K.-T. Feng, "Energy-efficient power allocation for distributed large-scale mimo cloud radio access networks," in *2014 IEEE Wireless Communications and Networking Conference (WCNC)*. IEEE, 2014, pp. 1856–1861.
- [16] S.-H. Park, O. Simeone, O. Sahin, and S. S. Shitz, "Fronthaul compression for cloud radio access networks: Signal processing advances inspired by network information theory," *IEEE Signal Processing Magazine*, vol. 31, no. 6, pp. 69–79, 2014.
- [17] D. W. K. Ng, E. S. Lo, and R. Schober, "Energy-efficient resource allocation for secure ofdma systems," *IEEE Transactions on Vehicular Technology*, vol. 61, no. 6, pp. 2572–2585, 2012.