# Cross-Layer Energy Efficient Resource Allocation in PD-NOMA based H-CRANs: Implementation via GPU

Ali Mokdad, Paeiz Azmi, *Senior Member, IEEE*, Nader Mokari, *Member, IEEE*, Mohammad Moltafet, and Mohsen Ghaffari-Miab, *Member, IEEE*

*Abstract*—In this paper, we propose a cross layer energy efficient resource allocation and remote radio head (RRH) selection algorithm for heterogeneous traffic in power domain - non-orthogonal multiple access (PD-NOMA) based heterogeneous cloud radio access networks (H-CRANs). The main aim is to maximize the EE of the elastic users subject to the average delay constraint of the streaming users and the constraints, RRH selection, subcarrier, transmit power and successive interference cancellation. The considered optimization problem is non-convex, NP-hard and intractable. To solve this problem, we transform the fractional objective function into a subtractive form. Then, we utilize successive convex approximation approach. Moreover, in order to increase the processing speed, we introduce a framework for accelerating the successive convex approximation for low complexity with the Lagrangian method on graphics processing unit. Furthermore, in order to show the optimality gap of the proposed successive convex approximation approach, we solve the proposed optimization problem by applying an optimal method based on the monotonic optimization. Studying different scenarios show that by using both PD-NOMA technique and H-CRAN, the system energy efficiency is improved.

*Index Terms*—Heterogeneous traffic, PD-NOMA, remote radio head selection, graphics processing unit.

## I. INTRODUCTION

### A. State of the Art

IN next cellular communication systems, power domain - non-orthogonal multiple access (PD-NOMA) is a novel multiple access scheme which is a promising candidate for the fifth generation (5G) cellular communication systems [1]. PD-NOMA multiplexes different users symbols by applying the superposition coding (SC) technique at the transmitter side, while at the receiver side the successive interference cancellation (SIC) technique is applied to recover back the multiplexed symbols [1].

Heterogeneous cloud radio access network (H-CRAN) is a novel architecture which is proposed as a promising technology for next cellular communication systems [2]. H-CRAN combines heterogeneous cellular network (HCN) with cloud radio access network (C-RAN). In addition, H-CRAN covers the advantages of C-RAN and HCN at the same time [2]. The main subsystems of the H-CRAN architecture are the baseband unit (BBU) pool, fiber links and remote radio heads

(RRHs) where one of the RRHs is a high power node (HPN) and the others are low power nodes (LPNs). Instead of the processing that is distributed at the base stations (BSs) in the HCN, a centralized signal processing is applied in the BBU pool which reduces the manufacturing and operating cost. Moreover, a cooperation between different RRHs is permitted due to the centralized signal processing, thus spectrum efficiency and link reliability are improved. The RRHs compress and forward the received signals from the user to the BBU pool via high bandwidth and low latency fiber links [2]. Therefore, H-CRANs improve the users quality of service (QoS), the spectral efficiency (SE) of the system and increase the network architecture flexibility. Moreover, H-CRANs decrease the power consumption of the system, and PD-NOMA technique improves the system throughput, SE, and energy efficiency (EE) of the fifth generation (5G) cellular communication systems. In order to cover the advantages of H-CRAN and PD-NOMA technique at the same time, we consider PD-NOMA based H-CRAN system.

Due to the enormous increase in mobile data traffic and the complexity of the proposed technologies including PD-NOMA and H-CRAN, a high computational processing is needed where the conventional methods can not tackle this issue. Therefore, we seek toward a new processing method which accelerates the processing time. Graphics Processing Unit (GPU), due to the advantage of its massive number of cores and its parallelism directives, handles the works with parallel data [3]–[7]. Accelerating applications and simulations with using GPUs has turned out to be progressively well-known from 2006 [8]. OpenACC is an open GPU directives standard which makes GPU programming simple and portable over the parallel multi-core processors [3]. In [9], a communication optimization for multi GPU implementation of Smith-Waterman Algorithm is investigated. In [6], stochastic finite-difference time domain method is investigated on GPU by employing OpenACC application program interface (API).

### B. Related Works

During the past decade, numerous energy efficient (EE), BS selection and cross layer resource allocation problems for OFDMA systems are investigated [2], [10]–[15]. Furthermore, different PD-NOMA systems are studied [16]–[20].

In [10], the EE orthogonal frequency division multiplexing (OFDM) relay system is developed where both the transmit

The authors are with the Department of Electrical and Computer Engineering, Tarbiat Modares University, Tehran, Iran (ali.mokdad@modares.ac.ir, pazmi@modares.ac.ir, nader.nmy@gmail.com, m.moltafet@modares.ac.ir, mghaffari@modares.ac.ir).

and circuit power consumptions are considered. The EE power allocation for OFDM based cognitive radio networks is investigated in [11]. BS or cell selection for the mobile user is investigated in [13], [14]. In [15], a cross layer resource allocation scheme for OFDMA systems is investigated. In [2], the EE resource allocation in H-CRANs is studied, where RRHs are basically utilized to supply high data rates for users with high quality of service (QoS) requirements, while HPN is created to ensure the coverage and serve users with low QoS requirements. In [2], the number of RRHs is supposed to be sufficiently large, then the considered overall EE optimization problem of the H-CRAN system is approximated to EE optimization problem for only one RRH.

In [21], a comprehensive overview of the latest NOMA research and innovations as well as their applications are summarized and discussed. In [16], the effect of user pairing on the performance of PD-NOMA systems is investigated. A power allocation in OFDM-NOMA system is studied in [17], where a single BS is taken into consideration. In [18], joint power and channel allocation for PD-NOMA in 5G downlink cellular systems by considering one BS is developed. In [19], [20], the radio resource allocation for HCNs based on PD-NOMA is studied. In [22], robust radio resource allocation for a cellular system based on PD-NOMA is investigated.

To the best of our knowledge, cross layer resource allocation and RRH selection problems neither for systems based on PD-NOMA technique nor for H-CRAN have been investigated yet. As well, resource allocation for H-CRAN systems neither based on PD-NOMA nor with heterogeneous traffic have been studied so far. Moreover, successive convex approximation for low complexity (SCALE) [23] with the Lagrangian method has not been analyzed on GPU using OpenACC API yet.

## C. Contributions

In our work, we consider a cross layer EE radio resource allocation and RRH selection problem for heterogeneous traffic in PD-NOMA based H-CRANs. In this formulation, two types of traffic are taken into account, elastic traffic and streaming traffic. In our design, first, the radio resources are assigned to the streaming traffic users in a way that the streaming users QoS constraints are satisfied. Thereafter, the remaining radio resources are assigned to the elastic traffic users. The optimization problem is to maximize the energy efficiency of the elastic users where the total power consumption is partitioned to three parts: 1) the power consumption in the fiber links depending on the active RRHs, 2) the power consumption of RRHs and 3) the circuit power consumption [2]. Moreover, due to utilizing the PD-NOMA technique more than one user can be allocated at the same subcarrier and each user can be served by only one RRH. The considered EE optimization problem is non-convex, intractable, and NP-hard. Therefore, we solve the considered optimization problem by applying the successive convex approximation (SCA) method. Therefore, in our paper, we focus on both resource allocation and remote radio head selection. Then, due to the different factors taken into account which are from power allocation, subcarrier allocation and remote radio head selection, and at the same

time the enormous increase in mobile data traffic, a high computational processing is needed where the conventional methods can not tackle this issue. Moreover, increasing the number of variables in the system which means increasing the number of parameters is beneficial since it makes the system more flexible in allocating the energy efficiency which helps in maximizing the energy efficiency of the system. Thus, to accelerate the processing speed, we introduce a framework for SCALE with the Lagrangian method on GPU and we run the proposed optimization problem on GPU by utilizing OpenACC API. Moreover, in order to evaluate the optimality gap of the proposed solution, we solve the considered optimization problem by applying an optimal algorithm based on the monotonic optimization [24]–[26]. Simulation results confirm that the energy efficiency performance of the H-CRAN based on the PD-NOMA method is approximately 14% more than the systems based on orthogonal multiple access (OMA) where only one user can be selected on each subcarrier. Moreover, simulation results show that the system energy efficiency in H-CRAN scenario is enhanced compared to the conventional, C-RAN, HCN and 1-tier HPN scenarios.

The key contributions of this paper are summarized as follows:

- We propose a cross layer EE radio resource allocation and RRH selection algorithm for heterogeneous traffic in PD-NOMA based H-CRANs.
- We prove the convergence of the SCA approach for the cross layer EE radio resource allocation and RRH selection in PD-NOMA based H-CRANs and we highlight on the performance improvements of the NOMA technique.
- We solve the considered optimization problem by applying the monotonic optimization method. First, we transform the optimization problem to a monotonic one in a canonical form, then we obtain the solution by applying the polyblock algorithm.
- We introduce a framework for accelerating SCALE with the Lagrangian method on GPU and we run the proposed optimization problem by using OpenACC API on GPU.

## D. Paper Organization

The reminder of this paper is organized as follows. In Section II, we describe the system model and problem formulation of our design. The transformation of the fractional objective function problem to a problem with an objective function with subtractive form is introduced in Section III. The proposed approaches to solve the equivalent cross layer EE resource allocation and RRH selection problem are presented in Section IV. Computational complexity of the proposed solution methods are studied in Section V. Distributed solution and signalling overhead of both the centralized and distributed solutions are investigated in Section VI. A framework for accelerating the general SCALE with the Lagrangian method using GPU is proposed in Section VII. The performance of the proposed algorithm and our system model through different numerical experiments are examined in Section VIII. Lastly, we conclude the paper in Section IX.

Fig. 1: A two-tier H-CRAN consisting of one HPN RRH and set of LPN RRHs.

## II. SYSTEM MODEL AND PROBLEM FORMULATION

### A. System Model

We consider a two tier downlink H-CRAN, where a typical illustration example of this network is presented in Fig. 1. As well, the proposed cross layer with RRH selection system in PD-NOMA H-CRANs is shown in Fig. 2. In this network, $M_f$ LPN RRHs and one HPN RRH cover the desired coverage area sharing the available radio spectrum. Table I summarizes the parameters and symbols used in the system model and problem formulation.

The RRHs set is denoted by $\mathcal{M} = \{0, 1, 2, ..., M_f\}$, where 0 is the index of the HPN RRH and $\mathcal{M}_f = \{1, 2, ..., M_f\}$ is the set of the LPN RRHs. $M = M_f + 1$ is the number of all RRHs. We denote the set of all users by $\mathcal{K} = \{1, 2, ..., K\}$. The users set is split into two sets: 1) streaming users set $\mathcal{K}^s = \{1, 2, ..., K^s\}$ and 2) elastic users set $\mathcal{K}^e = \{1, 2, ..., K^e\}$. The number of streaming users and elastic users are equal to $K^s = |\mathcal{K}^s|$ and $K^e = |\mathcal{K}^e|$, respectively. Therefore, $\mathcal{K} = \mathcal{K}^e \cup \mathcal{K}^s$ and the number of all users is $K = K^s + K^e$. Due to the PD-NOMA technique, over each subcarrier in RRH $m$, $l$ users can be allocated where $l \leq K$. In this system model, we suppose the system bandwidth is equal to $B$ partitioned to $N$ subcarriers with bandwidth $B_n = B/N$ and the subcarriers set is denoted by $\mathcal{N} = \{1, 2, ..., N\}$. $h_{m,k}^{(n)}$ denotes the channel gain from RRH $m$ to user $k$ over subcarrier $n$ and $\Gamma_{m,k}^{(n)} = |h_{m,k}^{(n)}|^2$. Due to using the PD-NOMA technique, signals of users with better channel condition is considered as noise while the signals of users with weaker channel condition can be successfully decoded and removed during the decoding process [27]–[29]. Then, the RRH $m$ transmits $\sum_{k \in \mathcal{K}} A_{m,k}^{(n)} \rho_{m,k}^{(n)} \sqrt{p_{m,k}^{(n)}} s_{m,k}^{(n)}$ over subcarrier $n$ where $s_{m,k}^{(n)}$ is the information signal for the $k^{th}$ user from RRH $m$ over subcarrier $n$, $p_{m,k}^{(n)}$ represents the transmit power from RRH $m$ to user $k$ over subcarrier $n$ and $\rho_{m,k}^{(n)}$ is a binary variable as user and subcarrier allocation indicator where $\rho_{m,k}^{(n)} = 1$ if

TABLE I: Table of symbols used in the system model.

| Symbol | Definition / Description |
|---|---|
| $M_f$ | Number of LPN RRHs |
| $\mathcal{M} = \{0, 1, 2, ..., M_f\}$ | RRHs set |
| $\mathcal{M}_f = \{1, 2, ..., M_f\}$ | LPN RRHs set |
| $M$ | Number of all RRHs |
| $\mathcal{K} = \{1, 2, ..., K\}$ | Users set |
| $\mathcal{K}^s = \{1, 2, ..., K^s\}$ | Streaming users set |
| $\mathcal{K}^e = \{1, 2, ..., K^e\}$ | Elastic users set |
| $K^s$ | Number of streaming users |
| $K^e$ | Number of elastic users |
| $K$ | Number of all users |
| $l$ | Number of users that can be allocated on each subcarrier |
| $B$ | System bandwidth |
| $N$ | Number of subcarriers |
| $B_n$ | Subcarrier bandwidth |
| $\mathcal{N} = \{1, 2, ..., N\}$ | Subcarriers set |
| $h_{m,k}^{(n)}$ | Channel gain from RRH $m$ to user $k$ over subcarrier $n$ |
| $s_{m,k}^{(n)}$ | Information signal for the $k^{th}$ user |
| $p_{m,k}^{(n)}$ | Transmit power from RRH $m$ to user $k$ over subcarrier $n$ |
| $\rho_{m,k}^{(n)}$ | User and subcarrier allocation indicator |
| $A_{m,k}$ | User and RRH allocation indicator |
| $\gamma_{m,k}^{(n)}$ | SINR of user $k$ on subcarrier $n$ in RRH $m$ |
| $\sigma_{m,k}^{(n)}$ | Noise power at user $k$ in RRH $m$ over subcarrier $n$ |
| $I_{m,k}^{(n)}$ | Received interference power from the multiplexed users and other RRHs |
| $r_{m,k}^{(n)}$ | Rate of user $k$ over subcarrier $n$ in RRH $m$ |
| $r_k$ | Full achievable rate of the user $k$ |
| $w_{m,k}$ | Priority weight of the user $k$ in RRH $m$ |
| $R$ | Total weighted sum rate of the elastic users |
| $P_f^L$ | LPN RRH fiber link power consumption |
| $P_f^H$ | HPN RRH fiber link power consumption |
| $\eta_m$ | Efficiency of the power amplifier in RRH $m$ |
| $P_c^L$ | LPN RRH circuit power consumption |
| $P_c^H$ | HPN RRH circuit power consumption |
| $P$ | Total power consumption of the elastic users |
| $E$ | Overall energy efficiency for the H-CRAN |
| $\lambda_k$ | Arrival rate |
| $T_k$ | Desired maximum delay requirement |
| $q_k$ | Average queue length |
| $p_m^{max}$ | RRH $m$ maximum allowable transmit power |
| $p_{m,k}^{(n),mask}$ | Transmit power spectral mask for user $k$ |
| $\overline{X_k}$ | Average time that user $k$ waits in the queue in addition to the service time |
| $\overline{X_k^2}$ | Second moment of the service time |
| $\overline{z}$ | Packet size |
| $\varrho_1, \varrho_2, \xi, \varpi_1$ and $\varpi_2$ | small positive numbers |
| $i$ | Index of the iterative algorithm |
| $\boldsymbol{\xi'}, \boldsymbol{\zeta'}, \boldsymbol{\vartheta}, \boldsymbol{\vartheta'}$ and $\boldsymbol{\tilde{\zeta}'}$ | Lagrangian multipliers vectors |
| $P_c^M$ | MBS static circuit power consumption |
| $P_c^P$ | PBS static circuit power consumption |
| $\eta_0$ | Power efficiency for each MBS or PBS |

user $k$ is allocated over the subcarrier $n$ in RRH $m$ and equal to zero otherwise.

As well, $A_{m,k}$ is a binary variable as user and RRH allocation indicator where $A_{m,k} = 1$ if user $k$ is served by RRH $m$ and equal to zero otherwise. Then we denote $\boldsymbol{\rho}_{m,k} = [\rho_{m,k}^{(1)}, \rho_{m,k}^{(2)}, ..., \rho_{m,k}^{(N)}]$, $\boldsymbol{\rho}_m = [\boldsymbol{\rho}_{m,1}, \boldsymbol{\rho}_{m,2}, ..., \boldsymbol{\rho}_{m,K}]$ and $\boldsymbol{\rho} = [\boldsymbol{\rho}_0, \boldsymbol{\rho}_1, ..., \boldsymbol{\rho}_{M_f}]$. Moreover, we denote $\mathbf{p}_m^{(n)} =$

Fig. 2: Scheduling model for the H-CRAN system.

$[p_{m,0}^{(n)}, p_{m,1}^{(n)}, ..., p_{m,K}^{(n)}]$, $\mathbf{p}^{(n)} = [\mathbf{p}_0^{(n)}, \mathbf{p}_1^{(n)}, ..., \mathbf{p}_{M_f}^{(n)}]$, $\mathbf{p}_{m,k} = [p_{m,k}^{(1)}, p_{m,k}^{(2)}, ..., p_{m,k}^{(N)}]$, $\mathbf{p}_m = [\mathbf{p}_{m,1}, \mathbf{p}_{m,2}, ..., \mathbf{p}_{m,K}]$ and $\mathbf{p} = [\mathbf{p}_0, \mathbf{p}_1, ..., \mathbf{p}_{M_f}]$.

As such, the signal to interference plus noise ratio (SINR) of user $k$ over subcarrier $n$ in RRH $m$ after performing SIC is $\gamma_{m,k}^{(n)} = \frac{p_{m,k}^{(n)} \Gamma_{m,k}^{(n)}}{\sigma_{m,k}^{(n)} + I_{m,k}^{(n)}}$ where $\sigma_{m,k}^{(n)}$ is the noise power at user $k$ in RRH $m$ over subcarrier $n$ and $I_{m,k}^{(n)} = \sum_{i \in \mathcal{K}, \Gamma_{m,k}^{(n)} \leq \Gamma_{m,i}^{(n)}, i \neq k} A_{m,i}^{(n)} \rho_{m,i}^{(n)} p_{m,i}^{(n)} \Gamma_{m,k}^{(n)} + \sum_{j \in \mathcal{M}/\{m\}} \sum_{i \in \mathcal{K}} A_{j,i}^{(n)} \rho_{j,i}^{(n)} p_{j,i}^{(n)} \Gamma_{j,k}^{(n)}$ is the received interference power from the multiplexed users at the same subcarrier and other RRHs.

Based on information theory, in a PD-NOMA based system, user $k$ can successfully detect the signals of user $k'$ which has less SINR than that of user $k$, if the SINR of user $k'$ at user $k$ is higher than its own SINR [16], [30]. Therefore, mathematically we have $\gamma_{m,k}^{(n)}(k') \geq \gamma_{m,k'}^{(n)}(k')$, where $\gamma_{m,k}^{(n)}(k')$ is the SINR of user $k'$ at user $k$ and $\gamma_{m,k'}^{(n)}(k')$ is the SINR of user $k'$. Consequently, from the SINR definition, we have $\frac{p_{m,k'}^{(n)} \Gamma_{m,k}^{(n)}}{\sigma_{m,k}^{(n)} + I_{m,k}^{(n)}} \geq \frac{p_{m,k'}^{(n)} \Gamma_{m,k'}^{(n)}}{\sigma_{m,k'}^{(n)} + I_{m,k'}^{(n)}}$, where it is equivalent to

$$\Omega_{m,k,k'}^{(n)}(\mathbf{A}, \boldsymbol{\rho}, \mathbf{p}) = \Gamma_{m,k'}^{(n)} \sigma_{m,k}^{(n)} - \Gamma_{m,k}^{(n)} \sigma_{m,k'}^{(n)} + \Gamma_{m,k'}^{(n)} \sum_{j \in \mathcal{M}/\{m\}} \sum_{i \in \mathcal{K}} A_{j,i}^{(n)} \rho_{j,i}^{(n)} p_{j,i}^{(n)} \Gamma_{j,k}^{(n)} - \Gamma_{m,k}^{(n)} \sum_{j \in \mathcal{M}/\{m\}} \sum_{i \in \mathcal{K}} A_{j,i}^{(n)} \rho_{j,i}^{(n)} p_{j,i}^{(n)} \Gamma_{j,k'}^{(n)} \leq 0.$$

(1)

The rate of user $k$ over subcarrier $n$ in RRH $m$ is adopted by $r_{m,k}^{(n)}(\mathbf{p}^{(n)}) = \log_2(1 + \gamma_{m,k}^{(n)}(\mathbf{p}^{(n)}))$. Then, the full achievable rate of the user $k$ is expressed as $r_k(\mathbf{A}, \boldsymbol{\rho}, \mathbf{p}) = \sum_{m \in \mathcal{M}} A_{m,k} w_{m,k} \sum_{n \in \mathcal{N}} \rho_{m,k}^{(n)} r_{m,k}^{(n)}(\mathbf{p}^{(n)})$, where $w_{m,k} \in [0,1]$ is a priority weight of the user $k$ in RRH $m$. By regulating these weights, the behavior of proportional fairness between users can be enforced and a trade-off between the user's rate can be adopted and different QoSs or importance levels can be placed by the operator [31]–[33]. Therefore, the total weighted sum rate

of the elastic users can be calculated by $R(\mathbf{A}, \boldsymbol{\rho}, \mathbf{p}) = \sum_{m \in \mathcal{M}} \sum_{k \in \mathcal{K}^e} A_{m,k} w_{m,k} \sum_{n \in \mathcal{N}} \rho_{m,k}^{(n)} r_{m,k}^{(n)}(\mathbf{p}^{(n)})$.

The effect of the data rate change on the power consumption of the fronthaul and the circuit power consumption is neglected since it is rather small compared with the transmit power of RRHs, circuit power consumption and the power consumption in the fiber links. Moreover, the energy consumption of air conditioning is avoided. Therefore, we suppose that the power consumption in the fiber links and the circuit power consumption are fixed to constant values [2], [34]–[37]. Thus, as mentioned before, the total power consumption of the system consists of three parts: 1) the power consumption of the fiber links where the power consumption of each LPN RRH and HPN RRH fiber links are equal to $P_f^L$ and $P_f^H$, respectively, 2) the power consumption at RRHs where the power consumption at each RRH $m$ is equal to $\eta_m \sum_{k \in \mathcal{K}} \sum_{n \in \mathcal{N}} \rho_{m,k}^{(n)} p_{m,k}^{(n)}$ where $\eta_m$ is the efficiency of the power amplifier in each RRH and 3) the circuit power consumption for each LPN RRH $m$ and HPN RRH is equal to $P_c^L$ and $P_c^H$, respectively [2]. Therefore, the total power consumption of the elastic users is expressed as $P(\mathbf{A}, \boldsymbol{\rho}, \mathbf{p}) = P_f^H + M_f P_f^L + \eta_m \sum_{m \in \mathcal{M}} \sum_{k \in \mathcal{K}^e} A_{m,k} \sum_{n \in \mathcal{N}} \rho_{m,k}^{(n)} p_{m,k}^{(n)} + M_f P_c^L + P_c^H$. Thus, the overall energy efficiency performance for the H-CRAN which consists of one HPN RRH and $M_f$ LPN RRHs is defined as $E = \frac{R(\mathbf{A}, \boldsymbol{\rho}, \mathbf{p})}{P(\mathbf{A}, \boldsymbol{\rho}, \mathbf{p})}$. Moreover, the packets for each user are first being queued temporarily where a separate queue is maintained for each user then passed to the radio resource allocator [38]–[40]. Thus, only one queue is required for each user. Therefore, corresponding to each user, we consider the M/G/1 queue model where it is sufficient for our work. This model contributes particular solutions that provides insights into the best model to be chosen for particular queuing situations [41], and as well, it is very reasonable for modelling different types of traffic with various QoS requirements and it is a single server queuing system with unlimited number of waiting positions, [15] and [42]. Hence, the QoS constraints are forced on the streaming users, where we assume that the arrival traffic for user $k \in \mathcal{K}_s$ has a Poisson distribution with arrival rate $\lambda_k$ and the desired maximum delay requirement of the streaming user $k \in \mathcal{K}_s$ is $T_k$. The maximum delay requirement corresponding to each packet arrival rate is $T_k = \frac{q_k}{\lambda_k}$ where $q_k$ is the average queue length, [43] and [15].

### B. Problem Formulation

The cross layer EE maximization resource allocation and RRH selection problem in the downlink H-CRAN can be

mathematically formulated as follows

$$\max_{\boldsymbol{\rho},\mathbf{p},\mathbf{A}} \quad O1 : E = \frac{R(\mathbf{A},\boldsymbol{\rho},\mathbf{p})}{P(\mathbf{A},\boldsymbol{\rho},\mathbf{p})},$$

$$s.t. \quad C1 : \sum_{k \in \mathcal{K}} \rho_{m,k}^{(n)} \le l, \forall m \in \mathcal{M}, n \in \mathcal{N},$$

$$C2 : \rho_{m,k}^{(n)} \in \{0,1\}, \forall m \in \mathcal{M}, n \in \mathcal{N}, k \in \mathcal{K},$$

$$C3 : \sum_{k \in \mathcal{K}} \sum_{n \in \mathcal{N}} A_{m,k} \rho_{m,k}^{(n)} p_{m,k}^{(n)} \le p_m^{\max}, \forall m \in \mathcal{M},$$

$$C4 : 0 \le p_{m,k}^{(n)} \le p_{m,k}^{(n),\mathrm{mask}}, \forall m \in \mathcal{M}, n \in \mathcal{N}, k \in \mathcal{K},$$

$$C5 : \sum_{m \in \mathcal{M}} A_{m,k} \le 1, \forall k \in \mathcal{K},$$

$$C6 : A_{m,k} \in \{0,1\}, \forall m \in \mathcal{M}, k \in \mathcal{K},$$

$$C7 : \overline{X_k} \le T_k, \forall k \in \mathcal{K}^s,$$

$$C8 : A_{m,k} A_{m,k'} \rho_{m,k}^{(n)} \rho_{m,k'}^{(n)} \Omega_{m,k,k'}^{(n)}(\mathbf{A},\boldsymbol{\rho},\mathbf{p}) \le 0,$$

$$\forall m \in \mathcal{M}, n \in \mathcal{N}, k, k' \in \mathcal{K}, \Gamma_{m,k'}^{(n)} \le \Gamma_{m,k}^{(n)}, k \ne k', \quad (2)$$

where $O1$ represents the total energy efficiency for the elastic users. The constraints $C1$ and $C2$ guarantee the PD-NOMA technique assumption on each subcarrier. The constraint $C1$ indicates that maximum $l$ users can be allocated at the same subcarrier. Therefore, when $l = 1$, the system will be equivalent to OFDMA system where at most one user can be allocated to each subcarrier. Then, for example if we have 3 users and $\rho_{m,1}^{(n)} = 1$, $\rho_{m,2}^{(n)} = 0$ and $\rho_{m,3}^{(n)} = 1$, then only the users 1 and 2 are allocated on subcarrier $n$ in RRH $m$. The constraints $C3$ and $C4$ represent the total transmit power limits for each RRH and the transmit power spectral masks for each user, respectively where $p_m^{\max}$ is the maximum allowable transmit power which can be transmitted by RRH $m$ and $p_{m,k}^{(n),\mathrm{mask}}$ is the transmit power spectral mask for user $k$ served by RRH $m$ on subcarrier $n$. Furthermore, the constraints $C5$ and $C6$ ensure the RRH selection assumption. Constraint $C5$ ensures that each user can be served by only one RRH because if $A_{m,k} = 1$ then $A_{m,k'}$ will be equal to zero for any user $k' \ne k$. Furthermore, each user can be allocated to various subcarriers where there is no constraint which limits that. The equation $C7$ defines the streaming users delay constraint where $\overline{X_k}$ is the average time that user $k$ waits in the queue in addition to the service time. Moreover, the constraint $C8$ ensures successful SIC if all $A_{m,k}, A_{m,k'}, \rho_{m,k}^{(n)}$ and $\rho_{m,k'}^{(n)}$ are equal to one. The constraints $C1 - C6$, and $C8$ are system constraints while $C7$ is a service constraint.

In order to solve the considered cross layer EE resource allocation and RRH selection optimization problem (2), we convert the delay constraint $C7$ into another constraint which is in terms of physical-layer parameters. The relationship between the scheduled streaming user $k$ rate and its traffic characteristic $(T_k, \lambda_k)$ is written as [43]

$$\overline{X_k} + \frac{\lambda_k \overline{X_k^2}}{2(1 - \lambda_k \overline{X_k})} \le T_k, \quad (3)$$

where $\overline{X_k}$ and $\overline{X_k^2}$ denote the average and second moment of the service time at the $k^{th}$ user, respectively [43].

Straightforward mathematical manipulation of (3) results in

$$\overline{X_k^2} \le \frac{2T_k - \overline{X_k}(2 + 2T_k\lambda_k) + 2\lambda_k(\overline{X_k})^2}{\lambda_k}. \quad (4)$$

Using the fact that $\overline{X_k^2} \ge (\overline{X_k})^2$ along with (4), we obtain

$$\lambda_k(\overline{X_k})^2 - \overline{X_k}(2 + 2T_k\lambda_k) + 2T_k \ge 0, \quad (5)$$

where the effect of the approximation $\overline{X_k^2} \ge (\overline{X_k})^2$ is tight and there is an ignorable gap between using $\overline{X_k^2}$ and $(\overline{X_k})^2$. Note that $\lambda_k > 0$, therefore, the polynomial in the left hand side of (5) is always greater than or equal to zero for $\overline{X_k} \ge (\overline{X_k}^*)_2$ and $\overline{X_k} \le (\overline{X_k}^*)_1$, where $(\overline{X_k}^*)_1 < (\overline{X_k}^*)_2$ are the roots of the left hand side polynomial in (5). The roots are

$$(\overline{X_k}^*)_{1,2} = \frac{(2 + 2\lambda_k T_k) \pm \sqrt{(2 + 2\lambda_k T_k)^2 - 8\lambda_k T_k}}{2\lambda_k}. \quad (6)$$

As it is seen, both roots are positive. Since we would like that the average service time, i.e., $\overline{X_k}$ to be small, we choose the smaller root. Therefore, holding the inequality in (5) requires that

$$\overline{X_k} \le \frac{(2 + 2\lambda_k T_k) - \sqrt{(2 + 2\lambda_k T_k)^2 - 8\lambda_k T_k}}{2\lambda_k}. \quad (7)$$

Let $\overline{z}$ be a random variable representing the packet size in bits, therefore, $\overline{X_k} = \frac{\overline{z}}{r_k \times B_n}$. Thus, (7) leads us to the following necessary condition [15]

$$C9 : r_k(\mathbf{A},\boldsymbol{\rho},\mathbf{p}) \ge \Psi(\overline{z}, T_k, \lambda_k)(\mathrm{bits/s/Hz}), \forall k \in \mathcal{K}^s, \quad (8)$$

where $\Psi(\overline{z}, T_k, \lambda_k) = \hat{\Psi}(\overline{z}, T_k, \lambda_k)/B_n$ and $\hat{\Psi}(\overline{z}, T_k, \lambda_k) = \frac{2\lambda_k \overline{z}}{(2 + 2\lambda_k T_k) - \sqrt{(2 + 2\lambda_k T_k)^2 - 8\lambda_k T_k}}$.

Thus, the considered optimization problem (2) is reformulated as

$$\max_{\boldsymbol{\rho},\mathbf{p},\mathbf{A}} \quad O1 : E = \frac{R(\mathbf{A},\boldsymbol{\rho},\mathbf{p})}{P(\mathbf{A},\boldsymbol{\rho},\mathbf{p})}, \quad s.t. \quad C1 - C6, C8, C9. \quad (9)$$

The optimization problem (9) is a non-linear program containing both continuous and integer variables. As well, the optimization problem (9) is a NP-hard problem. Therefore, we transform it into an optimization problem with only continuous variables.

Clearly, from $C5$ and $C6$, we obtain that if $A_{m,k} = 1$ then $A_{m',k} = 0 \ \forall m' \ne m$. Thus, if $p_{m,k}^{(n)} \ne 0$ then $p_{m',k}^{(n')} = 0 \ \forall m' \ne m$. Therefore, the RRH selection constraints $C5$ and $C6$ are equivalent to

$$p_{m,k}^{(n)} p_{m',k}^{(n')} = 0, \forall m, m' \in \mathcal{M}, n \in \mathcal{N}, n' \in \mathcal{N}, k \in \mathcal{K}, m \ne m'. \quad (10)$$

The constraint (10) ensures that each user can be at most served by one RRH, since if $p_{m,k}^{(n)} \ne 0$ for RRH $m$ then $p_{m',k}^{(n')} = 0$ for any RRH $m' \ne m$, but each user can be

allocated to various subcarriers in the same RRH because we may have $p_{m,k}^{(n)} \neq 0$ and $p_{m,k}^{(n')} \neq 0$ for $n \neq n'$ which means that user $k$ is allocated to subcarriers $n$ and $n'$, that is because constraint (10) holds only for different RRHs $m \neq m'$. As well, for simplicity we suppose that at most three users can be allocated on the same subcarrier, $l = 3$. Thus, from constraints $C1$ and $C2$, we obtain that if $p_{m,k}^{(n)} \neq 0$, $p_{m,i}^{(n)} \neq 0$ and $p_{m,j}^{(n)} \neq 0$ for users $k$, $i$ and $j$ then $p_{m,x}^{(n)} = 0 \; \forall x \in K$ and $x \neq k \neq i \neq j$. Therefore, the subcarrier allocation constraints $C1$ and $C2$ are equivalent to

$$p_{m,k}^{(n)} p_{m,i}^{(n)} p_{m,j}^{(n)} p_{m,x}^{(n)} = 0, \qquad (11)$$
$$\forall m \in \mathcal{M}, n \in \mathcal{N}, k, i, j, x \in \mathcal{K}, k \neq i \neq j \neq x.$$

Moreover, the constraints (10) and (11) are not compatible with the SCALE method, then the constraints (10) and (11) are replaced by the following constraints

$$C10 : p_{m,k}^{(n)} p_{m',k}^{(n')} \leq \varrho_1, \qquad (12)$$
$$\forall m, m' \in \mathcal{M}, n \in \mathcal{N}, n' \in \mathcal{N}, k \in \mathcal{K}, m \neq m',$$

and

$$C11 : p_{m,k}^{(n)} p_{m,i}^{(n)} p_{m,j}^{(n)} p_{m,x}^{(n)} \leq \varrho_2,$$
$$\forall m \in \mathcal{M}, n \in \mathcal{N}, k \in \mathcal{K}, i \in \mathcal{K}, j \in \mathcal{K}, x \in \mathcal{K}, k \neq i \neq j \neq x, \qquad (13)$$

where $\varrho_1$ and $\varrho_2$ are two small positive numbers. Therefore, the optimization problem (9) can be transformed to

$$\max_{\mathbf{p}} \quad O2 : \frac{R(\mathbf{p})}{P(\mathbf{p})}$$
$$s.t. \quad C4, C10, C11,$$
$$\qquad C12 : \sum_{k \in \mathcal{K}} \sum_{n \in \mathcal{N}} p_{m,k}^{(n)} \leq p_m^{\max}, \forall m \in \mathcal{M},$$
$$\qquad \forall m \in \mathcal{M}, n \in \mathcal{N}, k \in \mathcal{K},$$
$$\qquad C13 : r_k(\mathbf{p}) \geq \Psi(\overline{z}, T_k, \lambda_k), \forall k \in \mathcal{K}^s,$$
$$\qquad C14 : p_{m,k}^{(n)} p_{m,k'}^{(n)} \Omega_{m,k,k'}^{(n)}(\mathbf{p}) \leq 0,$$
$$\qquad \forall m \in \mathcal{M}, n \in \mathcal{N}, k, k' \in \mathcal{K}, \Gamma_{m,k'}^{(n)} \leq \Gamma_{m,k}^{(n)}, k \neq k', \qquad (14)$$

where $R(\mathbf{p}) = \sum_{m \in \mathcal{M}} \sum_{k \in \mathcal{K}^e} w_{m,k} \sum_{n \in \mathcal{N}} r_{m,k}^{(n)}(\mathbf{p}^{(n)})$, $P(\mathbf{p}) = P_f^H + M_f P_f^L + \eta_m \sum_{m \in \mathcal{M}} \sum_{k \in \mathcal{K}^e} \sum_{n \in \mathcal{N}} p_{m,k}^{(n)} + M_f P_c^L + P_c^H$, $r_k(\mathbf{p}) = \sum_{m \in \mathcal{M}} w_{m,k} \sum_{n \in \mathcal{N}} r_{m,k}^{(n)}(\mathbf{p}^{(n)})$, $r_{m,k}^{(n)}(\mathbf{p}^{(n)}) = \log_2(1 + \gamma_{m,k}^{''(n)})$, $\gamma_{m,k}^{''(n)} = \frac{p_{m,k}^{(n)} \Gamma_{m,k}^{(n)}}{\sigma_{m,k}^{(n)} + \overline{I}_{m,k}^{(n)}}$, $\overline{I}_{m,k}^{(n)} = \sum_{i \in \mathcal{K}, \Gamma_{m,k}^{(n)} \leq \Gamma_{m,i}^{(n)}, i \neq k} p_{m,i}^{(n)} \Gamma_{m,k}^{(n)} + \sum_{j \in \mathcal{M}/\{m\}} \sum_{i \in \mathcal{K}} p_{j,i}^{(n)} \Gamma_{j,k}^{(n)}$ and $\Omega_{m,k,k'}^{(n)}(\mathbf{p}) = \Gamma_{m,k'}^{(n)} \sigma_{m,k}^{(n)} - \Gamma_{m,k}^{(n)} \sigma_{m,k'}^{(n)} + \Gamma_{m,k'}^{(n)} \sum_{j \in \mathcal{M}/\{m\}} \sum_{i \in \mathcal{K}} p_{j,i}^{(n)} \Gamma_{j,k}^{(n)} - \Gamma_{m,k}^{(n)} \sum_{j \in \mathcal{M}/\{m\}} \sum_{i \in \mathcal{K}} p_{j,i}^{(n)} \Gamma_{j,k'}^{(n)}$. The objective function $O2$ is not a concave function and is a fractional function. Hence, the optimization problem (14) is a non-convex intractable NP-hard optimization problem. Thus, we transform the fractional objective function $O2$ into a non-fractional subtractive function and then solve the transformed optimization problem.

## III. OPTIMIZATION PROBLEM TRANSFORMATION

The optimization problem (14) is a non linear fractional programming problem which can be transformed by utilizing the well-known Dinkelbach method [44]. Let the optimal energy efficiency value of the optimization problem (14) be $E^* = \frac{R(\mathbf{p}^*)}{P(\mathbf{p}^*)}$.

**Theorem 1.** *The optimal energy efficiency value $E^*$ is achieved if and only if*

$$\max_{\mathbf{p}} \quad R(\mathbf{p}) - E^* P(\mathbf{p}) = R(\mathbf{p}^*) - E^* P(\mathbf{p}^*) = 0, \qquad (15)$$

*where $\mathbf{p}$ is any feasible solution to satisfy the constraints of the optimization problem (14).*

*Proof.* Theorem 1 is proved in two steps by establishing both the sufficient and necessary conditions

1) Clearly, we have $E^* = \frac{R(\mathbf{p}^*)}{P(\mathbf{p}^*)} \geq \frac{R(\mathbf{p})}{P(\mathbf{p})}$, where $\mathbf{p}^*$ is the optimal solution and $\mathbf{p}$ is a feasible solution, which satisfies the constraints of the optimization problem (14). Therefore, we have $R(\mathbf{p}) - E^* P(\mathbf{p}) \leq 0$ and $R(\mathbf{p}^*) - E^* P(\mathbf{p}^*) = 0$. Thus, we obtain that $\max_{\mathbf{p}} \quad R(\mathbf{p}) - E^* P(\mathbf{p}) = 0$ and it is achievable with the optimal solution $\mathbf{p}^*$. Hence, the sufficient condition of Theorem 1 is proved.

2) The objective function of the transformed optimization problem (14) is $R(\mathbf{p}) - E^* P(\mathbf{p})$ and we assume that $\mathbf{p}^{**}$ is the optimal solution of the transformed objective function. Therefore, $R(\mathbf{p}^{**}) - E^* P(\mathbf{p}^{**}) = 0$, then we have $R(\mathbf{p}) - E^* P(\mathbf{p}) \leq R(\mathbf{p}^{**}) - E^* P(\mathbf{p}^{**}) = 0$. Subsequently, $\frac{R(\mathbf{p})}{P(\mathbf{p})} \leq E^*$ and $\frac{R(\mathbf{p}^{**})}{P(\mathbf{p}^{**})} = E^*$. Thus, the optimal solution of the transformed objective function are also the optimal solution for the objective function of the optimization problem (14). Hence, the necessary condition of Theorem 1 is proved. $\square$

Consequently, the transformed optimization problem of the equivalent cross layer EE resource allocation and RRH selection optimization problem (14) is written as

$$\max_{\mathbf{p}} \quad O3 : R(\mathbf{p}) - E^* P(\mathbf{p}) \quad s.t. \quad C4, C10 - C14. \qquad (16)$$

Moreover, an equivalent optimization problem of the transformed optimization problem (16) is represented as

$$\max_{\mathbf{p}} \quad O4 : R(\mathbf{p}) - E P(\mathbf{p}) \quad s.t. \quad C4, C10 - C14, \qquad (17)$$

with the following Lemma.

**Lemma 1.** *for all feasible $\mathbf{p}$ and $E$, $\max_{\mathbf{p}} \quad R(\mathbf{p}) - E P(\mathbf{p})$ is: 1) strictly monotonic decreasing function with respect to $E$, 2) greater than or equal to zero.*

*Proof.* Lemma 1 is proved in two steps:

1) Let $E_1$ and $E_2$ be two optimal values for the two optimal solutions $\mathbf{p}_1$ and $\mathbf{p}_2$, respectively and $E_2 > E_1$. Then, we have $R(\mathbf{p}_1) - E_1 P(\mathbf{p}_1) > R(\mathbf{p}_2) - E_1 P(\mathbf{p}_2) > R(\mathbf{p}_2) - E_2 P(\mathbf{p}_2)$.

Therefore, $\max_{\mathbf{p}} \quad R(\mathbf{p}) - EP(\mathbf{p})$ is a strictly monotonic decreasing function with respect to $E$.

2) Let $\tilde{\mathbf{p}}$ be a feasible solution. Thus, $\tilde{E} = \frac{R(\tilde{\mathbf{p}})}{P(\tilde{\mathbf{p}})}$. Therefore, we have $\max_{\mathbf{p}} \quad R(\mathbf{p}) - \tilde{E}P(\mathbf{p}) \geq R(\tilde{\mathbf{p}}) - \tilde{E}P(\tilde{\mathbf{p}})$. Then $\max_{\mathbf{p}} \quad R(\mathbf{p}) - EP(\mathbf{p})$ is greater than or equal to zero. $\qquad \square$

## IV. SOLVING THE CROSS LAYER EE RESOURCE ALLOCATION AND RRH SELECTION PROBLEM

To solve the optimization problem (17), we apply the following iterative algorithm, where $E$ is updated in each iteration.

$$\overbrace{\mathbf{E}^0 \to \mathbf{p}^0}^{\text{initialization}} \longrightarrow \bullet \bullet \bullet \longrightarrow \overbrace{\mathbf{E}^i \to \mathbf{p}^i}^{\text{Iteration } i} \longrightarrow \bullet \bullet \bullet \longrightarrow \overbrace{\mathbf{E}^* \to \mathbf{p}^*}^{\text{Optimal Solution}}. \tag{18}$$

For utilizing this algorithm, firstly, we have to set an initial value for $E$ denoted by $E^0 = 0$ then find an initial feasible solution $\mathbf{p}^0$ which satisfies the constraints of the optimization problem (17). After that, for each iteration the value of $E$ is updated by $E^{i+1} = \frac{R(\mathbf{p}^i)}{P(\mathbf{p}^i)}$, where for each iteration $i$, $\mathbf{p}^i$, is obtained by solving the following optimization problem

$$\max_{\mathbf{p}} \quad O5: R(\mathbf{p}) - E^i P(\mathbf{p}) \quad s.t. \quad C4, C10 - C14, \tag{19}$$

The process of this algorithm ends when the number of iterations reaches a predefined value which is feasible for practice or $R(\mathbf{p}^i) - E^i P(\mathbf{p}^i) \leq \xi$. The output of the last iteration is the optimal solution of the considered optimization problem.

**Proposition 1.** *The iterative algorithm* (18) *converges to an optimal solution.*

*Proof.* To prove the Proposition 1, we assume that the energy efficiency of the iterations $i$ and $i+1$ are $E^i$ and $E^{i+1}$, respectively, where both of them are greater than zero and not equal to the optimal solution $E^*$ and $E^{i+1} = \frac{R(\mathbf{p}^i)}{P(\mathbf{p}^i)}$. As well, since $E^*$ is the maximum energy efficiency can be achieved then we have $E^{i+1} < E^*$. Moreover, from Lemma 1, we can clearly see that $R(\mathbf{p}) - EP(\mathbf{p}) > 0$ if $E$ is not the optimal value. Therefore, we have $R(\mathbf{p}^i) - E^i P(\mathbf{p}^i) = P(\mathbf{p}^i)\frac{R(\mathbf{p}^i)}{P(\mathbf{p}^i)} - E^i P(\mathbf{p}^i) = P(\mathbf{p}^i)(E^{i+1} - E^i) > 0$. Thus, since $P(\mathbf{p}^i)$ is always greater than zero then we have $E^{i+1} > E^i$. Therefore, after each iteration the energy efficiency $E$ increases. Moreover, according to Lemma 1, after each iteration, due to the increasing of $E$, $R(\mathbf{p}) - EP(\mathbf{p})$ decreases. Furthermore, when the updated value of $E$ increases to the achievable maximum value of $E^*$, the optimization problem (17), with $E^*$ and the optimal condition $R(\mathbf{p}^*) - E^*P(\mathbf{p}^*) = 0$ which is proved in Theorem 1, can be solved. Then, the optimal solution $\mathbf{p}^*$ for the optimization problem (17) is determined. The iterative algorithm updates $E$ to obtain the optimal value $E^*$. Moreover, when the number of iterations is adequately large it can be shown that $\max_{\mathbf{p}} \quad R(\mathbf{p}) - EP(\mathbf{p})$ converges to zero and the optimal condition as expressed in Theorem 1 is attained. Hence, the convergence to the global optimal solution of the outer iterative algorithm is proved [2]. $\qquad \square$

### A. Successive Convex Approximation

The considered optimization problem (30) is non convex. The SCALE method attempts to solve non convex problems by exploiting their underlying convexity which is an iterative algorithm that has low complexity. Therefore, the basic idea behind this approach is applying an inequality which achieves a convex tight lower bound for each non convex function. Thus, to obtain the convexity of this optimization problem, we use the SCALE approach [23] . It can be demonstrated analytically that the SCALE approach has a convergence to a local optimum point. We use the following lower bound [23]

$$\begin{aligned} & \hat{\alpha} \log_2 z + \hat{\beta} \leq \log_2(1+z), \\ & \hat{\alpha} = \frac{z_0}{1+z_0}, \hat{\beta} = \log_2(1+z_0) - \frac{z_0}{1+z_0} \log_2 z_0, \end{aligned} \tag{20}$$

where it is tight at $z = z_0$. Thus, user $k$ rate over subcarrier $n$ in RRH $m$ is approximated to $\hat{r}_{m,k}^{(n),t} = \hat{\beta}_{m,k}^{(n),t} + \hat{\alpha}_{m,k}^{(n),t} \log_2(\gamma_{m,k}''^{(n),t})$, where $\hat{\alpha}_{m,k}^{(n),t} = \frac{\gamma_{m,k}''^{(n),t-1}}{1+\gamma_{m,k}''^{(n),t-1}}$ and $\hat{\beta}_{m,k}^{(n),t} = \log_2(1 + \gamma_{m,k}''^{(n),t-1}) - \hat{\alpha}_{m,k}^{(n),t} \log_2(\gamma_{m,k}''^{(n),t-1})$. Therefore, the optimization problem (30) is rewritten as

$$\begin{aligned} \max_{\mathbf{p}} \quad & \hat{O}5: \sum_{m \in \mathcal{M}} \sum_{k \in \mathcal{K}^e} w_{m,k} \sum_{n \in \mathcal{N}} \hat{r}_{m,k}^{(n)}(\mathbf{p}^{(n)}) - E^i P(\mathbf{p}) \\ s.t. \quad & C4, C10 - C12, C14, \\ & \hat{C}13: \sum_{m \in \mathcal{M}} w_{m,k} \sum_{n \in \mathcal{N}} \hat{r}_{m,k}^{(n)}(\mathbf{p}^{(n)}) \geq \Psi(\overline{z}, T_k, \lambda_k), \\ & \forall k \in \mathcal{K}^s. \end{aligned} \tag{21}$$

The problem (21) is also non convex. Therefore, we apply the change of variable $\mathbf{p} = \exp(\hat{\mathbf{p}})$. Then, we have

$$\begin{aligned} \max_{\hat{\mathbf{p}}} \quad & \hat{O}5e: \sum_{m \in \mathcal{M}} \sum_{k \in \mathcal{K}^e} w_{m,k} \sum_{n \in \mathcal{N}} \hat{r}_{m,k}^{(n)}(e^{\hat{\mathbf{p}}^{(n)}}) - E^i P(e^{\hat{\mathbf{p}}}) \\ s.t. \quad & \hat{C}4e: 0 \leq e^{\hat{p}_{m,k}^{(n)}} \leq p_{m,k}^{(n),\text{mask}}, \forall m \in \mathcal{M}, n \in \mathcal{N}, k \in \mathcal{K}, \\ & \hat{C}10e: e^{\hat{p}_{m,k}^{(n)} + \hat{p}_{m',k}^{(n')}} \leq \varrho_1, \\ & \forall m, m' \in \mathcal{M}, n \in \mathcal{N}, n' \in \mathcal{N}, k \in \mathcal{K}, m \neq m', \\ & \hat{C}11e: e^{\hat{p}_{m,k}^{(n)} + \hat{p}_{m,i}^{(n)} + \hat{p}_{m,j}^{(n)} + \hat{p}_{m,x}^{(n)}} \leq \varrho_2, \\ & \forall m \in \mathcal{M}, n \in \mathcal{N}, k \in \mathcal{K}, i \in \mathcal{K}, j \in \mathcal{K}, x \in \mathcal{K}, \\ & k \neq i \neq j \neq x, \\ & \hat{C}12e: \sum_{k \in \mathcal{K}} \sum_{n \in \mathcal{N}} e^{\hat{p}_{m,k}^{(n)}} \leq p_m^{\max}, \forall m \in \mathcal{M}, \\ & \forall m \in \mathcal{M}, n \in \mathcal{N}, k \in \mathcal{K}, \\ & \hat{C}13e: \sum_{m \in \mathcal{M}} w_{m,k} \sum_{n \in \mathcal{N}} \hat{r}_{m,k}^{(n)}(e^{\hat{\mathbf{p}}^{(n)}}) \geq \Psi(\overline{z}, T_k, \lambda_k), \\ & \forall k \in \mathcal{K}^s, \\ & \hat{C}14e: e^{\hat{p}_{m,k}^{(n)}} e^{\hat{p}_{m,k'}^{(n)}} \Omega_{m,k,k'}^{(n)}(e^{\hat{\mathbf{p}}}) \leq 0, \\ & \forall m \in \mathcal{M}, n \in \mathcal{N}, k, k' \in \mathcal{K}, \Gamma_{m,k'}^{(n)} \leq \Gamma_{m,k}^{(n)}, k \neq k', \end{aligned} \tag{22}$$

Moreover, the optimization problem (21) is also non convex since the constraint $\hat{C}14e$ becomes a non convex function after the transformation $\mathbf{p} = \exp(\hat{\mathbf{p}})$. To obtain the convexity of the constraint $\hat{C}14e$, we apply the difference of two convex function method [45]. Therefore, at iteration $t$, the constraint $\hat{C}14e$ is replaced by

$$
\hat{C}14e' : e^{\hat{p}_{m,k}^{(n)}} e^{\hat{p}_{m,k'}^{(n)}} \Omega_{m,k,k'}^{(n)'}(e^{\hat{\mathbf{p}}}) = e^{\hat{p}_{m,k}^{(n)}} e^{\hat{p}_{m,k'}^{(n)}} (\Gamma_{m,k'}^{(n)} \sigma_{m,k}^{(n)} 
$$
$$
- \Gamma_{m,k}^{(n)} \sigma_{m,k'}^{(n)} + \Gamma_{m,k'}^{(n)} \sum_{j\in\mathcal{M}/\{m\}} \sum_{i\in\mathcal{K}} p_{j,i}^{(n)} \Gamma_{j,k}^{(n)}) -
$$
$$
g(\mathbf{p}^{t-1}) - \bigtriangledown g^T(\mathbf{p}^{t-1})(\mathbf{p}^t - \mathbf{p}^{t-1}) \leq 0,
\tag{23}
$$

where $g(\mathbf{p}) = \Gamma_{m,k}^{(n)} e^{\hat{p}_{m,k}^{(n)}} e^{\hat{p}_{m,k'}^{(n)}} \sum_{j\in\mathcal{M}/\{m\}} \sum_{i\in\mathcal{K}} e^{\hat{p}_{j,i}^{(n)}} \Gamma_{j,k'}^{(n)}$ and its gradient $\bigtriangledown g^T(\mathbf{p}^{t-1})$ is also its super-gradient. Therefore, the optimization problem (22) is transformed to

$$
\max_{\hat{\mathbf{p}}} \quad \hat{O}5e : \sum_{m\in\mathcal{M}} \sum_{k\in\mathcal{K}^e} w_{m,k} \sum_{n\in\mathcal{N}} \hat{r}_{m,k}^{(n)}(e^{\hat{\mathbf{p}}^{(n)}}) - E^i P(\boldsymbol{\rho}, e^{\hat{\mathbf{p}}})
$$
$$
s.t. \quad \hat{C}4e, \hat{C}10e, \hat{C}11e, \hat{C}12e, \hat{C}13e, \hat{C}14e'.
\tag{24}
$$

The optimization problem (24) is a convex approximation problem with respect to the variable $\hat{\mathbf{p}}$ [46], [47]. To solve the considered convex approximation problem (22) using its dual function and related Karush-Kuhn-Tucker (KKT) conditions, we suppose $\boldsymbol{\xi}', \boldsymbol{\zeta}', \boldsymbol{\vartheta}, \boldsymbol{\vartheta}'$ and $\tilde{\boldsymbol{\zeta}}'$ are the Lagrangian multipliers of the approximation problem (22). Thus, after applying the Lagrangian method, the transmit power of each elastic user $k$ over subcarrier $n$ from RRH $m$ is found using (25), where

$\hat{\psi}_{m,k}^{'(n)} = \sum_{i\in\mathcal{K}^e, \Gamma_{m,k}^{(n)} > \Gamma_{m,i}^{(n)}} w_{m,l} \hat{\alpha}_{m,l}^{(n)} \frac{\gamma_{m,l}^{''(n)}}{p_{m,l}^{(n)} \ln(2)}$,

$\overline{\psi}_{m,k}^{'(n)} = \sum_{m'\in\mathcal{M}/\{m\}} \sum_{l\in\mathcal{K}^e} w_{m',l} \hat{\alpha}_{m',l}^{(n)} \frac{\Gamma_{m,l}^{(n)} \gamma_{m',l}^{''(n)}}{p_{m',l}^{(n)} \Gamma_{m',l}^{(n)} \ln(2)}$,

$\tilde{\psi}_{m,k}^{(n)} = \sum_{m'\in\mathcal{M}/\{m\}} \sum_{n'\in\mathcal{N}} 2\vartheta_{mm'knn'} p_{m',k}^{(n')}$,

$\tilde{\psi}_{m,k}^{(n)} = \sum_{i\in\mathcal{K}/\{k,j,x\}} \sum_{j\in\mathcal{K}/\{k,i,x\}} \sum_{x\in\mathcal{K}/\{k,i,j\}} 4\vartheta'_{mnkijx} p_{m,i}^{(n)} p_{m,j}^{(n)} p_{m,x}^{(n)}$,

$\hat{\psi}_{m,k}^{'(n)} = -\sum_{k'\in\mathcal{K}^e, \Gamma_{m,k'}^{(n)} \leq \Gamma_{m,k}^{(n)}, k\neq k'} \tilde{\zeta}_{mnkk'} \Gamma_{m,k}^{(n)} \sum_{j\in\mathcal{M}/\{m\}}$ $\sum_{i\in\mathcal{K}} (p_{m,k'}^{(n)} p_{j,i}^{(n)} \Gamma_{j,k'}^{(n)}) - \sum_{k''\in\mathcal{K}^e, \Gamma_{m,k}^{(n)} \leq \Gamma_{m,k''}^{(n)}, k''\neq k} \tilde{\zeta}_{mnk''k}$ $\Gamma_{m,k''}^{(n)} \sum_{j\in\mathcal{M}/\{m\}} \sum_{i\in\mathcal{K}} (p_{m,k''}^{(n)} p_{j,i}^{(n)} \Gamma_{j,k}^{(n)})$ $-\sum_{m'\in\mathcal{M}/\{m\}} \sum_{k''\in\mathcal{K}^e}$ $\sum_{k'\in\mathcal{K}^e, \Gamma_{m',k'}^{(n)} \leq \Gamma_{m',k''}^{(n)}, k''\neq k'} \tilde{\zeta}'_{m'nk''k'} \Gamma_{m',k''}^{(n)} \Gamma_{m,k'}^{(n)} p_{m',k'}^{(n)}$ $p_{m',k''}^{(n)}, \overline{\tilde{\psi}}_{m,k}^{'(n)} = -\sum_{k'\in\mathcal{K}^e, \Gamma_{m,k'}^{(n)} \leq \Gamma_{m,k}^{(n)}, k\neq k'} \tilde{\zeta}_{mnkk'} \Gamma_{m,k'}^{(n)}$ $\sum_{j\in\mathcal{M}/\{m\}} \sum_{i\in\mathcal{K}} ((p_{m,k'}^{(n)} p_{m,k}^{(n)} p_{j,i}^{(n)})^{t-1} \Gamma_{j,k}^{(n)}) -$ $\sum_{k''\in\mathcal{K}^e, \Gamma_{m,k}^{(n)} \leq \Gamma_{m,k''}^{(n)}, k''\neq k} \tilde{\zeta}_{mnk''k} \Gamma_{m,k}^{(n)}$ $\sum_{j\in\mathcal{M}/\{m\}} \sum_{i\in\mathcal{K}} ((p_{m,k''}^{(n)} p_{m,k}^{(n)} p_{j,i}^{(n)})^{t-1} \Gamma_{j,k''}^{(n)}) -$ $\sum_{m'\in\mathcal{M}/\{m\}} \sum_{k''\in\mathcal{K}^e} \sum_{k'\in\mathcal{K}^e, \Gamma_{m',k'}^{(n)} \leq \Gamma_{m',k''}^{(n)}, k''\neq k'} \tilde{\zeta}'_{m'nk''k'}$ $\Gamma_{m',k'}^{(n)} \Gamma_{m,k''}^{(n)} (p_{m',k'}^{(n)} p_{m',k''}^{(n)} p_{m,k}^{(n)})^{t-1}$.

As well, the transmit power for each streaming user is found by (26), where

$\hat{\psi}_{m,k}^{'(n)} = \sum_{i\in\mathcal{K}^s, \Gamma_{m,k}^{(n)} > \Gamma_{m,i}^{(n)}} w_{m,l} \zeta'_l \hat{\alpha}_{m,l}^{(n)} \frac{\gamma_{m,l}^{''(n)}}{p_{m,l}^{(n)} \ln(2)}$,

$\overline{\psi}_{m,k}^{'(n)} = \sum_{m'\in\mathcal{M}/\{m\}} \sum_{l\in\mathcal{K}^s} w_{m',l} \zeta'_l \hat{\alpha}_{m',l}^{(n)} \frac{\Gamma_{m,l}^{(n)} \gamma_{m',l}^{''(n)}}{p_{m',l}^{(n)} \Gamma_{m',l}^{(n)} \ln(2)}$,

$\hat{\psi}_{m,k}^{'(n)} = -\sum_{k'\in\mathcal{K}^s, \Gamma_{m,k}^{(n)} \leq \Gamma_{m,k}^{(n)}, k\neq k'} \tilde{\zeta}'_{mnkk'} \Gamma_{m,k}^{(n)} \sum_{j\in\mathcal{M}/\{m\}}$ $\sum_{i\in\mathcal{K}} (p_{m,k'}^{(n)} p_{j,i}^{(n)} \Gamma_{j,k'}^{(n)}) - \sum_{k''\in\mathcal{K}^s, \Gamma_{m,k}^{(n)} \leq \Gamma_{m,k''}^{(n)}, k''\neq k} \tilde{\zeta}_{mnk''k}$ $\Gamma_{m,k''}^{(n)} \sum_{j\in\mathcal{M}/\{m\}} \sum_{i\in\mathcal{K}} (p_{m,k''}^{(n)} p_{j,i}^{(n)} \Gamma_{j,k}^{(n)})$ $-\sum_{m'\in\mathcal{M}/\{m\}} \sum_{k''\in\mathcal{K}^s} \sum_{k'\in\mathcal{K}^s, \Gamma_{m,k'}^{(n)} \leq \Gamma_{m,k''}^{(n)}, k''\neq k'} \tilde{\zeta}'_{m'nk''k'}$ $\Gamma_{m',k''}^{(n)} \Gamma_{m,k'}^{(n)} p_{m',k'}^{(n)} p_{m',k''}^{(n)}$,

$\overline{\tilde{\psi}}_{m,k}^{'(n)} = -\sum_{k'\in\mathcal{K}^s, \Gamma_{m,k'}^{(n)} \leq \Gamma_{m,k}^{(n)}, k\neq k'} \tilde{\zeta}_{mnkk'} \Gamma_{m,k'}^{(n)} \sum_{j\in\mathcal{M}/\{m\}}$ $\sum_{i\in\mathcal{K}} ((p_{m,k'}^{(n)} p_{m,k}^{(n)} p_{j,i}^{(n)})^{t-1} \Gamma_{j,k}^{(n)})$ $-\sum_{k''\in\mathcal{K}^s, \Gamma_{m,k}^{(n)} \leq \Gamma_{m,k''}^{(n)}, k''\neq k} \tilde{\zeta}'_{mnk''k}$ $\Gamma_{m,k}^{(n)} \sum_{j\in\mathcal{M}/\{m\}} \sum_{i\in\mathcal{K}} ((p_{m,k''}^{(n)} p_{m,k}^{(n)} p_{j,i}^{(n)})^{t-1} \Gamma_{j,k''}^{(n)})$ $-\sum_{m'\in\mathcal{M}/\{m\}} \sum_{k''\in\mathcal{K}^s} \sum_{k'\in\mathcal{K}^s, \Gamma_{m,k'}^{(n)} \leq \Gamma_{m,k''}^{(n)}, k''\neq k'} \tilde{\zeta}'_{m'nk''k'}$ $\Gamma_{m',k'}^{(n)} \Gamma_{m,k''}^{(n)} (p_{m',k'}^{(n)} p_{m',k''}^{(n)} p_{m,k}^{(n)})^{t-1}$, where the Lagrangian multipliers are updated by applying the sub-gradient method. Algorithm 1 portrays the transmit power allocation algorithm procedures for each iteration in the iterative algorithm where the output is $\mathbf{p}^{t+1}$, for the input $\mathbf{p}^t$ of iteration $t$. The process of Algorithm 1 ends when a predefined threshold $S$ is accessed or if $||\mathbf{p}^{t,s} - \mathbf{p}^{t,s-1}|| < \varpi_2$.

---

**Algorithm 1** Transmit Power Allocation Algorithm

| | |
|---|---|
| 1 | INITIALIZE $s = 0$, $\mathbf{p}^{t,s} = \mathbf{p}^t$, $\hat{\alpha}_{m,k}^{(n),s} = 1$ and |
| 2 | $\hat{\beta}_{m,k}^{(n),s} = 0$, $\forall m \in \mathcal{M}, k \in \mathcal{K}, n \in \mathcal{N}$ |
| 3 | (a simple high-SIR approximation) |
| 4 | REPEAT |
| 5 | $\quad$ Initialize $v = 0$, $\mathbf{p}^{t,s,v} = \mathbf{p}^{t,s}$ and calculate |
| 6 | $\quad$ $\boldsymbol{\xi}'^v, \boldsymbol{\zeta}'^v, \boldsymbol{\vartheta}^v, \boldsymbol{\vartheta}'^v$ and $\tilde{\boldsymbol{\zeta}}'^v$; |
| 7 | $\quad$ Repeat |
| 8 | $\quad\quad$ • Update $\mathbf{p}^{t,s,v}$ using (25) and (26) |
| 9 | $\quad\quad$ • Update $\boldsymbol{\xi}'^v, \boldsymbol{\zeta}'^v, \boldsymbol{\vartheta}^v, \boldsymbol{\vartheta}'^v$ and $\tilde{\boldsymbol{\zeta}}'^v$ |
| 10 | $\quad\quad$ by applying the sub-gradient method, |
| 11 | $\quad\quad$ • $v = v + 1$ |
| 12 | $\quad$ Until $||\mathbf{p}^{t,s,v} - \mathbf{p}^{t,s,v-1}|| < \varpi_1$ |
| 13 | $\quad$ $\mathbf{p}^{t,s} = \mathbf{p}^{t,s,v}$ |
| 14 | $\quad$ Update $\hat{\alpha}_{m,k}^{(n),s+1}$ and $\hat{\beta}_{m,k}^{(n),s+1}$ $\forall m \in \mathcal{M}, k \in \mathcal{K}$, |
| 15 | $\quad$ $n \in \mathcal{N}$ at $(\mathbf{p}^{t,s})$ |
| 16 | $\quad$ $s = s + 1$ |
| 17 | UNTIL $||\mathbf{p}^{t,s} - \mathbf{p}^{t,s-1}|| < \varpi_2$ or $s = S$ |
| 18 | OUTPUT $\mathbf{p}^{t+1} = \mathbf{p}^{t,s}$ |

---

**Proposition 2.** *The Successive Convex Approximation (SCA) with the SCALE approach, creates a sequence of enhanced solutions that converges to a local optimum.*

*Proof.* Let $R_k^{\text{target}} = \Psi(\overline{z}, T_k, \lambda_k)$ and $\hat{r}_k(\mathbf{p}) = \sum_{m\in\mathcal{M}} w_{m,k} \sum_{n\in\mathcal{N}} \hat{r}_{m,k}^{(n)}(\mathbf{p}^{(n)})$. After the first iteration, $t = 1$, because of the high-SIR assumption, we have a feasible solution $\mathbf{p}^1$ [46] and Theorem 1 in [23]. Meanwhile, for every streaming user $k \in \mathcal{K}^s$ and for each iteration $t > 1$,

$$p_{m,k}^{(n)} = \left[ \frac{w_{m,k}\hat{\alpha}_{m,k}^{(n)}\frac{1}{\ln(2)} + \overline{\tilde{\psi}}_{m,k}^{'(n)}}{E^i\eta_m + \xi_m' + \hat{\psi}_{m,k}^{'(n)} + \overline{\psi}_{m,k}^{'(n)} + \tilde{\psi}_{m,k}^{(n)} + \tilde{\psi}_{m,k}^{',(n)} + \hat{\tilde{\psi}}_{m,k}^{'(n)}} \right]_0^{p_{m,k}^{(n),\text{mask}}}, \tag{25}$$

$$p_{m,k}^{(n)} = \left[ \frac{\zeta_k' w_{m,k}\hat{\alpha}_{m,k}^{(n)}\frac{1}{\ln(2)} + \overline{\tilde{\underline{\psi}}}_{m,k}^{'(n)}}{\xi_m' + \hat{\underline{\psi}}_{m,k}^{'(n)} + \overline{\underline{\psi}}_{m,k}^{'(n)} + \tilde{\psi}_{m,k}^{(n)} + \tilde{\psi}_{m,k}^{',(n)} + \hat{\underline{\tilde{\psi}}}_{m,k}^{'(n)}} \right]_0^{p_{m,k}^{(n),\text{mask}}}, \tag{26}$$

we have

$$\begin{aligned} R_k^{\text{target}} &\overset{(i)}{=} \hat{r}_k(\mathbf{p}^{t-1}; \hat{\boldsymbol{\alpha}}^{t-1}, \hat{\boldsymbol{\beta}}^{t-1}) \\ &\overset{(ii)}{\leq} \hat{r}_k(\mathbf{p}^{t-1}) \overset{(iii)}{\leq} \hat{r}_k(\mathbf{p}^{t-1}; \hat{\boldsymbol{\alpha}}^t, \hat{\boldsymbol{\beta}}^t). \end{aligned} \tag{27}$$

In (27), the equality $(i)$ follows from that all the target rate constraints $\hat{C}18e$ are active at the optimal solution of the optimization problem (22), Lemma 2 in [23]. The inequality $(ii)$ follows from the bound in (20) and the equality $(iii)$ follows from the update step of $\hat{\alpha}$ and $\hat{\beta}$ in the transmit power allocation algorithm, [46] and Theorem 1 in [23]. Therefore, it is proved that the solution after each iteration $t-1$, is a feasible solution at iteration $t$.

Additionally, let $\hat{R}(\mathbf{p}; \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}) - E^i\hat{P}(\mathbf{p}) = \sum_{m\in\mathcal{M}}\sum_{k\in\mathcal{K}^e} w_{m,k}\sum_{n\in\mathcal{N}} \hat{r}_{m,k}^{(n)}(\mathbf{p}^{(n)}) - E^i(P_f^H + M_f P_f^L + \eta_m \sum_{m\in\mathcal{M}}\sum_{k\in\mathcal{K}_m^e}\sum_{n\in\mathcal{N}} p_{m,k}^{(n)} + P_c^H + M P_c^L)$. Therefore, we have

$$\begin{aligned} &\hat{R}(\mathbf{p}^t; \hat{\boldsymbol{\alpha}}^t, \hat{\boldsymbol{\beta}}^t) - E^i\hat{P}(\mathbf{p}^t) = \max_{\mathbf{p}} \hat{R}(\mathbf{p}; \hat{\boldsymbol{\alpha}}^t, \hat{\boldsymbol{\beta}}^t) - E^i\hat{P}(\mathbf{p}) \\ &\geq \hat{R}(\mathbf{p}^{t-1}; \hat{\boldsymbol{\alpha}}^t, \hat{\boldsymbol{\beta}}^t) - E^i\hat{P}(\mathbf{p}^{t-1}) = \hat{R}(\mathbf{p}^{t-1}) - E^i\hat{P}(\mathbf{p}^{t-1}) \\ &\geq \hat{R}(\mathbf{p}^{t-1}; \hat{\boldsymbol{\alpha}}^{t-1}, \hat{\boldsymbol{\beta}}^{t-1}) - E^i\hat{P}(\mathbf{p}^{t-1}). \end{aligned} \tag{28}$$

Thus, it is demonstrated that the objective function value, after each iteration $t$, either increases or stays unaltered as that at iteration $t-1$. Therefore, the SCA converges to the last feasible solution acquired due to the compact of the feasible region of the optimization problem. Moreover, according to [48] and [23], the last feasible solution satisfies the necessary KKT conditions of the optimization problem (30). $\square$

### B. Optimal Solution

In order to find the global optimal solution of our system model, we utilize a global optimization framework named monotonic optimization method. Monotonic optimization method takes advantage of the monotonicity or hidden monotonicity in the constraints and the objective function to reduce the computational complexity and provide a guaranteed convergence [24]–[26].

**Definition 1.** *(Monotonicity). For $\mathbf{y}_1 \succeq \mathbf{y}_2$, if $f(\mathbf{y}_1) \geq f(\mathbf{y}_2)$, then, any function $f$ is monotonically increasing.*

**Definition 2.** *(Hyper-rectangle). If $\boldsymbol{b}_1 \preceq \boldsymbol{b}_2$ and $\boldsymbol{b}_1 \preceq \mathbf{y}_1 \preceq \boldsymbol{b}_2$, then, the set of all $\mathbf{y}_1$ is a hyper-rectangle in $[\boldsymbol{b}_1, \boldsymbol{b}_2]$.*

**Definition 3.** *(Normal set). A set $\Upsilon_1$ is a normal set if $\forall \mathbf{y}_1 \in \Upsilon_1$, then the hyper-rectangle $[\boldsymbol{0}, \mathbf{y}_1] \in \Upsilon_1$.*

**Definition 4.** *(Co-normal set). A set $\Upsilon_2$ is a co-normal set in $[\boldsymbol{0}, \boldsymbol{b}_2]$ if $\forall \mathbf{y}_1 \in \Upsilon_2$, then $[\mathbf{y}_1, \boldsymbol{b}_2] \subset \Upsilon_2$.*

**Definition 5.** *(Monotonic optimization). A monotonic optimization problem in canonical form is defined as*

$$\max_{\mathbf{y}_1} f(\mathbf{y}_1) \ s.t. \ \mathbf{y}_1 \in \Upsilon_1 \cap \Upsilon_2, \tag{29}$$

*where $\Upsilon_1 \subset [\boldsymbol{0}, \boldsymbol{b}_2]$ is a normal set with non-empty interior, $\Upsilon_2$ is a closed co-normal set in $[\boldsymbol{0}, \boldsymbol{b}_2]$ and $f$ is an increasing function.*

The considered optimization problem is

$$\max_{\mathbf{p}} \ O5 : R(\mathbf{p}) - E^iP(\mathbf{p}) \quad s.t. \quad C4, C10 - C14, \tag{30}$$

Problem (30) is a non-monotonic problem due to the objective function and the constraints $C13$ and $C14$. Therefore, in order to globally solve the optimization problem (30), we first write the considered optimization problem as a monotonic optimization problem in canonical form, then, we apply the polyblock algorithm [24]–[26]. Thus, let $r_{m,k}^{(n)}(\mathbf{p}) = q_{m,k}^{(n)+}(\mathbf{p}) - q_{m,k}^{(n)-}(\mathbf{p})$ and $p_{m,k}^{(n)}p_{m,k'}^{(n)}\Omega_{m,k,k'}^{(n)}(\mathbf{p}) = \hat{q}_{m,k,k'}^{(n)+}(\mathbf{p}) - \hat{q}_{m,k,k'}^{(n)-}(\mathbf{p})$, where $q_{m,k}^{(n)+}(\mathbf{p}) = \log_2(\sigma_{m,k}^{(n)} + \bar{I}_{m,k}^{(n)} + p_{m,k}^{(n)}\Gamma_{m,k}^{(n)})$, $q_{m,k}^{(n)-}(\mathbf{p}) = \log_2(\sigma_{m,k}^{(n)} + \bar{I}_{m,k}^{(n)})$, $\hat{q}_{m,k,k'}^{(n)+}(\mathbf{p}) = p_{m,k}^{(n)}p_{m,k'}^{(n)}(\Gamma_{m,k'}^{(n)}\sigma_{m,k}^{(n)} - \Gamma_{m,k}^{(n)}\sigma_{m,k'}^{(n)} + \Gamma_{m,k'}^{(n)}\sum_{j\in\mathcal{M}/\{m\}}\sum_{i\in\mathcal{K}} p_{j,i}^{(n)}\Gamma_{j,k}^{(n)})$ and $\hat{q}_{m,k,k'}^{(n)-}(\mathbf{p}) = p_{m,k}^{(n)}p_{m,k'}^{(n)}(\Gamma_{m,k}^{(n)}\sum_{j\in\mathcal{M}/\{m\}}\sum_{i\in\mathcal{K}} p_{j,i}^{(n)}\Gamma_{j,k'}^{(n)})$. Therefore, $R(\mathbf{p}) = \sum_{m\in\mathcal{M}}\sum_{k\in\mathcal{K}^e} w_{m,k}\sum_{n\in\mathcal{N}}(q_{m,k}^{(n)+}(\mathbf{p}) - q_{m,k}^{(n)-}(\mathbf{p}))$ and $r_k(\mathbf{p}) = \sum_{m\in\mathcal{M}} w_{m,k}\sum_{n\in\mathcal{N}}(q_{m,k}^{(n)+}(\mathbf{p}) - q_{m,k}^{(n)-}(\mathbf{p}))$.

The objective function $O5 : R(\mathbf{p}) - E^iP(\mathbf{p})$ can be equivalently rewritten as a difference of two increasing functions

$$R(\mathbf{p}) - E^iP(\mathbf{p}) = q^+(\mathbf{p}) - q^-(\mathbf{p}, E^i), \tag{31}$$

where $q^+(\mathbf{p}) = \sum_{m\in\mathcal{M}}\sum_{k\in\mathcal{K}^e} w_{m,k}\sum_{n\in\mathcal{N}} q_{m,k}^{(n)+}(\mathbf{p})$ and $q^-(\mathbf{p}, E^i) = \sum_{m\in\mathcal{M}}\sum_{k\in\mathcal{K}^e} w_{m,k}\sum_{n\in\mathcal{N}} q_{m,k}^{(n)-}(\mathbf{p}) + E^iP(\mathbf{p})$. Moreover, The set of constraints in $C13$ can be equivalently rewritten as the following single constraint:

$$\min_{\forall k\in\mathcal{K}^s} [q_k^+(\mathbf{p}) - q_k^-(\mathbf{p}) - \Psi(\overline{z}, T_k, \lambda_k)] \geq 0, \tag{32}$$

where $q_k^+(\mathbf{p}) = \sum_{m \in \mathcal{M}} w_{m,k} \sum_{n \in \mathcal{N}} q_{m,k}^{(n)+}(\mathbf{p})$ and $q_k^-(\mathbf{p}) = \sum_{m \in \mathcal{M}} w_{m,k} \sum_{n \in \mathcal{N}} q_{m,k}^{(n)-}(\mathbf{p})$. Then, $\min_{\forall k \in \mathcal{K}^s}[q_k^+(\mathbf{p}) - q_k^-(\mathbf{p}) - \Psi(\overline{z}, T_k, \lambda_k)] = \min_{\forall k \in \mathcal{K}^s}[q_k^+(\mathbf{p}) - (\sum_{\forall k' \in \mathcal{K}^s} q_k^-(\mathbf{p}) - \sum_{\forall k' \in \mathcal{K}^s, k' \neq k} q_k^-(\mathbf{p})) - \Psi(\overline{z}, T_k, \lambda_k)] = \min_{\forall k \in \mathcal{K}^s}[q_k^+(\mathbf{p}) + \sum_{\forall k' \in \mathcal{K}^s, k' \neq k} q_k^-(\mathbf{p}) - \Psi(\overline{z}, T_k, \lambda_k)] - \sum_{\forall k' \in \mathcal{K}^s} q_k^-(\mathbf{p}) \geq 0$, where it is a difference of two increasing functions, $\tilde{q}_k^+(\mathbf{p}) = \min_{\forall k \in \mathcal{K}^s}[q_k^+(\mathbf{p}) + \sum_{\forall k' \in \mathcal{K}^s, k' \neq k} q_k^-(\mathbf{p}) - \Psi(\overline{z}, T_k, \lambda_k)]$ and $\tilde{q}_k^-(\mathbf{p}) = \sum_{\forall k' \in \mathcal{K}^s} q_k^-(\mathbf{p})$. By introducing the auxiliary variables $s_1$, $s_2$, and $\mathbf{s}_3$, the problem formulation (30) is reformulated as [24]–[26]:

$$
\begin{aligned}
\max_{\mathbf{p}, s_1, s_2, \mathbf{s}_3} \quad & O6 : q^+(\mathbf{p}) + s_1, \\
s.t. \quad & C4, C10 - C12, \\
& C15 : 0 \leq s_1 + q^-(\mathbf{p}, E^i) \leq q^-(\mathbf{p}^{\text{mask}}, E^i), \\
& C16 : 0 \leq s_1 \leq q^-(\mathbf{p}^{\text{mask}}, E^i) - q^-(\mathbf{0}, E^i), \\
& C17 : 0 \leq s_2 \leq \tilde{q}_k^-(\mathbf{p}^{\text{mask}}) - \tilde{q}_k^-(\mathbf{0}), \\
& C18 : \tilde{q}_k^-(\mathbf{p}) + s_2 \leq \tilde{q}_k^-(\mathbf{p}^{\text{mask}}), \\
& C19 : \tilde{q}_k^+(\mathbf{p}) + s_2 \geq \tilde{q}_k^-(\mathbf{p}^{\text{mask}}), \\
& C20 : \hat{q}_{m,k,k'}^{(n)+}(\mathbf{p}) + s_{3,m,k,k'}^{(n)} \leq \hat{q}_{m,k,k'}^{(n)+}(\mathbf{p}^{\text{mask}}), \\
& \forall m \in \mathcal{M}, n \in \mathcal{N}, k, k' \in \mathcal{K}, \Gamma_{m,k'}^{(n)} \leq \Gamma_{m,k}^{(n)}, k \neq k', \\
& C21 : \hat{q}_{m,k,k'}^{(n)-}(\mathbf{p}) + s_{3,m,k,k'}^{(n)} \geq \hat{q}_{m,k,k'}^{(n)+}(\mathbf{p}^{\text{mask}}), \\
& \forall m \in \mathcal{M}, n \in \mathcal{N}, k, k' \in \mathcal{K}, \Gamma_{m,k'}^{(n)} \leq \Gamma_{m,k}^{(n)}, k \neq k', \\
& C22 : 0 \leq s_{3,m,k,k'}^{(n)} \leq \hat{q}_{m,k,k'}^{(n)+}(\mathbf{p}^{\text{mask}}) - \hat{q}_{m,k,k'}^{(n)+}(\mathbf{0}), \\
& \forall m \in \mathcal{M}, n \in \mathcal{N}, k, k' \in \mathcal{K}, \Gamma_{m,k'}^{(n)} \leq \Gamma_{m,k}^{(n)}, k \neq k'.
\end{aligned}
\tag{33}
$$

The feasible set of Problem (33) is described by the intersection of the following two sets:

$$
\begin{aligned}
\Upsilon_1 = \{(s_1, s_2, \mathbf{s}_3, \mathbf{P}) : \mathbf{P} \preceq \mathbf{P}^{\text{mask}}, C10, C11, C12, \\
s_1 + q^-(\mathbf{p}, E^i) \leq q^-(\mathbf{p}^{\text{mask}}, E^i), C18, C20\},
\end{aligned}
\tag{34}
$$

and

$$
\Upsilon_2 = \{(s_1, s_2, \mathbf{s}_3, \mathbf{P}) : \mathbf{P} \succeq \mathbf{0}, s1 \geq 0, C19, C21\},
\tag{35}
$$

where $\Upsilon_1$ and $\Upsilon_2$ are the normal and co-normal sets, respectively, in the following hyper-rectangle [24]–[26]

$$
[0, q^-(\mathbf{p}^{\text{mask}}, E^i) - q^-(\mathbf{0}, E^i)] \times [0, \tilde{q}_k^-(\mathbf{p}^{\text{mask}}) - \tilde{q}_k^-(\mathbf{0})] \times
\tag{36}
$$

$$
[0, \hat{q}_{m,k,k'}^{(n)+}(\mathbf{p}^{\text{mask}}) - \hat{q}_{m,k,k'}^{(n)+}(\mathbf{0})] \times [\mathbf{0}, \mathbf{P}^{\text{mask}}].
$$

Problem (33) fulfills Definition 5. Then, Problem (33) is a monotonic optimization problem in a canonical form [24]–[26]. After that, problem (33) is solved by applying the polyblock algorithm.

## V. COMPUTATIONAL COMPLEXITY

In this section, the computational complexity of the proposed optimization problem for both the the solution global optimal approach and suboptimal approach are studied. In this work, in order to find the global optimal solution, we applied the monotonic optimization approach by utilizing the polyblock algorithm.

The polyblock algorithm consists of four main steps as:

- Obtaining the best vertex which its projection belongs to the normal set
- Obtaining the projection of selected vertex
- Removing the improper vertexes
- Obtaining the new vertex set

We consider that the dimensions of the proposed problem is $\overline{T}_1$, the projection of each vertex is given by the bisection algorithm with $\overline{T}_2$ iterations and after $\overline{T}_3$ iterations the polyblock algorithm converges. Then, a simplified complexity order can be given by [26]

$$
\mathcal{O}(\overline{T}_3(\overline{T}_3 \times \overline{T}_1 + \overline{T}_2)).
$$

Moreover, to find the suboptimal solution we applied the SCALE method. To solve the optimization problem (30), one step is applied to determine the power allocation through iterative approach. The power allocation values are obtained by solving (25) and (26). Therefore, in each iteration, the power allocation values are obtained with computational complexity equal to $O(M \times K \times N)$. Moreover, in each iteration, the dual variables are computed with computational complexity equal to $O(M(1 + N \times K + N \times K^4 + N \times K^2 + M \times K \times N^2) + K^s)$ [45]. Thus, for each iteration, the total computational complexity is equal to $O(M \times K \times N)(M(1 + N \times K + N \times K^4 + N \times K^2 + M \times K \times N^2) + K^s)$.

## VI. DISTRIBUTED SOLUTION AND SIGNALLING OVERHEAD DISCUSSION

In this section, at first the distributed solution is explained, and then, the signalling overhead for both centralized and distributed solution are investigated. In order to solve the proposed optimization problem, in a distributed network, at first each RRH initializes the corresponding parameters (power of the assigned users and Lagrangian multipliers) and broadcasts them to the other RRHs. Then, with the received parameters, each RRH calculates the power of the assigned users in addition to updating the corresponding Lagrangian multipliers, and broadcasts them to the other RRHs. Calculation of user power, updating the Lagrangian multipliers, and broadcasting the results is continued until the convergence is achieved. The main steps of distributed solution are summarized as follows:

- Initialize the power of its assigned user and initialize the corresponding Lagrangian multipliers
- Broadcast the initialized parameters
- **Repeat**
  - Receive the broadcasted parameters from the other RRHs
  - Update the corresponding Lagrangian multipliers
  - Calculate the power of its assigned users
  - Check the convergence condition
  - Broadcast the calculated power and Lagrangian multipliers
- end

TABLE II: Quantization of variables

| Feedback variable | Number of bits |
|---|---|
| Each entry of matrices $\boldsymbol{\zeta'}, \boldsymbol{\vartheta}, \boldsymbol{\vartheta'}, \tilde{\boldsymbol{\zeta'}}$ | 3 |
| Each entry of matrices $\boldsymbol{\rho}, \mathbf{p}, \mathbf{A}$ | 3 |
| $h_{m,k}^{(n)}$ | 3 |
| $\sum_{l \in \mathcal{K}^s} w_{m',l} \zeta_l' \hat{\alpha}_{m',l}^{(n)} \frac{\Gamma_{m,l}^{(n)} \gamma_{m',l}''^{(n)}}{p_{m',l}^{(n)} \Gamma_{m',l}^{(n)} \ln(2)}$ | 3 |
| $\sum_{i \in \mathcal{K}} (p_{m,k''}^{(n)} p_{j,i}^{(n)} \Gamma_{j,k}^{(n)})$ | 3 |
| $\sum_{k'' \in \mathcal{K}^s} \sum_{k' \in \mathcal{K}^s, \Gamma_{m,k'}^{(n)} \leq \Gamma_{m,k''}^{(n)}, k'' \neq k'} \tilde{\zeta}_{m'nk''k'}'$ $\Gamma_{m',k''}^{(n)} \Gamma_{m,k'}^{(n)} p_{m',k'}^{(n)} p_{m',k''}^{(n)} \rho_{b_i,n_i,k}, x_{b_i,n_i,k}$ | 3 |
| $\sum_{i \in \mathcal{K}} ((p_{m,k'}^{(n)} p_{m,k}^{(n)} p_{j,i}^{(n)})^{t-1} \Gamma_{j,k}^{(n)})$ | 3 |
| $\sum_{i \in \mathcal{K}} ((p_{m,k''}^{(n)} p_{m,k}^{(n)} p_{j,i}^{(n)})^{t-1} \Gamma_{j,k''}^{(n)})$ | 3 |
| $\sum_{k'' \in \mathcal{K}^s} \sum_{k' \in \mathcal{K}^s, \Gamma_{m,k'}^{(n)} \leq \Gamma_{m,k''}^{(n)}, k'' \neq k'} \tilde{\zeta}_{m'nk''k'}'$ $\Gamma_{m',k'}^{(n)} \Gamma_{m,k''}^{(n)} (p_{m',k'}^{(n)} p_{m',k''}^{(n)} p_{m,k}^{(n)})^{t-1}$ | 3 |
| $\sum_{l \in \mathcal{K}^e} w_{m',l} \hat{\alpha}_{m',l}^{(n)} \frac{\Gamma_{m,l}^{(n)} \gamma_{m',l}''^{(n)}}{p_{m',l}^{(n)} \Gamma_{m',l}^{(n)} \ln(2)}$ | 3 |
| $\sum_{n' \in \mathcal{N}} 2 \vartheta_{mm'knn'} p_{m',k}^{(n')}$ | 3 |
| $\sum_{i \in \mathcal{K}} (p_{m,k'}^{(n)} p_{j,i}^{(n)} \Gamma_{j,k'}^{(n)})$ | 3 |
| $\sum_{i \in \mathcal{K}} (p_{m,k''}^{(n)} p_{j,i}^{(n)} \Gamma_{j,k}^{(n)})$ | 3 |
| $\sum_{k'' \in \mathcal{K}^e} \sum_{k' \in \mathcal{K}^e, \Gamma_{m',k'}^{(n)} \leq \Gamma_{m',k''}^{(n)}, k'' \neq k'}$ $\tilde{\zeta}_{m'nk''k'}' \Gamma_{m',k''}^{(n)} \Gamma_{m,k'}^{(n)} p_{m',k'}^{(n)} p_{m',k''}^{(n)}$ | 3 |



Fig. 3: Signalling overhead for the centralized and distributed approaches.

In the following, the signalling overhead of the centralized and distribution solutions are plotted versus the number of users. The number of bits used for the quantization of the different variables are summarized in Table II. The signalling overhead for the centralized and distributed approaches is shown in Fig. 3. As can be seen, the signalling overhead of the centralized solution is more than that of the distributed solution.

## VII. A FRAMEWORK FOR ACCELERATING THE GENERAL SCALE WITH LAGRANGIAN METHOD USING GPU

In next generation of cellular systems, high computational processing is required which calls for sophisticated method. Thus, in order to tackle this issue, we design a parallel framework for accelerating the general SCALE with the Lagrangian method on GPU using OpenACC API [49]. The OpenACC

API creates high-level heterogeneous programs employing a set of compiler directives to appoint the code's parallel regions in standard C, C++, and Fortran in order to be offloaded from a host central processing unit (CPU) to an attached GPU accelerator [49]. OpenACC directives, facilitate the process of converting an existing serial code into a parallel one in a productive way without substantially exchanging the code. The important task in this work, is to determine the parallel regions of the code.

Algorithm 2 describes all the steps of the SCALE with the Lagrangian method where $\hat{\alpha}$ and $\hat{\beta}$ are the values obtained when applying the lower bound of (20) and $\mathbf{y}$ is the vector of the optimization variables. In each iteration, the optimization variables, Lagrangian multipliers, $\hat{\alpha}$ and $\hat{\beta}$ can be updated independently. Therefore, the parallel regions in the algorithm that have the most calculations are 1) updating the optimization variables, 2) updating the Lagrangian multipliers and 3) updating $\hat{\alpha}$ and $\hat{\beta}$ which can be accelerated using OpenACC API. A few lines to the Fortran code (the highlighted lines in Algorithm (2) have to be added in order to offload the code from the host CPU to the GPU accelerator. These added lines indicate the OpenACC data clause and the kernels loop. The OpenACC data clause imports the data needed for the GPU and as well returns the code output to the host CPU. The kernels loop directive identifies the loops that can be parallelized for the compiler to be executed in parallel on the GPU.

---

**Algorithm 2** SCALE with the Lagrangian Algorithm Using OpenACC Programming Model

---

```
1   INITIALIZE s = 0, y^s = y^initial, α̂ and β̂
2   !$ acc data copyin(input-list) copyout(output-list)
3   REPEAT
4       Initialize v = 0, y^{s,v} = y^s and the Lagrangian
5   multipliers;
6       Repeat
7   !$ acc kernels loop independent
8           • Update y^{s,v}
9   !$ acc kernels loop independent
10          • Update the Lagrangian multipliers,
11          • v = v + 1
12      Until convergence
13      y^s = y^{s,v}
14  !$ acc kernels loop independent
15      Update α̂ and β̂ at (y^s)
16      s = s + 1
17  UNTIL convergence
18  OUTPUT y* = y^s
19  !$ acc end data
```

---

## VIII. SIMULATION RESULTS

In this section, the performance of the proposed algorithm is presented with different numerical experiments. We consider several LPN RRHs located in the coverage of one HPN RRH with 1 Km diameter. The maximum allowable transmit power of the HPN RRH is $p_0^{\max} = 42$ dBm while the maximum allowable transmit power of each LPN RRH is $p_m^{\max} = 23$ dBm, $\forall m \in \mathcal{M}/\{0\}$. Whereas, the spectral mask of each user over each subcarrier is $p_{m,k}^{(n),\text{mask}} = \frac{p_m^{\max}}{N}$ and the predefined

**Algorithm 3** Transmit Power Allocation pseudo code

```
 1 │ INITIALIZE s = 0, pᵗ,ˢ = pᵗ,
 2 │ α̂⁽ⁿ⁾,ˢ_{m,k} = 1 and β̂⁽ⁿ⁾,ˢ_{m,k} = 0 ∀m ∈ M, k ∈ K, n ∈ N
 3 │ (a simple high-SIR approximation)
 4 │ !$ acc data copyin(input-list) copyout(output-list)
 5 │ REPEAT
 6 │     Initialize v = 0, pᵗ,ˢ,ᵛ = pᵗ,ˢ and the
 7 │      Lagrangian multipliers;
 8 │     Repeat
 9 │         do ! Elastic users loop
10 │ !$ acc kernels loop independent
11 │             do ! RRHs loop
12 │ !$ acc loop independent
13 │                 do ! Subcarriers loop
14 │                     Compute the transmit power values
15 │                     of the elastic users using (25).
16 │                 end do
17 │             end do
18 │         end do
19 │         do ! Streaming users loop
20 │ !$ acc kernels loop independent
21 │             do ! RRHs loop
22 │ !$ acc loop independent
23 │                 do ! Subcarriers loop
24 │                     Compute the transmit power values
25 │                     of the streaming users using (26).
26 │                 end do
27 │             end do
28 │         end do
29 │         Update the Lagrangian multipliers by
30 │         applying the sub-gradient method,
31 │         v = v + 1
32 │     Until ||pᵗ,ˢ,ᵛ − pᵗ,ˢ,ᵛ⁻¹|| < ϖ₁
33 │     pᵗ,ˢ = pᵗ,ˢ,ᵛ
34 │
35 │ !$ acc kernels loop independent
36 │     Update α̂⁽ⁿ⁾,ˢ⁺¹_{m,k} and β̂⁽ⁿ⁾,ˢ⁺¹_{m,k} ∀m ∈ M, k ∈ K,
37 │     n ∈ N at (pᵗ,ˢ)
38 │     s = s + 1
39 │ UNTIL ||pᵗ,ˢ − pᵗ,ˢ⁻¹|| < ϖ₂ or s = S
40 │ OUTPUT pᵗ⁺¹ = pᵗ,ˢ
41 │ !$ acc end data
```



Fig. 4: Processing time speed comparison between the serial MATLAB code and the parallel Fortran code implemented on the GPU for different number of parameters.

the parallel fortran pseudo code procedures for each iteration in the iterative algorithm of problem (30) where the output is $\mathbf{p}^{t+1}$, for the input $\mathbf{p}^t$ of iteration $t$. It is worth noting that the loops for updating the transmit power variables, $\hat{\boldsymbol{\alpha}}$ and $\hat{\boldsymbol{\beta}}$ are independent in each iteration. Hence, in order to reduce the processing time, some lines are added to the code as described in Section VII using the Fortran programming language and offloaded from the host CPU to the GPU. Then the variables are updated at the same time by the streaming multi-core processors of the GPU.

In Figs. 4 and 5, we compare the processing time speed between the serial MATLAB code and the Fortran parallel code implemented on the GPU using OpenACC API for different number of parameters where $K = 20$. Fig. 4 shows the processing time speed difference for different number of subcarriers and RRHs where a wide range of values is considered. In Fig. 5, the number of RRHs is fixed to 10. These figures show that in the worst case, by implementing simulations on GPU using OpenACC API, the processing time speed-up of about 255 times with respect to the serial MATLAB code and in the best case the processing time speed-up of about 1058 times is achieved. The hosting CPU used for our simulation is Intel Core i7-4790 with 4 cores and clock speed of 3.6 GHz and the GPU card is NVIDIA GeForce GTX 760. The GPUs architecture is Kepler GK104 with 6 streaming multiprocessor each having 192 stream processors (SPs) thus having the total of 1152 SPs or Compute Unified Device Architecture (CUDA) cores. The GPU works at clock rate of 1150 MHz with memory bandwidth of 192.3 GB/s. It is worth mentioning that if we implement the simulations on a GPU card with different specifications then the processing time speed-up may differ. It is important to note that the significant speed-up is achieved while using a GPU card which is at the same price range as the hosting CPU that is utilized for our simulation.

In the simulations shown in Figs. 6, 7, 8, and 9, we suppose that there are two LPN RRHs installed in the coverage area of the HPN RRH and the total number of subcarriers in each RRH is $N = 32$. Moreover, in the simulations of Figs. 6, 7, 9, and 10, we consider that the packet arrival rate of the streaming

value to end the process of the iterative algorithm is $\xi = 0.01$. The noise power density and the weight of each user are $-174$ dBm/Hz and $w_{m,k} = 1$, respectively. Moreover, $h_{m,k}^{(n)} = \chi_{m,k}^{(n)} d_{m,k}^{-\psi}$ where $d_{m,k}$ is the distance between the RRH $m$ and the user $k$, $\chi_{m,k}^{(n)}$ is an exponential random variable, i.e., representing the Rayleigh fading and $\psi = 3$ is the path loss exponent.

We suppose that the static circuit power consumption is $P_c^L = 0.1$ W and $P_c^H = 3$ W for each LPN RRH and HPN RRH, respectively. Moreover, we assume the power efficiency of each LPN RRH and the HPN RRH to be $\eta_m = 2$, $\forall m \neq 0$ and $\eta_0 = 4$, respectively. Furthermore, the fiber link power consumption between each LPN RRH and the BBU pool is $P_f^L = 1$ W and between the HPN RRH and the BBU pool is $P_f^H = 3$ W. The packet size is 1024 bits and the average queue length, $q_k$, is set to 25 packets.

We simulate the cross layer EE resource allocation problem solution using OpenACC compiler directives on GPU. Algorithm 3 portrays the transmit power allocation algorithm and

Fig. 5: Processing time speed comparison between the serial MATLAB code and the parallel Fortran code implemented on the GPU for different number of subcarriers, $M = 10$.



Fig. 7: The elastic users EE versus the total number of users for different number of streaming users.



Fig. 6: The elastic users EE versus the total number of users for different architectures.

users is 125 packets/s. Then, according to (8), the minimum rate requirement to each streaming user is 4.18 bits/s/Hz and the maximum delay requirement corresponding to each packet arrival rate is $T_k = 0.2$ s.

In Fig. 6, we compare the energy efficiency of H-CRANs with different conventional, 1-tier C-RAN, 2-tier HCN and 1-tier HPN scenarios. In the 1-tier C-RAN scenario, three LPN RRHs are considered. In the 2-tier HCN, one Micro BS (MBS) and two Pico BSs (PBSs) are considered where the static circuit power consumption for the MBS and each PBS are $P_c^M = 10$ W and $P_c^P = 6.8$ W, respectively and the power efficiency for each MBS or PBS is $\eta_0 = 4$. Furthermore, in the 1-tier HPN scenario two MBSs are considered [2]. From Fig. 6, it is shown that the worst energy efficiency is in the 1-tier HPN scenario while energy efficiency in the 2-tier HCN scenario is better than that in the 1-tier HPN scenario since lower transmit power is required and higher sum rate is achieved. Moreover, due to the coverage limitation in the 1-tier C-RAN scenario, the energy efficiency in the 1-tier C-RAN scenario is slightly worse than the 2-tier H-CRAN scenario where the best energy efficiency is reached in the 2-tier H-CRAN scenario due to the advantages of the 1-tier C-RAN and the 2-tier HCN architectures.

The energy efficiency of the elastic users versus the total number of users for various number of streaming users is plotted in Fig. 7. As it is seen, the energy efficiency of the elastic users increases by increasing the total number of users which means increasing the number of elastic users since the number of the streaming users is fixed and that is due to multi-user diversity gain [45] and [50]. As well, In Fig. 7, the effect of the streaming traffic is analyzed. It is observed that the energy efficiency of the elastic users decreases by increasing the number of streaming users. That is because by increasing the number of streaming users, more rate is required for the streaming users, then, less rate will be allocated to the elastic users which will affect the energy efficiency of the elastic users.

In Fig. 8, the effect of the packet arrival rate of streaming traffic is evaluated. The energy efficiency of the elastic users versus the total number of users for different packet arrival rates of streaming users is plotted where the number of streaming users is fixed to 6. By increasing the packet arrival rate of streaming traffic, the minimum required rate of the streaming users is increased then more rate is allocated to the streaming users, therefore, the energy efficiency of the elastic users is affected. Thus, due to what is just described, in Fig. 8, the energy efficiency of the elastic users decreases by increasing the packet arrival rate of the streaming users.

Furthermore, In Figs. 9 and 10, we compare the PD-NOMA and OMA based systems where in OMA based system at most one user can be allocated on a subcarrier. In Fig. 9, the elastic users energy efficiency versus the total number of users is evaluated where the number of the streaming users is fixed to 6. In Fig. 10, the energy efficiency of the elastic users versus the number of LPN RRHs is plotted where the total number of users is 12 which is divided equally between streaming users and elastic users. Clearly, it is observed that the system energy efficiency based on the PD-NOMA technique is better than that based on OMA. Moreover, from Fig. 10, it is seen that by increasing the number of LPN RRHs till $M_f \leq 3$ the energy efficiency of the elastic users increases but when the number of the LPN RRH is $M_f > 3$ both the total sum rate and the power consumption of the elastic users increase approximately in a linear way. Hence, the energy efficiency almost stays

Fig. 8: The elastic users EE versus the total number of users for different packet arrival rate of streaming users.



Fig. 10: The elastic users EE versus number of femtocell RRHs for both PD-NOMA and OMA systems.



Fig. 9: The elastic users EE versus the total number of users for both PD-NOMA and OMA systems.

stable. Moreover, the proposed suboptimal solution with low complexity is perfectly close to the optimal solution.

## IX. CONCLUSION

In this work, we analyzed the performance of the cross layer energy efficiency of PD-NOMA H-CRANs with RRH selection for heterogeneous traffic. In particular, we jointly optimized the RRH selection, subcarrier allocation and transmit power allocation subject to the QoS constraints of streaming users, in addition to the subcarrier and transmit power limitations. In the proposed method, the resources are allocated first to the streaming users and the remaining resources, if exist, are assigned to the elastic users. To solve the considered optimization problem, we utilized the SCA method. Moreover, we obtained the optimal solution of the proposed optimization problem by transforming it to monotonic optimization problem of the canonical form and then applying the polyblock algorithm. Furthermore, we introduced a framework for accelerating SCALE with the Lagrangian method over GPU and we run the proposed particular optimization problem by utilizing OpenACC API. Simulation results showed that the processing time by using OpenACC API on GPU increased for about 1500 times with respect to that by using MATLAB. As well, numerical experiments confirmed that systems based on

the PD-NOMA technique outperforms those based on OMA. Moreover, the energy efficiency in the H-CRAN scenario is shown to perform better than that in the traditional scenarios such as C-RAN, HCN and 1-tier HPN.

## REFERENCES

[1] N. DOCOMO, "5G radio access: Requirements, concept and technologies," *White Paper, July*, 2014.

[2] M. Peng, K. Zhang, J. Jiang, J. Wang, and W. Wang, "Energy-efficient resource assignment and power allocation in heterogeneous cloud radio access networks," *IEEE Transactions on Vehicular Technology*, vol. 64, no. 11, pp. 5275–5287, 2015.

[3] D. Poole, "Introduction to OpenACC directives," in *NVIDIA GPU technology conference*, vol. 12, 2012.

[4] S. Feki, A. Al-Jarro, A. Clo, and H. Bagci, "Porting an explicit time-domain volume-integral-equation solver on GPUs with OpenACC [open problems in cem]," *IEEE Antennas and Propagation Magazine*, vol. 56, no. 2, pp. 265–277, 2014.

[5] T. T. Zygiridis, "High-order error-optimized FDTD algorithm with GPU implementation," *IEEE Transactions on Magnetics*, vol. 49, no. 5, pp. 1809–1812, 2013.

[6] K. Masumnia-Bisheh, M. Ghaffari-Miab, and B. Zakeri, "Evaluation of different approximations for correlation coefficients in stochastic FDTD to estimate SAR variance in a human head model," *IEEE Transactions on Electromagnetic Compatibility*, 2016.

[7] S. R. M. Rostami and M. Ghaffari-Miab, "Fast computation of finite difference generated time-domain green's functions of layered media using OpenAcc on graphics processors," in *proc. IEEE Iranian Conference on Electrical Engineering (ICEE)*, pp. 1596–1599, 2017.

[8] G.-J. Van Den Braak, C. Nugteren, B. Mesman, and H. Corporaal, "GPU-vote: a framework for accelerating voting algorithms on GPU," in *Proc. European Conference on Parallel Processing Springer*, pp. 945–956, 2012.

[9] S. Kumar and P. Baruah, "Communication optimization for multi GPU implementation of smith-waterman algorithm," *matrix*, vol. 80, no. 12, 2013.

[10] C. Sun, Y. Cen, and C. Yang, "Energy efficient OFDM relay systems," *IEEE Transactions on Communications*, vol. 61, no. 5, pp. 1797–1809, 2013.

[11] Y. Wang, W. Xu, K. Yang, and J. Lin, "Optimal energy-efficient power allocation for OFDM-based cognitive radio networks," *IEEE Communications Letters*, vol. 16, no. 9, pp. 1420–1423, 2012.

[12] L. Gao, X. Wang, G. Sun, and Y. Xu, "A game approach for cell selection and resource allocation in heterogeneous wireless networks," in *Proc. IEEE Communications Society Conference on Sensor, Mesh and Ad Hoc Communications and Networks (SECON)*, pp. 530–538, 2011.

[13] D. Amzallag, R. Bar-Yehuda, D. Raz, and G. Scalosub, "Cell selection in 4G cellular networks," *IEEE Transactions on mobile computing*, vol. 12, no. 7, pp. 1443–1455, 2013.

[14] J.-H. Chu, K.-T. Feng, and T.-S. Chang, "Energy-efficient cell selection and resource allocation in LTE-A heterogeneous networks," in *Proc. IEEE 25th Annual International Symposium on Personal, Indoor, and Mobile Radio Communication (PIMRC)*, pp. 976–980, 2014.

[15] N. Mokari, M. R. Javan, and K. Navaie, "Cross-layer resource allocation in OFDMA systems for heterogeneous traffic with imperfect CSI," *IEEE transactions on vehicular technology*, vol. 59, no. 2, pp. 1011–1017, 2010.

[16] Z. Ding, P. Fan, and H. V. Poor, "Impact of user pairing on 5G non-orthogonal multiple access downlink transmissions," *IEEE Transactions on Vehicular Technology*, vol. 65, no. 8, pp. 6010–6014, 2016.

[17] P. Parida and S. S. Das, "Power allocation in OFDM based NOMA systems: A DC programming approach," in *Proc. IEEE Globecom Workshops (GC Wkshps)*, pp. 1026–1031, 2014.

[18] L. Lei, D. Yuan, C. K. Ho, and S. Sun, "Joint optimization of power and channel allocation with non-orthogonal multiple access for 5G cellular systems," in *Proc. IEEE Global Communications Conference (GLOBECOM)*, pp. 1–6, 2015.

[19] A. Mokdad, P. Azmi, and N. Mokari, "Radio resource allocation for heterogeneous traffic in GFDM-NOMA heterogeneous cellular networks," *IET Communications*, vol. 10, pp. 1444–1455, 2016.

[20] M. Moltafet, A. Mokdad, P. Azmi, and N. Mokari, "Radio resource allocation in PD-NOMA based HCN system considering CoMP technology," in *Proc. IEEE International Conference on Electrical and Electronic Engineering, Telecommunication Engineering and Mechatronics (EEETEM)*, 2017.

[21] Z. Ding, X. Lei, G. K. Karagiannidis, R. Schober, J. Yuan, and V. K. Bhargava, "A survey on non-orthogonal multiple access for 5g networks: Research challenges and future trends," *IEEE Journal on Selected Areas in Communications*, vol. 35, no. 10, pp. 2181–2195, 2017.

[22] A. Mokdad, M. Moltafet, P. Azmi, and N. Mokari, "Robust radio resource allocation for heterogeneous traffic in PD-NOMA-based cellular systems," in *Proc. IEEE Iranian Conference on Electrical Engineering (ICEE)*, 2017.

[23] J. Papandriopoulos and J. S. Evans, "Scale: a low-complexity distributed protocol for spectrum balancing in multiuser DSL networks," *IEEE Transactions on Information Theory*, vol. 55, no. 8, pp. 3711–3724, 2009.

[24] A. Zappone, E. Björnson, L. Sanguinetti, and E. Jorswieck, "A framework for globally optimal energy-efficient resource allocation in wireless networks," in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 3616–3620, IEEE, 2016.

[25] A. Zappone, E. Björnson, L. Sanguinetti, and E. Jorswieck, "Globally optimal energy-efficient power control and receiver design in wireless networks," *IEEE Transactions on Signal Processing*, vol. 65, no. 11, pp. 2844–2859, 2017.

[26] M. Moltafet, P. Azmi, N. Mokari, M. R. Javan, and A. Mokdad, "Optimal and fair energy efficient resource allocation for energy harvesting enabled-pd-noma based hetnets," *IEEE Transactions on Wireless Communications*, vol. 17, no. 3, pp. 2054–2067, 2018.

[27] Z. Ding, M. Peng, and H. V. Poor, "Cooperative non-orthogonal multiple access in 5G systems," *IEEE Communications Letters*, vol. 19, no. 8, pp. 1462–1465, 2015.

[28] Y. Liu, Z. Ding, M. Eikashlan, and H. V. Poor, "Cooperative non-orthogonal multiple access in 5G systems with SWIPT," in *Proc. IEEE European Signal Processing Conference (EUSIPCO)*, pp. 1999–2003, 2015.

[29] Z. Ding, Z. Yang, P. Fan, and H. V. Poor, "On the performance of non-orthogonal multiple access in 5G systems with randomly deployed users," *IEEE Signal Processing Letters*, vol. 21, no. 12, pp. 1501–1505, 2014.

[30] D. Tse and P. Viswanath, *Fundamentals of wireless communication*. Cambridge university press, 2005.

[31] I. C. Wong and B. L. Evans, "Optimal downlink OFDMA resource allocation with linear complexity to maximize ergodic rates," *IEEE Transactions on Wireless Communications*, vol. 7, no. 3, pp. 962–971, 2008.

[32] G. Song and Y. G. Li, "Cross-layer optimization for OFDM wireless networks-part II: algorithm development," *IEEE Transactions on Wireless Communications*, vol. 4, no. 2, pp. 625–634, 2005.

[33] Y. Liu and E. Knightly, "Opportunistic fair scheduling over multiple wireless channels," in *Proc. IEEE Societies INFOCOM*, vol. 2, pp. 1106–1115, 2003.

[34] M. Peng, Y. Yu, H. Xiang, and H. V. Poor, "Energy-efficient resource allocation optimization for multimedia heterogeneous cloud radio access networks," *IEEE transactions on Multimedia*, vol. 18, no. 5, pp. 879–892, 2016.

[35] K. Wang, W. Zhou, and S. Mao, "On joint BBU/RRH resource allocation in heterogeneous cloud-RANs," *IEEE Internet of Things Journal*, vol. 4, no. 3, pp. 749–759, 2017.

[36] J. Li, M. Peng, Y. Yu, and Z. Ding, "Energy-efficient joint congestion control and resource optimization in heterogeneous cloud radio access networks," *IEEE Transactions on Vehicular Technology*, vol. 65, no. 12, pp. 9873–9887, 2016.

[37] M. A. Schimuneck, M. Kist, J. Rochol, A. C. Ribeiro-Teixeira, and C. B. Both, "Adaptive monte carlo algorithm to global radio resources optimization in H-CRAN," in *Proc. IEEE International Conference on Communications (ICC)*, pp. 1–6, 2017.

[38] S. Shakkottai, T. S. Rappaport, and P. C. Karlsson, "Cross-layer design for wireless networks," *IEEE Communications magazine*, vol. 41, no. 10, pp. 74–80, 2003.

[39] A. Todini, M. Moretti, A. Valletta, and A. Baiocchi, "Wlc46-1: A modular cross-layer scheduling and resource allocation architecture for ofdma systems," in *Proc. IEEE Global Telecommunications Conference (GLOBECOM)*, pp. 1–6, 2006.

[40] D. S. W. Hui, V. K. N. Lau, and W. H. Lam, "Cross-layer design for OFDMA wireless systems with heterogeneous delay requirements," *IEEE Transactions on Wireless Communications*, vol. 6, no. 8, 2007.

[41] C. C. Zarakovitis, Q. Ni, D. E. Skordoulis, and M. G. Hadjinicolaou, "Power-efficient cross-layer design for OFDMA systems with heterogeneous QoS, imperfect CSI, and outage considerations," *IEEE Transactions on Vehicular Technology*, vol. 61, no. 2, pp. 781–798, 2012.

[42] D. S. W. Hui and V. K. Lau, "Distributive delay-sensitive cross-layer design for OFDMA systems," in *Proc. IEEE International Conference on Communications*, pp. 3542–3546, 2008.

[43] L. Kleinrock, *Queuing Systems*. Hoboken, NJ: Wiley, 1975.

[44] W. Dinkelbach, "On nonlinear fractional programming," *Management Science*, vol. 13, no. 7, pp. 492–498, 1967.

[45] N. Mokari, F. Alavi, S. Parsaeefard, and T. Le-Ngoc, "Limited-feedback resource allocation in heterogeneous cellular networks," *IEEE Transactions on Vehicular Technology*, vol. 65, no. 4, pp. 2509–2521, 2016.

[46] D. T. Ngo, S. Khakurel, and T. Le-Ngoc, "Joint subchannel assignment and power allocation for OFDMA femtocell networks," *IEEE Transactions on Wireless Communications*, vol. 13, no. 1, pp. 342–355, 2014.

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

[48] L. Venturino, N. Prasad, and X. Wang, "Coordinated scheduling and power allocation in downlink multicell OFDMA networks," *IEEE Transactions on Vehicular Technology*, vol. 58, no. 6, pp. 2835–2848, 2009.

[49] OpenACC Working Group, *The OpenACC Application Programming Interface*. Version 2.5, 2015.

[50] A. Goldsmith, *Wireless communications*. Cambridge university press, 2005.

**Ali Mokdad** received the B.Eng. degree in computer and communication engineering from Islamic University of Lebanon, Beirut, Lebanon, in 2008, the M.Eng. degree in Electrical Engineering - Communication Systems from Shahed University, Tehran, Iran, in 2013 and the Ph.D degree in Electrical Engineering - Communication Systems from Tarbiat Modares University, Tehran, Iran, in 2017. His research interests include wireless communications, radio resource allocations and spectrum sharing.

**Paeiz Azmi** (M'05-SM'10) received the B.Sc., M.Sc., and Ph.D. degrees in electrical engineering from Sharif University of Technology (SUT), TehranIran, in 1996, 1998, and 2002, respectively. Since September 2002, he has been with the Electrical and Computer Engineering Department of Tarbiat Modares University, Tehran-Iran, where he became an associate professor on January 2006 and he is a full professor now. His current research interests include modulation and coding techniques, digital signal processing, wireless communications, and estimation and detection theories.

**Nader Mokari** received the Ph.D. degree in electrical engineering from Tarbiat Modares University, Tehran, Iran, in 2014. He joined the Department of Electrical and Computer Engineering, Tarbiat Modares University, as an Assistant Professor, in 2015. He has been involved in a number of large scale network design and consulting projects in the telecom industry. His research interests include design, analysis, and optimization of communication networks.

**Mohammad Moltafet** received his M.Sc. degree from Tarbiat Modares University, Tehran, Iran, in Electrical and Computer Engineering in 2015. He is currently working toward the PhD degree in the Department of Electrical and Computer Engineering, Tarbiat Modares University, Tehran, Iran. His current research interests include wireless communication networks with emphasis on non-orthogonal multiple access (NOMA), and radio resource allocation.

**Mohsen Ghaffari-Miab** (S'06-M'13) received the B.Sc., M.S., and Ph.D. degrees all in electrical engineering from the University of Tehran, Tehran, Iran, in 2005, 2007, and 2012, respectively. From 2010 to 2011, he was a Visiting Scholar at the University of Michigan, Ann Arbor, MI, USA. From 2012 to 2013 he was a Postdoctoral Fellow at the University of Tehran. From 2013 to 2014, he was an Assistant Professor at the Department of Engineering Science, University of Tehran. In 2014, he joined the Department of Electrical and Computer Engineering, Tarbiat Modares University as Assistant Professor. His research interests include theoretical and computational electromagnetics, with focus on frequency- and time-domain integral equation-based methods, finite difference-based methods, GPU-based parallel computing, analysis of layered media, scattering and antenna analysis.