

Hierarchical Multiple Access (HiMA) for Fog-RAN: Protocol Design and Resource Allocation

Vasilis K. Papanikolaou, *Graduate Student Member, IEEE*, Nikos A. Mitsiou, *Student Member, IEEE*, Panagiotis D. Diamantoulakis, *Senior Member, IEEE*, Zhiguo Ding, *Fellow, IEEE*, and George K. Karagiannidis, *Fellow, IEEE*

Abstract—We introduce a set of multiple access protocols, called *hierarchical multiple access (HiMA)*, which are based on non-orthogonal multiple access (NOMA) and time-division multiple access (TDMA), optimized for the hierarchical network scenario. The proposed protocols can be efficiently utilized in various network configurations with an hierarchical form, such as relay networks, cloud-radio access networks (C-RANs), and fog-radio access networks (F-RANs). In particular, C-RANs and, more recently, F-RANs are regarded as promising paradigms to fully utilize the edge of the networks. Therefore, the implementation of novel multiple access protocols to properly exploit these configurations is critical for the fifth generation and beyond of wireless access. Furthermore, the resource allocation problem is formulated for each protocol with respect to the timeslot duration and power. As a result two fairness metrics are optimized, namely max-min rate fairness and proportional fairness. Finally, numerical results reveal the effectiveness of the joint design in the hierarchical network and an interesting trade-off is identified between fairness and achievable rate. Interestingly, despite NOMA being a very promising alternative to conventional multiple access schemes, the protocol that is solely based on NOMA does not always outperform the rest.

Index Terms—non-orthogonal multiple access (noma), hierarchical noma, asynchronous tdma, convex optimization, resource allocation.

I. INTRODUCTION

The development of efficient methods to utilize the radio resources is of great importance in the design of the current and next generations of wireless networks (5G and beyond), where an exponential growth of mobile traffic and connected devices is expected. In the last years, non-orthogonal multiple access (NOMA) has been proposed as a capacity-achieving multiple

access scheme for the Gaussian broadcast channel, which can overcome the limitations of the orthogonal multiple access (OMA) [1]. In NOMA, users' messages are superimposed in a single resource block through multiplexing in the power domain. By applying advanced signal processing techniques at the receiver side, such as successive interference cancellation (SIC) and multi-user detection (MUD), the interference is mitigated, increasing the system's spectral efficiency [2].

On the other side, small-cell architecture is an attractive solution to increase the area capacity, i.e., the total throughput per unit area, which is a fundamental concept in 5G networks. On that matter, the design of smaller cells that operate at higher frequency bands, e.g., mmWave, has made the use of relay networks very attractive, due to higher channel attenuation [3]. Also, a promising alternative to the conventional cellular network architecture is the Cloud-Radio Access Network (C-RAN). C-RAN essentially divides the base station (BS) in two remote parts: the remote radio head (RRH), which grants wireless access to the end users and the centralized pool of baseband units (BBU), at which resource management and large-scale signal processing take place. Capacity limited fronthaul links connect each RRH to the BBU pool. However, mandatory centralized processing and transport over capacity-limited fronthaul links increase the delay, making C-RANs unsuitable for delay-sensitive applications. As an alternative, fog-radio access networks (F-RANs), also known as mobile edge computing (MEC), can perform these tasks in a distributed manner. More specifically, instead of an RRH, in F-RAN a fog access point (FAP) is employed, that can also perform computational tasks. As far as the architecture is concerned, the aforementioned networks share similar traits, inciting the use of a common term to describe them, *hierarchical networks*.

A. Literature

Hierarchical networks like C-RAN and F-RAN have been extensively studied in recent years, mostly in terms of rate, resource allocation, and user scheduling. Especially, NOMA's applicability has been examined in these networks, yielding better spectral efficiency compared to OMA schemes [1], [4], [5]. In [6], stochastic geometry tools was utilized to obtain an expression of the outage probability for NOMA based downlink C-RAN, where RRH are uniformly distributed and they simultaneously serve two paired users. In [7], distributed NOMA was proposed for the uplink of C-RAN, focusing on the optimization of the subset of messages that is decoded

V. K. Papanikolaou, N. A. Mitsiou, P. D. Diamantoulakis, and G. K. Karagiannidis are with the Wireless Communications and Information Processing (WCIP) Group, Department of Electrical and Computer Engineering, Aristotle University of Thessaloniki, 54636, Thessaloniki, Greece (e-mail: vpapanikk@auth.gr; nmitsiou@ece.auth.gr; padiaman@ieee.org; geokarag@auth.gr).

Z. Ding is with the School of Electrical and Electronic Engineering, the University of Manchester, Manchester, UK (e-mail: zhiguo.ding@manchester.ac.uk).

The work of V. K. Papanikolaou and G. K. Karagiannidis has been co-financed by the European Union and Greek national funds through the Competitiveness, Entrepreneurship and Innovation Operational Program (EPAnEK), under the special actions aquaculture - industrial materials - open innovation in culture (project code: T6YBP-00134).

The work of N. A. Mitsiou and P. D. Diamantoulakis has been supported by the European Union's Horizon 2020 research and innovation programme under grant agreement No 957406.

The work of Z. Ding was supported by the UK EPSRC under grant number EP/P009719/2 and by H2020-MSCA-RISE under grant number 10100641.

Digital Object Identifier xx.xxxx/xxxx.xxxx.xxxxxxx

by each RRH and the corresponding decoding order with the aim to enhance the achievable rate region, assuming that the RRHs can exchange digital information with the BBU via ideal feedback links. Also, in [8] the outage probability of distributed uplink NOMA was derived, considering fixed transmission rates and the use of non-ideal feedback links.

Research activity on resource allocation in this type of networks has yielded some interesting results, as well. More specifically, in [9], a F-RAN with NOMA was considered, where capacity-limited fronthaul links connect the BBU pool and the FAPs. In that paper the user assignment and resource allocation are jointly optimized for the users, by maximizing the weighted sum rate. Furthermore, the NOMA strategy seems to be superior, compared to conventional orthogonal multiple access (OMA), in terms of fairness. In [10], a NOMA based C-RAN system was investigated and resource allocation optimized in terms of weighted sum rate, while it was shown that NOMA can achieve higher spectral efficiency than conventional OMA schemes. In [11], the authors have optimized the fronthaul rate allocation of the uplink of a C-RAN multi-cell NOMA system to ensure high-throughput in cell-edge regions. In order to decrease the system complexity, users are partitioned in groups based on a novel and iterative algorithm. Furthermore, a C-RAN with NOMA was also investigated in [12], where a heuristic algorithm for user scheduling and power allocation was proposed, that shows the superiority of NOMA in terms of sum-rate. In [13], a multi-tier heterogeneous C-RAN with NOMA was studied. The energy efficiency was investigated for various environments and the maximum number of cells that can be supported in the network was shown. More details about the energy efficiency of heterogeneous NOMA C-RANs can be found in [14], where key promising technologies were presented. Furthermore, in [15], the energy efficiency was optimized for a NOMA C-RAN, where a sub-6 GHz link was assumed as a fronthaul and the wireless access is provided via mmWave links. Also, in [16], a C-RAN with content multicast based on file popularity was studied. In that paper, the authors solved a maximization problem of the minimum delivery rate of files that are requested by the users. It was shown that NOMA with content-centric multicast outperforms the OMA schemes in terms of delivery rate. In [17], an energy efficient maximization problem in an F-RAN with subchannel reuse assignment was modeled as a Stackelberg game, showing that low latency can be achieved at the users with low complexity algorithms. In [18], a many-to-one matching game was used to jointly optimize resource allocation for the weighted sum rate with NOMA. Moreover, in [19], the problem of minimizing the average delay, while finding the optimal cache placement was studied. For that non-convex problem, the McCormick envelopes and the Lagrange partial relaxation method were used to study the optimal subchannel assignment and power allocation to maximize the sum rate in three different transmission modes. In [20], a pricing two stage Stackelberg game was developed to minimize the energy through power allocation, while taking into account the inter-cell interference. For that purpose, a matching algorithm is used to extract the optimal subchannel allocation.

Cooperative networks, which can offer high reliability such as low latency or low outage probability can also be considered as hierarchical networks. In [21] the authors investigated the use of NOMA in a cooperative V2X network for low latency broadcasting or multicasting. Other notable contributions on cooperative NOMA networks are [21]–[27]. However, cooperative NOMA networks, as those assumed in these works, do not employ multiple relays each of which can serve a set of users, as in hierarchical networks. Instead, the stronger users, having already decoded the weaker users' messages, due to SIC, can act as a decode and forward relay.

B. Motivation and Contributions

Inspired by the above contributions and the increasing interest for the C-RAN and F-RAN architectures, in this paper we propose the use of a generalized cooperative framework, termed as *hierarchical multiple access (HiMA)*. HiMA comprises three multiple access schemes to be utilized by hierarchical networks. Different from most cooperative NOMA scenarios, an hierarchical network includes a central node that communicates with dedicated relay nodes, where each of them serves their own users. A close work is the one presented in [9]. However, a fixed capacity fronthaul link is adopted in [9], whereas in this contribution the fronthaul is optimized as well. In general, as opposed to [9], the main scope of this work is to introduce and optimize appropriate HiMA protocols, that do not assume a fixed capacity fronthaul link, while also considering the time and energy constraints for both hops of all cooperative links. Moreover, in the present contribution an arbitrary number of FAPs and UEs can be considered in the hierarchical network, in contrast to most works on NOMA that study only scenarios with a pair of users. Additionally, in this contribution it is revealed that, in contrast to conventional networks, NOMA is not always the best option when user fairness is concerned. Following this, novel multiple access schemes are proposed to be utilized in the hierarchical network, as well as a version of NOMA optimized for use in this type of networking scenario. The contributions of this paper are summarized below:

- Novel multiple access schemes, termed as Hierarchical Multiple Access (HiMA) are introduced for the hierarchical networks. The three schemes are, namely, *Hierarchical NOMA (HiNOMA)*, *Asynchronous TDMA (A-TDMA)*, and *mixed A-TDMA/HiNOMA*. These multiple access schemes are described in terms of energy and power consumption, and user scheduling.
- An optimization framework is proposed, that guarantees fair allocation of the available resources between the network nodes and is jointly performed for both hops of the hierarchical network, i.e., the fronthaul link is not assumed to be perfect.
- Two different fairness metrics are optimized, maximization of the minimum rate among the users and maximization of the proportional fairness. The first metric guarantees the fair rate allocation at the cost of the sum rate, while proportional fairness offers a trade-off between fairness and total throughput.

- Simulation results are presented to validate the proposed analysis and compare the proposed schemes. Interesting remarks are offered concerning the operation of the hierarchical network. The proposed schemes are compared with a benchmark TDMA, where resource allocation is also optimized. Finally, it is shown that NOMA is not always superior to OMA and the mixed solution can outperform both NOMA and OMA.

C. Structure

In section II, a comprehensive system model is presented. In section III, the proposed protocols for the hierarchical multiple access are presented, namely hierarchical NOMA, asynchronous TDMA and a mixed protocol. Following that, in section IV, resource allocation is optimized for different fairness metrics for each protocol design. Also, in section V, numerical results are presented and discussed. Finally, in section VI, some conclusions of this work are drawn.

II. SYSTEM MODEL

We consider the downlink transmission of an *hierarchical* F-RAN, with a BBU pool, multiple fog access points (FAPs), and multiple user equipment nodes (UEs). In line with the main principles of F-RAN, this work focuses on the use of multiple lightweight access points, which communicate wirelessly with both the BBU and the UEs. In more detail, the network operates in a resource block where the BBU serves $|\mathcal{N}| = N$ RNs and in N resource blocks where each FAP, n , serves a total of $|\mathcal{M}_n| = M_n$ users in their respective resource block, where $\mathcal{N} = \{1, \dots, n, \dots, N\}$, $\mathcal{M}_n = \{1, \dots, m, \dots, M_n\}$, respectively, and the operator $|\mathcal{A}|$ denotes the cardinality of set \mathcal{A} . Additionally $M = \sum_n M_n$ denotes the total number of UEs. Each UE is served by only one FAP and we further assume that both FAP and UE are equipped with single antennas.

Moreover, the FAPs are assumed to perform half-duplex (HD) decode and forward (DF) out-band relaying [28], [29] since full duplex relaying is avoided due to the aggregated interference and increased complexity, while pure time domain multiplexing used for in-band relaying cannot provide perfect loop interference protection [28]. As such, there is no inter-cell interference between UEs served by different FAPs. The BBU utilizes an orthogonal frequency resource block B_0 to serve the FAPs and each FAP uses a corresponding B_n block to serve its users. HD relaying dictates that the FAPs listen to the BBU in the first hop with duration τ_1 and transmit to the UEs during the second hop, which lasts τ_2^n , for each FAP. The total timeslot duration is required to follow $\tau_1 + \tau_2^n = 1$, which is true according to the normalization of the timeslot duration. Moreover, it is assumed that full channel state information (CSI) is available at the BBU, which also has the computational capabilities to perform the orchestration of the F-RAN. To this end, for the second hop, each subgroup of UEs obtains the CSI via the use of pilot symbols transmitted by the FAPs that serves them. The FAPs obtain the CSI of the first hop and the second hop through pilot symbols transmitted by the BBU and feedback from the UEs, respectively. On the

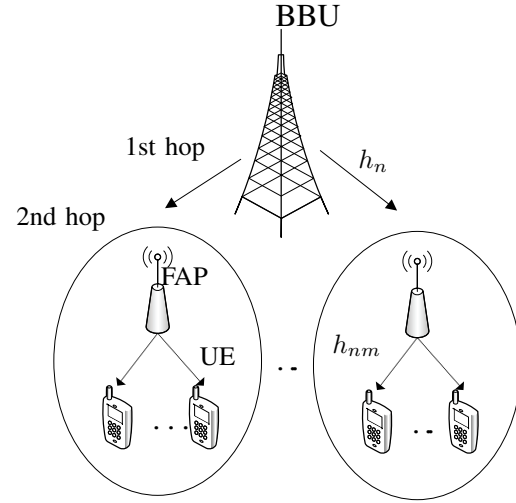


Fig. 1. System Model of an Hierarchical Network

other hand, the channel estimation can also happen during the uplink, taking advantage of the channel reciprocity.

In addition, it is assumed that all nodes consume energy solely for information transmission. This is because the circuit power consumption of the BBU and the FAPs during downlink is much lower than the transmit power and is therefore considered negligible. Furthermore, the BBU has a limited supply of energy that can spend in a single timeslot (i.e., the average power that is consumed). This could potentially lead to cases where a relatively high amount of energy is used for a very small timeslot, while the average power stays within the limits. However, hardware issues and safety regulations limit the maximum transmitted power as well as the average power. Therefore, a set of two constraints regarding the power need to hold for the BBU:

$$\tau_1^n P_{\tau_1} \leq E_{\text{BBU}}, \quad (1)$$

where τ_1^n is the timeslot duration of the first hop for FAP n and P_{τ_1} is the total transmitted power of the BBU at the specific timeslot τ_1^n , and

$$P_{\tau_1} \leq P_{\text{BBU}}, \quad (2)$$

where P_{BBU} denotes the maximum allowed power of the BBU. In a similar manner, the energy and power constraints for the FAPs can be expressed as

$$\tau_2^n P_{\tau_2} \leq E_{\text{FAP}}, \quad (3)$$

where τ_2^n is the timeslot duration of the second hop for FAP n and P_{τ_2} is the total transmitted power of the FAP at the specific timeslot τ_2^n , and

$$P_{\tau_2} \leq P_{\text{FAP}}, \quad (4)$$

where P_{FAP} is the maximum allowed power for the FAPs, which are assumed to have the same limitations concerning power and energy.

Finally, in order for the FAPs to be able to support their respective UEs, the achievable data rate of FAP n needs to be

greater or equal than the sum of the required rates of its UEs, i.e.,

$$R_n \geq \sum_{m=1}^{M_n} R_{nm}, \forall n \in \mathcal{N}. \quad (5)$$

III. PROPOSED MULTIPLE ACCESS PROTOCOLS FOR HIERARCHICAL NETWORKS

In this section, we propose three novel protocols for HiMA, namely hierarchical NOMA (HiNOMA), asynchronous TDMA (A-TDMA), and Mixed A-TDMA/HiNOMA.

A. Hierarchical NOMA (HiNOMA)

In *hierarchical* NOMA (HiNOMA) both the first and the second hop in the cooperative network are carried out by power domain NOMA. Since HD relaying is used, the two hops are orthogonal in time, while the first hop has a duration of τ_1 , with $0 \leq \tau_1 \leq 1$ of the total timeslot duration and the second hop durates τ_2 , with $0 \leq \tau_2 \leq 1$, respectively. Moreover, τ_1 and τ_2 are normalized to the total duration of the timeslot as

$$\tau_1 + \tau_2 \leq 1. \quad (6)$$

During the transmission phase of the first hop, N signals are transmitted to each FAP from the BBU. The baseband equivalent of the received signal y_n at FAP n is given by

$$y_n = h_n \sum_{i=1}^N \sqrt{p_i} s_i + w_n, \quad (7)$$

where h_n denotes the channel gain coefficient between the n -th FAP and the BBU, p_i represents the allocated power for the i -th FAP, s_i denotes the message sent from the BBU to the n -th FAP and w_n is the additive Gaussian white noise (AWGN) at the receiver of the n -th FAP. Since in NOMA messages for the users in a group are transmitted simultaneously, the total energy and power consumption constraints for the timeslot need to apply for the sum of the allocated powers. Therefore for the energy requirement, (1) is formulated as

$$\tau_1 \sum_{n=1}^N P_n \leq E_{\text{BBU}}, \quad (8)$$

while (2) is expressed as

$$\sum_{n=1}^N P_n \leq P_{\text{BBU}}. \quad (9)$$

The FAP decodes the received message and transmits with power domain NOMA to its users (decode and forward relaying). In the second hop, each FAP transmits in non-interfering resource blocks the message to its respective users. Therefore, the received signal y_{nm} of UE m that is served by FAP n is given by

$$y_{nm} = h_{nm} \sum_{j=1}^{M_n} \sqrt{q_{nj}} s_j + w_{nm}, \quad (10)$$

where q_j is the power coefficient of the j -th UE that is served by FAP n . Identically to the first hop, the total power consumption constraint for the timeslot needs to apply for

the sum of the allocated powers in the second hop as well, therefore (4) is expressed as

$$\sum_{m=1}^{M_n} q_{nm} \leq P_{\text{FAP}}, \forall n \in \mathcal{N}. \quad (11)$$

Similarly, for the energy requirement, (1) is formulated as

$$\tau_2 \sum_{m=1}^{M_n} q_{nm} \leq E_{\text{FAP}}, \forall n \in \mathcal{N}. \quad (12)$$

The achievable data rate in downlink NOMA is determined opportunistically by the FAP's channel condition, therefore SIC is successful if the FAP's are optimally ordered based on their channel conditions, i.e., $|h_m|^2 \geq |h_j|^2$, for $m > j$. Then, each user decodes the messages of all the users prior to them, and the users with better channel conditions than them remain as interfering terms. Ultimately the ascending decoding order is the optimal for downlink NOMA, since it is the only one that achieves the capacity limits of the broadcast channel [30]. Furthermore SIC is always perfect since adaptive rate is assumed for all nodes, which occur from the optimization problems, meaning all data rates belong in the capacity region of downlink NOMA [1], [8], [31]. Therefore, the achievable rate of FAP n can be expressed as

$$R_n = \tau_1 B_0 \log_2(1 + \gamma_n), \quad (13)$$

while the achievable data rate of UE m that is served by FAP n is given by

$$R_{nm} = \tau_2 B_n \log_2(1 + \gamma_{nm}), \quad (14)$$

where γ_n denotes the signal-to-interference plus noise ratio (SINR) after successive interference cancellation (SIC) at the FAP n , written as

$$\gamma_n = \frac{|h_n|^2 p_n}{|h_n|^2 \sum_{i=n+1}^N p_i + \sigma^2}, \quad (15)$$

where σ^2 denotes the variance of the AWGN at the receiver. Similarly, γ_{nm} denotes the SINR after SIC at the UE m and can be expressed as

$$\gamma_{nm} = \frac{|h_{nm}|^2 q_{nm}}{|h_{nm}|^2 \sum_{i=m+1}^{M_n} q_{ni} + \sigma^2}. \quad (16)$$

Furthermore, it should be noted that the rate expressed in (13) and (14) is the maximum achievable rate without considering jointly the two hops.

In order for the FAPs to be non-interfering, an orthogonal frequency resource block B_n is allocated to each of them, while the BBU transmits to the FAPs using a separate resource block B_0 . This is utilized for the time where the first hop is active, so the total time-frequency resource block (TFRB) for the first hop is $\tau_1 B_0$. For the second hop, the TFRB block is the sum of the respective resource blocks used by all FAPs, i.e., $\sum_{n=1}^N \tau_2^n B_n$, where τ_2^n is the time that the second hop is active for the FAP n . Therefore, HiNOMA needs a total of $\tau_1 B_0 + \tau_2 \sum_{n=1}^N B_n$ TFRBs.

TABLE I
RESOURCE COORDINATION OF THE PROPOSED PROTOCOLS

	HiNOMA		A-TDMA			Mixed A-TDMA/HiNOMA		
BBU at B_0	τ_1		τ_1^1	τ_1^2	τ_1^3	τ_1^1	τ_1^2	τ_1^3
FAP ₁ at B_1		τ_2		$\sum_{m=1}^{M_1} \tau_2^{1m}$			τ_2^1	
FAP ₂ at B_2		τ_2			$\sum_{m=1}^{M_2} \tau_2^{2m}$			τ_2^2
FAP ₃ at B_3		τ_2		$\sum_{m=1}^{M_3} \tau_2^{3m}$			τ_2^3	

B. Asynchronous TDMA (A-TDMA)

Next, we present an OMA scheme, designed and optimized for hierarchical networks, which, without loss of generality, is considered to be a follow up to the time-division multiple access (TDMA) concept. In this scheme, a timeslot τ_1^n is assigned to FAP n in order to receive data from the BBU, i.e., first hop, and then for the remainder of the timeslot the FAP transmits to its assigned UEs (second hop). In this way, an FAP that has received its message from the BBU does not have to wait for the BBU to finish transmitting to the rest of the FAPs and it can transmit immediately afterwards to its UEs. We call this OMA scheme *asynchronous* TDMA (A-TDMA). TDMA is also considered for the second hop of each FAP, with each UE assigned to a specific timeslot τ_2^{nm} . Similar to HiNOMA, the FAPs transmit at non-interfering orthogonal blocks. So, the following constraints hold for the A-TDMA

$$\sum_{n=1}^N \tau_1^n \leq 1, \quad (17)$$

which ensures that the duration of the first hop is lower or equal to the timeslot duration (normalized). Also, given the HD operation of the FAPs, the following constraint holds, similar to the HiNOMA case

$$\tau_1^n + \sum_{m=1}^{M_n} \tau_2^{nm} \leq 1, \quad \forall n \in \mathcal{N}. \quad (18)$$

This ensures that the duration of the two hops is lower or equal to the timeslot duration for each FAP. Due to the use of TDMA, a constraint similar to (17) is needed for the second hop as well, however this constraint is already integrated in (18).

Furthermore, we consider the same energy and maximum power constraints as in the previous section, which limit the total power and energy transmitted to the maximum allowed. In A-TDMA these constraints can be expressed as

$$\sum_{n=1}^N \tau_1^n P_n \leq E_{\text{BBU}}, \quad (19)$$

with $P_n \leq P_{\text{BBU}}$ for the BBU and

$$\sum_{m=1}^{M_n} \tau_2^{nm} q_{nm} \leq E_{\text{FAP}}, \quad (20)$$

with $q_{nm} \leq P_{\text{FAP}}, \forall n \in \mathcal{N}$.

The resource blocks needed for the A-TDMA are the bandwidth of the first hop for the total time the BBU transmits, i.e., $B_0 \sum_{n=1}^N \tau_1^n$. Also, the non-interfering FAPs use different frequency resource blocks, B_n , each for the specific

time needed for the second hop of that FAP. Therefore, for the second hop the total TFRBs are $\sum_{n=1}^N B_n \sum_{m=1}^{M_n} \tau_2^{nm}$. The total TFRBs needed for this system are $B_0 \sum_{n=1}^N \tau_1^n + \sum_{n=1}^N B_n \sum_{m=1}^{M_n} \tau_2^{nm}$. Finally, the achievable data rate of FAP n is given by

$$R_n = \tau_1^n B_0 \log_2 \left(1 + \frac{|h_n|^2 P_n}{\sigma^2} \right). \quad (21)$$

Similarly, for the m -th UE that is served by FAP n , the achievable data rate is

$$R_{nm} = \tau_2^{nm} B_n \log_2 \left(1 + \frac{|h_{nm}|^2 q_{nm}}{\sigma^2} \right). \quad (22)$$

C. Mixed A-TDMA/HiNOMA

In this subsection, we combine the two aforementioned protocols and propose a mixed asynchronous TDMA - hierarchical NOMA system. In the mixed system, the first hop utilizes the A-TDMA scheme and each FAP for the second hop serves the assigned UEs to it with NOMA. It deserves to be noted that the proposed mixed protocol is different to a hybrid NOMA/OMA scheme, as the protocol spans multiple hops, while hybrid NOMA/OMA schemes are employed in a single hop. Specifically, A-TDMA is employed in the first hop to benefit from its asynchronous nature, while NOMA is utilized in the second due to its higher spectral efficiency. The constraints of each hop in the mixed system are identical to the ones required in the respective employed protocol. Therefore, for the first hop (17), (19), and $P_n \leq P_{\text{BBU}}$ hold. For the second hop, the constraints from NOMA need to hold, specifically (11) and (12). Due to the use of half duplex FAPs, the following constraint needs to hold for the timeslots of the two hops

$$\tau_1^n + \tau_2^n \leq 1, \quad (23)$$

which reflects the asynchronous nature of the protocol but differs from the A-TDMA since each FAP n transmits simultaneously to their assigned UEs during τ_2^n .

The TFRBs that are required by the mixed protocol are the bandwidth used for the first hop, B_0 for the total timeslot duration that the first hop is active, i.e., $\sum_{n=1}^N \tau_1^n$ and the bandwidth B_n of each non-interfering FAP that each is used for τ_2^n duration. So the total TFRBs that are required are $B_0 \sum_{n=1}^N \tau_1^n + \sum_{n=1}^N B_n \tau_2^n$.

Finally, in this protocol the achievable data rate by FAP n is given by

$$R_n = \tau_1^n B_0 \log_2 \left(1 + \frac{|h_n|^2 P_n}{\sigma^2} \right). \quad (24)$$

Similarly, for the m -th UE that is served by FAP n , the achievable data rate is given by

$$R_{nm} = \tau_2^n B_n \log_2 \left(1 + \frac{|h_{nm}|^2 q_{nm}}{|h_{nm}|^2 \sum_{i=m+1}^{M_n} q_{ni} + \sigma^2} \right). \quad (25)$$

IV. OPTIMAL RESOURCE ALLOCATION

In this section, a resource allocation problem is formulated for the hierarchical network to ensure fairness among the UEs. User fairness is an important metric in the hierarchical network that needs to be taken into account in the resource allocation of the system, since nodes with weak channel conditions can end up with low data rates if more resources are allocated to stronger users in order to achieve higher total throughput. To this end, in this paper we optimize two of the most common fairness metrics in communications, the minimum rate of users and the proportional fairness [5], [32]. Therefore, we formulate two optimization problems for each of the proposed protocols in the hierarchical network. Since the constraints are the same in each problem, we denote the function f as the objective function in the following problems. Furthermore, since the signaling requirements are identical among the proposed protocols, the optimization and comparison among the proposed protocols focuses on the communications phase to explore the advantages of HiMA. In addition, considering that all of the proposed schemes need full CSI for the optimal resource allocation, the associated cost in terms of energy and time is the same and thus, the comparison between the protocols is fair.

The maximization of the minimum rate of the system leads to the maximum data rate that all users can achieve simultaneously. Hence, f in this problem is defined as

$$f = \min (R_{nm}(\mathbf{q}, \boldsymbol{\tau})), \forall m \in \mathcal{M}_n, \forall n \in \mathcal{N}. \quad (26)$$

The proportional fairness metric is defined by the sum-log-rate of the UEs. This is a fairness metric because this utility function shrinks rapidly for low values of rates. This is due to the logarithm's ability to tend to negative infinity when its argument tends to zero. Therefore, solutions that offer very low data rates to some UEs yield significantly lower proportional fairness, thus it's very unlikely that such a solution will be deemed the optimal. Moreover, proportional fairness is an increasing function with the achieved data rate. So, its maximization also increases the spectral efficiency in the system. Consequently, proportional fairness is a trade-off between user fairness and sum throughput in the system. Then f is defined for the case of proportional fairness optimization as

$$f = \sum_{n=1}^N \sum_{m=1}^{M_n} \log R_{nm}(\mathbf{q}, \boldsymbol{\tau}). \quad (27)$$

A. Hierarchical NOMA

When HiNOMA is used in the network, the following problem can be formulated for optimizing fairness according

to an objective function f .

$$\begin{aligned} & \mathbf{max}_{\boldsymbol{\tau}, \mathbf{P}, \mathbf{q}} \quad f \\ & \mathbf{s.t.} \quad C_1 : \sum_{m=1}^{M_n} R_{nm}(\mathbf{q}, \tau_2) \leq R_n(\mathbf{P}, \tau_1), \forall n \in \mathcal{N}, \\ & \quad C_2 : \tau_1 + \tau_2 \leq 1, \quad C_3 : \sum_{n=1}^N P_n \leq P_{\text{BBU}}, \\ & \quad C_4 : \sum_{m=1}^{M_n} q_{nm} \leq P_{\text{FAP}}, \forall n \in \mathcal{N}, \quad C_5 : \tau_1 \sum_{n=1}^N P_n \leq E_{\text{BBU}}, \\ & \quad C_6 : \tau_2 \sum_{m=1}^{M_n} q_{nm} \leq E_{\text{FAP}}, \quad \forall n \in \mathcal{N}, \\ & \quad C_7 : 0 \leq p_n \leq P_{\text{BBU}}, \quad \text{and} \quad 0 \leq q_{nm} \leq P_{\text{FAP}}. \end{aligned} \quad (28)$$

First, due to the nature of the constraints, it deserves to be observed that the aforementioned problem is always feasible, regardless of the channel gains and the power and energy limitations, while the same holds for all the considered problems in this section. More specifically, C_1 is a fundamental constraint for the operation of this network; it denotes the fact that the link between the BBU and each FAP must have greater capacity than the rate the UEs served by that FAP require. C_2 is necessary since the total duration of the two hops cannot be larger than the timeslot duration. The rest of the constraints follow the power and energy limitations of the HiNOMA system as they were described in the Section III, i.e., (9), (11), (8), (12). Note that the optimization of HiNOMA is complicated due to the interference terms appearing in the logarithm term of constraints C_1 and C_2 and the coupling between the power allocation at both the BBU and the fog access points and the time that is allocated at each hop.

1) *Minimum Rate*: Our aim is to optimize resource allocation, given by the power allocation coefficients \mathbf{p} and \mathbf{q} for every FAP and every UE respectively, and the timeframe duration of the two hops, $\boldsymbol{\tau}$, so that the minimum rate is maximized in the hierarchical NOMA system. We can express the max-min problem in its epigraph form, using a hypograph variable R_{\min} , which is expressed as

$$R_{nm}(\mathbf{q}, \tau_2) \geq R_{\min}. \quad (29)$$

Hence, the achievable rate at the UEs is at least R_{\min} for every UE. Therefore, according to the first constraint, each FAP needs to achieve a rate greater or equal to that of its respective UEs. Since, all UEs can achieve at least R_{\min} and there are M_n UEs for FAP n then $R_n \geq M_n R_{\min}$. Therefore, C_1 can be expressed in two parts as

$$C_{1a} : R_{nm}(\mathbf{q}, \tau_2) \geq R_{\min}, \quad (30)$$

$$C_{1b} : R_n(\mathbf{P}, \tau_1) \geq M_n R_{\min}. \quad (31)$$

Thus, the following optimization problem is formed as:

$$\begin{aligned}
 & \max_{\mathbf{p}, \mathbf{q}, \tau, R_{\min}} R_{\min} \\
 \text{s.t. } & C_{1a} : \tau_2 B_n \log_2 \left(1 + \frac{|h_{nm}|^2 q_{nm}}{|h_{nm}|^2 \sum_{i=m+1}^{M_n} q_{ni} + \sigma^2} \right) \geq R_{\min}, \\
 & \quad \forall m \in \mathcal{M}_n \quad \text{and} \quad \forall n \in \mathcal{N}, \\
 & C_{1b} : \tau_1 B_0 \log_2 \left(1 + \frac{|h_n|^2 P_n}{|h_n|^2 \sum_{i=n+1}^N P_i + \sigma^2} \right) \geq M_n R_{\min}, \\
 & \quad \forall n \in \mathcal{N}, \\
 & (28).C_2, (28).C_3, (28).C_4, (28).C_5, (28).C_6, (28).C_7.
 \end{aligned} \tag{32}$$

It is noted that the optimization problem in (32) is non-convex. The main reasons of non-convexity are the expressions of the capacity in both hops, i.e., C_{1a} and C_{1b} . More specifically, the term of the interference in the SINR leads to the inclusion of the power variable in the denominator, so the function is non-concave. Additionally, the inclusion of timeslot duration as a variable that is multiplied with the logarithm function causes the function to be non-concave, as well. Moreover, C_5 and C_6 are not convex because of the multiplication of τ_1 with P_n and τ_2 with q_{nm} . Constraints C_2 , C_3 , C_4 , and C_7 are linear in their present form. Therefore the complexity to solve this problem is high, mainly due to the relation of the rates with the power allocation variables. Thus, it is important to prove, that the problem in (32) can be transformed to a convex one; so, the process to find a global maximum can be solved in polynomial time.

Proposition 1: The optimization problem in (32) can be formulated as a convex one and is given in (33)

Following Proposition 1, the equivalent convex problem of (32) can be expressed as follows:

$$\begin{aligned}
 & \max_{\tilde{\mathbf{P}}, \tilde{\mathbf{q}}, \tilde{\tau}, \tilde{R}_{\min}} \tilde{R}_{\min} \\
 \text{s.t. } & C_{1a} : -\tilde{q}_{nm} - \log(|h_{nm}|^2) + \log \left(2^{\frac{\exp(\tilde{R}_{\min} - \tilde{\tau}_2)}{\tilde{B}_n}} - 1 \right) \\
 & + \log \left(\sigma^2 + |h_{nm}|^2 \sum_{i=m+1}^{M_n} e^{\tilde{q}_{ni}} \right) \leq 0, \forall n \in \mathcal{N}, \forall m \in \mathcal{M}_n, \\
 & C_{1b} : -\tilde{P}_n - \log(|h_n|^2) + \log \left(2^{\frac{M_n}{B_0} \exp(\tilde{R}_{\min} - \tilde{\tau}_1)} - 1 \right) \\
 & + \log \left(\sigma^2 + |h_n|^2 \sum_{i=n+1}^N e^{\tilde{P}_i} \right) \leq 0, \quad \forall n \in \mathcal{N}, \\
 & C_2 : e^{\tilde{\tau}_1} + e^{\tilde{\tau}_2} \leq 1, \quad C_3 : \sum_{n=1}^N e^{\tilde{P}_n} \leq P_{\text{BBU}}, \\
 & C_4 : \sum_{m=1}^{M_n} e^{\tilde{q}_{nm}} \leq P_{\text{FAP}}, \forall n \in \mathcal{N}, \quad C_5 : \sum_{n=1}^N e^{\tilde{P}_n + \tilde{\tau}_1} \leq E_{\text{BBU}}, \\
 & C_6 : \sum_{m=1}^{M_n} e^{\tilde{q}_{nm} + \tilde{\tau}_2} \leq E_{\text{FAP}}, \quad \forall n \in \mathcal{N}.
 \end{aligned} \tag{33}$$

Proof: The proof is provided in Appendix A. ■

Due to its convexity, the optimization problem in (33) can be solved by a decomposition method. Even though, mainly

due to $C_3 - C_6$, the problem could be decoupled to two smaller subproblems, one for each hop, and choose the R_{\min} based on the minimum of the two solutions, constraint C_2 hinders the effectiveness of this method, since the two hops cannot be decoupled completely, leading to a suboptimal solution. As such, primal decomposition methods are not appropriate for this problem. To this end the following considerations emerge:

- A dual decomposition method is required to solve this problem since the constraints cannot be decoupled, due to the two hops sharing a common timeslot. This is shown in constraint C_2 .
- As it is analytically shown in the proof of Proposition 1, the maximization problem in (33) is convex, since the objective function is concave with respect to all the optimization variables, the left terms of the constraints are convex and it satisfies the Slater's constraint qualification. Thus, the duality gap between the dual and the primal solution is zero [33]. Therefore, the solution of the dual problem leads to the optimal solution of the original problem.

Based on the above, problem (33) is solved with the Lagrange dual decomposition method. Note that the proposed method of Lagrange dual decomposition converges to the optimal solution in polynomial time for a convex problem. The solution of problem (33) is provided in Appendix B. The optimal solution to the proposed problem is found algorithmically rather than analytically, because of the multi-dimensionality of the problem, which is in part due to the arbitrary number of FAPs and UEs in the system.

2) *Proportional Fairness:* Next, we study the optimization of the proportional fairness metric based on the same constraints as before. The objective function f in the formulated problem then is defined as in (27) with regards to the resource allocation in the system. More specifically, given the power allocation coefficients, \mathbf{p} and \mathbf{q} for every FAP and every UE respectively, and the timeframe duration of the two hops, τ , the proportional fairness is maximized in the hierarchical NOMA system. Thus, the following optimization problem is formed as:

$$\begin{aligned}
 & \max_{\tau, \mathbf{p}, \mathbf{q}} \sum_{n=1}^N \sum_{m=1}^{M_n} \log(R_{nm}(\mathbf{q}, \tau_2)) \\
 \text{s.t. } & C_1 : \sum_{m=1}^{M_n} R_{nm}(\mathbf{q}, \tau_2) \leq R_n(\mathbf{P}, \tau_1), \quad \forall n \in \mathcal{N}, \\
 & (28).C_2, (28).C_3, (28).C_4, (28).C_5, (28).C_6, (28).C_7.
 \end{aligned} \tag{34}$$

Problem (34) is non-convex. The main reason behind this is the expression of the achievable rates by both the FAPs and the UEs. More specifically, due to the interference in the SINR, terms of power appear in the denominator in the logarithm. Moreover, the objective function is the logarithm of the rate and therefore it is non-concave. More importantly, in constraint C_1 a difference of logarithms appears causing the function to be non-convex. Finally, as in the case of the previous problem, C_5 and C_6 are non-convex because of the multiplication of τ_1 with P_n and τ_2 with q_{nm} . Therefore, the complexity to solve this problem is high and, in order to find a global maximum

in polynomial time, it is vital to transform problem (34) into an equivalent convex one.

Proposition 2: The optimization problem in (34) can be formulated as a convex one and is expressed as in (35).

Following Proposition 2, the equivalent convex problem of (34) can be expressed as follows:

$$\begin{aligned}
 & \max_{\tilde{\mathbf{P}}, \tilde{\mathbf{q}}, \tau, \tilde{r}_{nm}, \tilde{r}_n} \sum_{n=1}^N \sum_{m=1}^{M_n} \tilde{r}_{nm} \\
 \text{s.t. } & C_1 : \sum_{m=1}^{M_n} e^{\tilde{r}_{nm} - \tilde{r}_n} - 1 \leq 0, \quad \forall n \in \mathcal{N}, \\
 & (33).C_2, (33).C_3, (33).C_4, (33).C_5, (33).C_6, \\
 & C_7 : -\tilde{q}_{nm} - \log(|h_{nm}|^2) + \log\left(2^{\frac{\exp(\tilde{r}_{nm} - \tilde{r}_2)}{B_n}} - 1\right) \\
 & \quad + \log\left(\sigma^2 + |h_{nm}|^2 \sum_{k>m}^{M_n} \exp(\tilde{q}_{nk})\right) \leq 0, \\
 & C_8 : -\tilde{P}_n - \log(|h_n|^2) + \log\left(2^{\frac{1}{B_0} \exp(\tilde{r}_n - \tilde{r}_1)} - 1\right) \\
 & \quad + \log\left(\sigma^2 + |h_n|^2 \sum_{k>n}^{M_n} \exp(\tilde{P}_k)\right) \leq 0,
 \end{aligned} \tag{35}$$

where the last two constraints hold $\forall m \in \mathcal{M}_n$ and $\forall n \in \mathcal{N}$.

Proof: The proof is provided in Appendix C. ■

The problem in (35) can be solved in the same manner as problem (33). Since it is a convex problem, standard convex optimization methods, such as the subgradient method [34] that was used in the previous problem, or the interior-point method, can be utilized to efficiently solve this problem in polynomial time.

B. Asynchronous TDMA

Next, we move to develop the optimization framework for the asynchronous TDMA protocol. The constraints are adjusted according to the protocol analysis in the previous section and the following problem can be formulated and solved:

$$\begin{aligned}
 & \max_{\mathbf{p}, \mathbf{q}, \tau} f \\
 \text{s.t. } & C_1 : \sum_{m=1}^{M_n} R_{nm}(\mathbf{q}, \tau) \leq R_n(\mathbf{p}, \tau), \quad \forall n \in \mathcal{N}, \\
 & C_2 : \tau_1^n + \sum_{m=1}^{M_n} \tau_2^{nm} \leq 1, \quad \forall n \in \mathcal{N}, \quad C_3 : \sum_{n=1}^N \tau_1^n \leq 1, \\
 & C_4 : \sum_{n=1}^N \tau_1^n P_n \leq E_{\text{BBU}}, \quad C_5 : \sum_{m=1}^{M_n} \tau_2^{nm} q_{nm} \leq E_{\text{FAP}}, \quad \forall n \in \mathcal{N}, \\
 & C_6 : p_n \leq P_{\text{BBU}}, \quad \text{and} \quad q_{nm} \leq P_{\text{FAP}}, \quad \forall m \in \mathcal{M}_n, \quad \forall n \in \mathcal{N}.
 \end{aligned} \tag{36}$$

1) *Minimum Rate:* First, in the A-TDMA case as well, we solve the max-min problem. Following a similar procedure as in the previous subsection for the HiNOMA protocol, the max-min problem in A-TDMA can be expressed in its epigraph form as follows:

$$\begin{aligned}
 & \max_{\mathbf{p}, \mathbf{q}, \tau, R_{\min}} R_{\min} \\
 \text{s.t. } & C_{1a} : \tau_2^{nm} B_n \log_2 \left(1 + \frac{|h_{nm}|^2 q_{nm}}{\sigma^2}\right) \geq R_{\min}, \\
 & \forall m \in \mathcal{M}_n, \quad \forall n \in \mathcal{N}, \\
 & C_{1b} : \tau_1^n B_0 \log_2 \left(1 + \frac{|h_n|^2 P_n}{\sigma^2}\right) \geq M_n R_{\min}, \\
 & \forall n \in \mathcal{N}, \\
 & (36).C_2, (36).C_3, (36).C_4, (36).C_5, (36).C_6
 \end{aligned} \tag{37}$$

where R_{\min} is the hypograph variable of the problem when it is expressed in its epigraph form. Problem (37) is non-convex due to the multiplication of variables τ with a logarithm of power terms, i.e., P_n or q_{nm} . Moreover, C_4 and C_5 are non-convex as well, due to the multiplication of τ with P_n and q_{nm} , respectively. In order to solve this problem in a tractable manner and in polynomial time, the following proposition is necessary.

Proposition 3: Problem (37) can be formulated as an equivalent convex problem and is expressed as in (38).

Following Proposition 3, the equivalent convex problem of (37) can be expressed as follows:

$$\begin{aligned}
 & \max_{\mathbf{E}, \tau, R_{\min}} R_{\min} \\
 \text{s.t. } & C_{1a} : \tau_2^{nm} B_n \log_2 \left(1 + \frac{|h_{nm}|^2 E_{nm}}{\tau_2^{nm} \sigma^2}\right) - R_{\min} \geq 0, \\
 & \forall m \in \mathcal{M}_n, \quad \forall n \in \mathcal{N}, \\
 & C_{1b} : \tau_1^n B_0 \log_2 \left(1 + \frac{|h_n|^2 E_n}{\tau_1^n \sigma^2}\right) - M_n R_{\min} \geq 0, \quad \forall n \in \mathcal{N}, \\
 & (36).C_2, (36).C_3, (36). \\
 & C_4 : \sum_{n=1}^N E_n \leq E_{\text{BBU}}, \quad C_5 : \sum_{m=1}^{M_n} E_{nm} \leq E_{\text{FAP}}, \quad \forall n \in \mathcal{N}, \\
 & C_6 : E_n - \tau_1^n P_{\text{BBU}} \leq 0, \quad \text{and} \quad E_{nm} - \tau_2^{nm} P_{\text{FAP}} \leq 0, \\
 & \forall m \in \mathcal{M}_n, \quad \forall n \in \mathcal{N},
 \end{aligned} \tag{38}$$

where E_{nm} denotes the energy of UE m served by FAP n , given by $E_{nm} = q_{nm}/\tau_2^{nm}$. Respectively, $E_n = P_n/\tau_1^n$ is the energy of FAP n .

Proof: By exploiting the definition of energy consumption coefficients E_n and E_{nm} , the constraints C_{1a} and C_{1b} (37) can be rewritten as (38). The Hessian of C_{1a}, C_{1b} is proven to be negative semi-definite since its eigenvalues are equal or less than zero, so these constraints are in concave form. Constraints C_2, C_3 in (36) are linear. Also, the rest of the constraints, i.e., $C_4 - C_6$, are also transformed into linear ones, using the energy transformation as previously. So, problem (37) can be modeled as an equivalent convex problem. ■

2) *Proportional Fairness:* Next, we study the optimization of the proportional fairness metric based on the same constraints as before for the A-TDMA protocol. The objective function f in the formulated problem then is defined as in (27) with regards to the resource allocation in the system. More specifically, given the power allocation coefficients, \mathbf{p} and \mathbf{q} for every FAP and every UE respectively, and the timeframe

duration assigned to each FAP and its respective UEs, τ , the proportional fairness is maximized in the asynchronous TDMA system. Thus, the following optimization problem is formed as:

$$\begin{aligned} \max_{\tau, \mathbf{P}, \mathbf{q}} \quad & \sum_{n=1}^N \sum_{m=1}^{M_n} \log(R_{nm}(q_{nm}, \tau_2^{nm})) \\ \text{s.t.} \quad & C_1 : \sum_{m=1}^{M_n} \tau_2^{nm} B_n \log_2 \left(1 + \frac{|h_{nm}|^2 q_{nm}}{\sigma^2} \right) \\ & \leq \tau_1^n B_0 \log_2 \left(1 + \frac{|h_n|^2 P_n}{\sigma^2} \right), \forall n \in \mathcal{N}, \\ & (36).C_2, (36).C_3, (36).C_4, (36).C_5, (36).C_6. \end{aligned} \quad (39)$$

The problem in (39) is non-convex, mainly due to the first constraint, C_1 , where the difference of logarithms is not a convex function. Moreover, the multiplication of τ variables and power variables, P_n and q_{nm} in the objective function and in the constraints leads the problem to be classified as non-convex. Once again, a transformation is necessary to solve this problem in a tractable manner.

Proposition 4: Problem (39) can be transformed to an equivalent convex problem and is expressed as in (40).

Following Proposition 4, the equivalent convex problem of (39) can be expressed as follows:

$$\begin{aligned} \max_{\tau, \mathbf{E}, \mathbf{r}_n, \mathbf{r}_{nm}} \quad & \sum_{n=1}^N \sum_{m=1}^{M_n} \log(r_{nm}) \\ \text{s.t.} \quad & C_1 : \sum_{m=1}^{M_n} r_{nm} - r_n \leq 0 \quad \forall n \in \mathcal{N}, \\ & (38).C_2, (38).C_3, (38).C_4, (38).C_5, (38).C_6, \\ & C_7 : \tau_2^{nm} B_n \log_2 \left(1 + \frac{|h_{nm}|^2 E_{nm}}{\tau_2^{nm} \sigma^2} \right) - r_{nm} \geq 0, \\ & \forall n \in \mathcal{N}, \forall m \in \mathcal{M}_n, \\ & C_8 : \tau_1^n B_0 \log_2 \left(1 + \frac{|h_n|^2 E_n}{\tau_1^n \sigma^2} \right) - r_n \geq 0, \\ & \forall n \in \mathcal{N}. \end{aligned} \quad (40)$$

Proof: Two auxiliary variables r_n and r_{nm} are introduced as in Appendix C for the proof of Proposition 2 and leading to constraints C_1 , C_7 , C_8 . Then, following the same procedure and similar algebraic manipulation as in the proof of Proposition 3, the proof is completed. ■

C. Mixed A-TDMA/HiNOMA

In the same manner we can formulate the optimization framework for the mixed protocol, where A-TDMA is used in the first hop, whereas each FAP serves its users in the second hop with HiNOMA. The formulated problem is based on the two aforementioned approaches, where each protocol imposes its own constraints depending on the hop. Therefore,

according to the previous section, the optimization problem can be expressed as:

$$\begin{aligned} \max_{\mathbf{P}, \mathbf{q}, \tau} \quad & f \\ \text{s.t.} \quad & C_1 : \sum_{m=1}^{M_n} R_{nm}(\mathbf{q}, \tau) \leq R_n(\mathbf{P}, \tau), \forall n \in \mathcal{N}, \\ & C_2 : \tau_1^n + \tau_2^n \leq 1, \forall n \in \mathcal{N}, \quad C_3 : \sum_{n=1}^N \tau_1^n \leq 1, \\ & C_4 : P_n \leq P_{\text{BBU}}, \forall n \in \mathcal{N}, \\ & C_5 : \sum_{m=1}^{M_n} q_{nm} \leq P_{\text{FAP}}, \forall m \in \mathcal{M}_n, \forall n \in \mathcal{N}, \\ & C_6 : \sum_{n=1}^N \tau_1^n P_n \leq E_{\text{BBU}}, \quad C_7 : \tau_2^n \sum_{m=1}^{M_n} q_{nm} \leq E_{\text{FAP}}, \forall n \in \mathcal{N}, \end{aligned} \quad (41)$$

In problem (41), the objective function remains the same as before. In C_1 the throughput constraint of the FAPs is expressed, since the achievable rates of the UEs cannot be greater than the throughput of their respective FAP that serves them. In C_2 , the time constraint is realized, since both hops occupy the same timeslot. The protocol uses the first part of the timeslot τ_1^n to transmit to each FAP n via A-TDMA and during the τ_2^n , FAP n can transmit via NOMA to their UEs. C_3 ensures the sum of first hop timeslots for the A-TDMA does not exceed the total timeslot duration. The rest of the constraints, i.e., C_4 - C_7 , are the total energy and power constraints used for A-TDMA and HiNOMA, as explained above.

1) *Minimum Rate:* For this problem as well, by utilizing (41), we can express the max-min problem in its epigraph form using the hypograph variable R_{\min} . Therefore, C_1 from (41) is divided into two constraints, like in the aforementioned protocols, C_{1a} and C_{1b} . The max-min problem then is formulated as:

$$\begin{aligned} \max_{\tau, \mathbf{P}, \mathbf{q}, R_{\min}} \quad & R_{\min} \\ \text{s.t.} \quad & C_{1a} : \tau_2^n B_n \log_2 \left(1 + \frac{|h_{nm}|^2 q_{nm}}{|h_{nm}|^2 \sum_{i=m+1}^{M_n} q_{ni} + \sigma^2} \right) \geq R_{\min}, \\ & \forall m \in \mathcal{M}_n, \forall n \in \mathcal{N}, \\ & C_{1b} : \tau_1^n B_0 \log_2 \left(1 + \frac{|h_n|^2 P_n}{\sigma^2} \right) \geq M_n R_{\min}, \forall n \in \mathcal{N}, \\ & (41).C_2, (41).C_3, (41).C_4, (41).C_5, (41).C_6, (41).C_7. \end{aligned} \quad (42)$$

The problem in (42) is non-convex. By following a similar procedure and utilizing the same transformations as in A-TDMA and HiNOMA for the first and second hop in respect, the problem in (42) can be formulated as an equivalent convex

problem as shown below.

$$\begin{aligned}
 & \mathbf{max}_{\mathbf{E}, \tilde{\mathbf{q}}, \tilde{\tau}_2, \tau_1, \tilde{R}_{\min}} \tilde{R}_{\min} \\
 \mathbf{s.t.} \quad & C_{1a} : -\tilde{q}_{nm} - \log(|h_{nm}|^2) + \log\left(2^{\frac{\exp(\tilde{R}_{\min} - \tilde{\tau}_2^n)}{B_n}} - 1\right) \\
 & + \log\left(\sigma^2 + |h_{nm}|^2 \sum_{i>m}^{M_n} \exp(\tilde{q}_{ni})\right) \leq 0, \forall n \in \mathcal{N}, \forall m \in \mathcal{M}_n, \\
 & C_{1b} : \tau_1^n B_0 \log_2\left(1 + \frac{|h_n|^2 E_n}{\tau_1^n \sigma^2}\right) - M_n e^{\tilde{R}_{\min}} \geq 0, \forall n \in \mathcal{N}, \\
 & C_2 : \tau_1^n + e^{\tilde{\tau}_2^n} \leq 1, \quad \forall n \in \mathcal{N}, \quad C_3 : \sum_{n=1}^N \tau_1^n \leq 1, \\
 & C_4 : E_n - \tau_1^n P_{\text{BBU}} \leq 0, \forall n \in \mathcal{N}, \\
 & C_5 : \sum_{m=1}^{M_n} e^{\tilde{q}_{nm}} \leq P_{\text{FAP}}, \forall m \in \mathcal{M}_n, \forall n \in \mathcal{N}, \\
 & C_6 : \sum_{n=1}^N E_n \leq E_{\text{BBU}}, \\
 & C_7 : e^{\tilde{\tau}_2^n} \sum_{m=1}^{M_n} e^{\tilde{q}_{nm}} \leq E_{\text{FAP}}, \forall n \in \mathcal{N}.
 \end{aligned} \tag{43}$$

2) *Proportional Fairness*: Similarly, the problem of maximizing the proportional fairness in the mixed A-TDMA/HiNOMA protocol can be expressed as:

$$\begin{aligned}
 & \mathbf{max}_{\mathbf{P}, \mathbf{q}, \tau, R_{\text{nm}}, R_n} \sum_{n=1}^N \sum_{m=1}^{M_n} \log(R_{nm}(\mathbf{q}, \tau_2)) \\
 \mathbf{s.t.} \quad & C_1 : \sum_{m=1}^{M_n} \tau_2^n B_n \log_2\left(1 + \frac{|h_{nm}|^2 q_{nm}}{|h_{nm}|^2 \sum_{i=m+1}^{M_n} q_{ni} + \sigma^2}\right) \\
 & \leq \tau_1^n B_0 \log_2\left(1 + \frac{|h_n|^2 P_n}{\sigma^2}\right), \forall n \in \mathcal{N}, \\
 & (41).C_2, (41).C_3, (41).C_4, (41).C_5, (41).C_6, (41).C_7.
 \end{aligned} \tag{44}$$

Utilizing the same set of transformations and auxiliary variables as in the cases of HiNOMA and A-TDMA above, an equivalent convex problem can be formulated for the mixed protocol, which is given below:

$$\begin{aligned}
 & \mathbf{max}_{\mathbf{E}, \tilde{\mathbf{q}}, \tilde{\tau}_2, \tau_1, \tilde{r}_{\text{nm}}, r_n} \sum_{n=1}^N \sum_{m=1}^{M_n} \tilde{r}_{nm} \\
 \mathbf{s.t.} \quad & C_1 : \sum_{m=1}^{M_n} e^{\tilde{r}_{nm}} - r_n \leq 0 \forall n \in \mathcal{N}, \\
 & (43).C_2, (43).C_3, (43).C_4, (43).C_5, (43).C_6, (43).C_7. \\
 & C_8 : \tau_1^n B_0 \log_2\left(1 + \frac{|h_n|^2 E_n}{\tau_1^n \sigma^2}\right) - r_n \geq 0, \\
 & C_9 : -\tilde{q}_{nm} - \log(|h_{nm}|^2) + \log\left(2^{\frac{\exp(\tilde{r}_{nm} - \tilde{\tau}_2^n)}{B_n}} - 1\right) \\
 & + \log\left(\sigma^2 + |h_{nm}|^2 \sum_{k>m}^{M_n} \exp(\tilde{q}_{nk})\right) \leq 0,
 \end{aligned} \tag{45}$$

V. NUMERICAL RESULTS AND DISCUSSION

In this section, Monte Carlo simulation results with 10^5 iterations are presented for the hierarchical network with the proposed protocols. Rayleigh fading is assumed for the links between the BBU and the FAPs, as well as between the FAPs and their assigned UEs, i.e., $h_n, h_{nm} \sim CN(0, 1)$. Moreover, the FAPs are assumed to transmit with a power lower by a factor of 20 compared to the transmit power of the BBU, i.e., $P_{\text{BBU}} = 20P_{\text{FAP}}$, approximately 13dB lower. This is a practical assumption, since most BBU pools have greater power capabilities than their respective FAPs in the network. Additionally, $E_{\text{BBU}}/P_{\text{BBU}} = 1$ and $E_{\text{FAP}}/P_{\text{FAP}} = 1$ are taken into account for the presented simulation results. In Fig. 2 and Fig. 3, a total of $N = 2$ FAPs are deployed, with a total of $M_1 = 2$ and $M_2 = 3$ UEs, respectively. The effect of larger number of FAPs or UEs in the network is showcased in Figs. 4-7. The SNR presented in the figures is defined as the average received SNR when the maximum allowed power is allocated at the BBU, i.e., P_{BBU} . Therefore, without loss of generality, the path loss is included in the SNR term.

For the sake of comparison, a benchmark TDMA scheme is also considered to compare the proposed protocols with conventional OMA solutions. In this TDMA scheme, the total timeslot is based on the division of the frame it two times slots and the use of one slot for each hop, similar with the case of HiNOMA. During the first timeslot, τ_1 , the BBU transmits information to the FAPs with TDMA as the protocol. In the second timeslot, τ_2 , each FAP transmits to their assigned UEs also via a TDMA protocol. Resource allocation is also optimized in the benchmark TDMA scheme in the same way as the proposed schemes for fair comparison.

In Fig. 2, the maximum minimum rate achieved by each protocol is presented versus the SNR at the BBU. It is obvious that every proposed protocol outperforms the benchmark TDMA scheme. It is notable that out of the proposed schemes, HiNOMA falls behind from the rest of the protocols, following a slight lower increasing trend with the SNR compared to A-TDMA and the Mixed protocol. More specifically, at a minimum rate of 2 bps/Hz, HiNOMA has a 3dB better performance than the benchmark TDMA. A-TDMA has a gain of 8dB from HiNOMA at the same rate, while the Mixed protocol offers approximately 10dB improvement compared to HiNOMA. The main difference between HiNOMA and the mixed protocol is the use of A-TDMA for the first hop, where due to the half-duplex operation, NOMA cannot provide as high data rates to its weakest FAPs.

Similarly, in Fig. 3, the sum rate is presented for each protocol versus the SNR. Once more, the proposed protocols outperform the benchmark TDMA with a gain of over 10dB for a sum rate of 11 bps/Hz. HiNOMA, showing worse performance than the rest of the proposed protocols, outperforms the benchmark TDMA scheme with a gain of approximately 4dB. On the other hand, mixed A-TDMA/HiNOMA prevails with a gain of 4dB from the A-TDMA protocol. The reasoning behind these results is that more TFRBs are effectively utilized with the A-TDMA and the mixed protocols, since for these protocols both hops operate at the same time for different

TABLE II
COMPARISON BETWEEN THE HIMA PROTOCOLS

Criterion	HiNOMA	A-TDMA	A-TDMA/HiNOMA
Data Rate	Achieves the lowest rates in both maximized minimum and proportional fairness	Higher data rates than HiNOMA, lower than mixed A-TDMA/HiNOMA	Highest data rate
Utilization of Resource Blocks	Spends the least TFRBs	Spends more TFRBs than HiNOMA	Spends more TFRBs than HiNOMA
Receiver Complexity	Due to SIC in both hops	Due to synchronization	Due to synchronization for the first hop and SIC for the second.

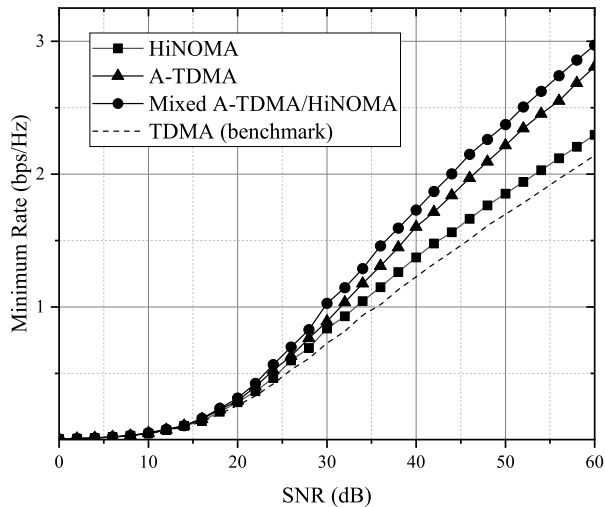


Fig. 2. Maximized R_{\min} vs SNR

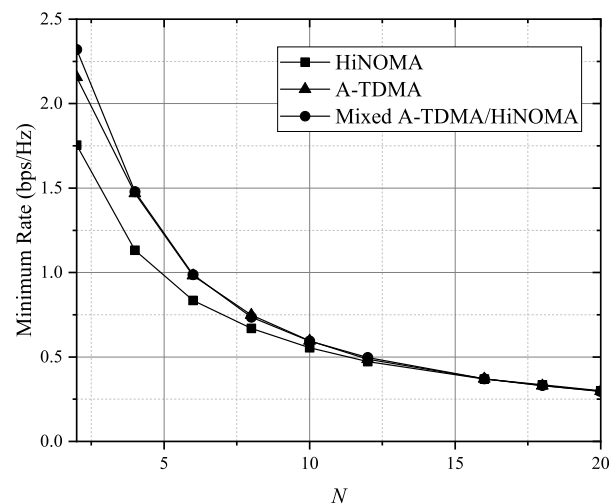


Fig. 4. Maximized R_{\min} vs total number of FAPs N with transmit SNR = 50dB.

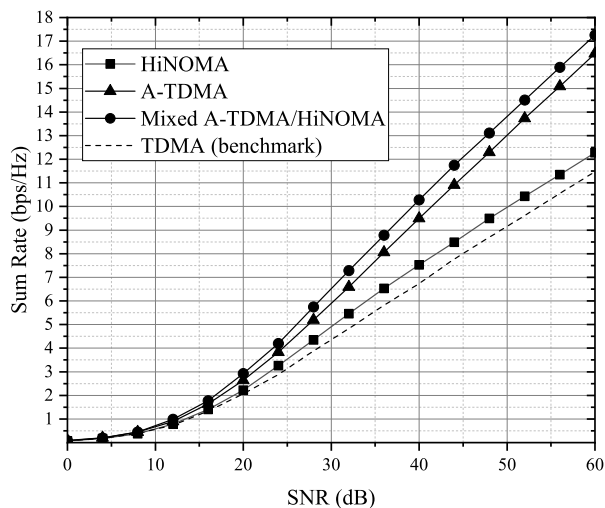


Fig. 3. Sum Rate vs SNR when proportional fairness is maximized.

FAPs. It is important to note here that this performance gain of A-TDMA and mixed A-TDMA/HiNOMA over HiNOMA is due to their ability to utilize more TFRBs in the system. Ultimately, HiMA's design is tailored to the hierarchical networks' structure for a more efficient communication between the BBU and the UEs.

In Figs. 4 and 5, the effect of the number of FAPs is

presented on the two metrics. It is assumed that each FAP served two UEs in each case. It can be observed that the maximum minimum rate of the system drops as the number of FAPs (and UEs) increases. For a low N , the mixed protocol outperforms the rest, although as the number of FAPs increases all protocols reach the same ceiling, due to the bottleneck in the first hop. On the other hand, when proportional fairness is maximized, the sum rate of A-TDMA and mixed A-TDMA/HiNOMA remain unchanged, therefore the average rate per user decreases. This is not the case for HiNOMA, which for low numbers falls behind the other two protocols, since the sum rate with HiNOMA increases as more FAPs are introduced in the system, ultimately surpassing the other two protocols before reaching a ceiling. This is attributed to NOMA's capability to offer high connectivity, compared to the OMA-based protocols that reach a bottleneck as more nodes join the network. More specifically, with more and more FAPs joining the network, the available TFRBs are increasing, closing the gap between the utilized TFRBs of the A-TDMA and the mixed protocols and HiNOMA. In this setting, HiNOMA can outperform the rest of the protocols.

Additionally, in Figs. 6 and 7, the effect of additional UEs in the system is investigated. The UEs are equally distributed among 2 FAPs in the system. The observed results are similar to the ones presented for additional FAPs in the system. However, it can be seen that the mixed protocol offers higher

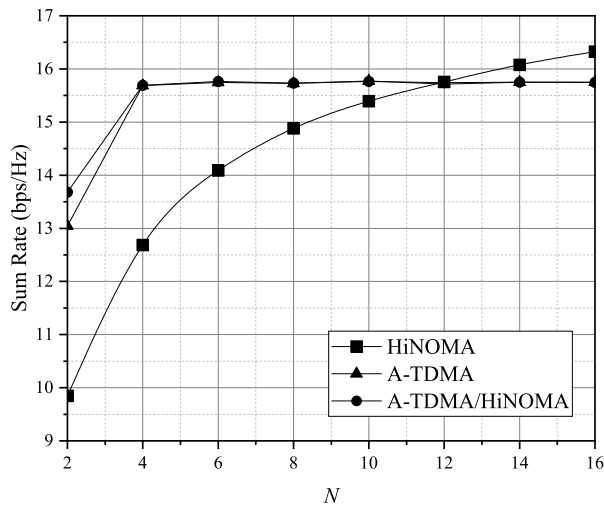


Fig. 5. Sum Rate vs total number of FAPs N when proportional fairness is maximized with transmit SNR = 50dB.

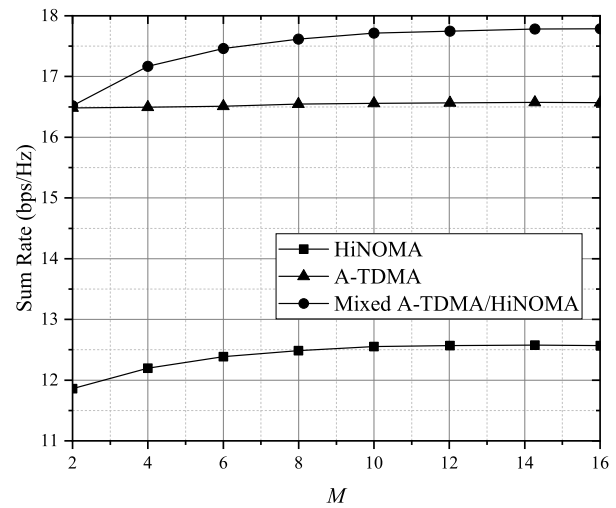


Fig. 7. Sum Rate vs total number of UEs M when proportional fairness is maximized with transmit SNR = 60dB.

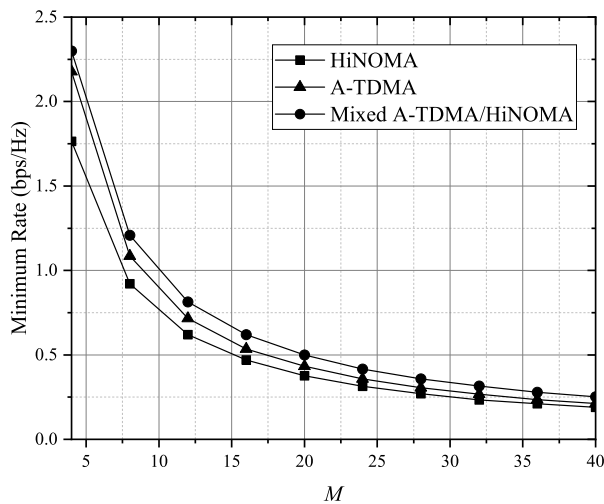


Fig. 6. Maximized R_{\min} vs total number of UEs M with transmit SNR = 50dB.

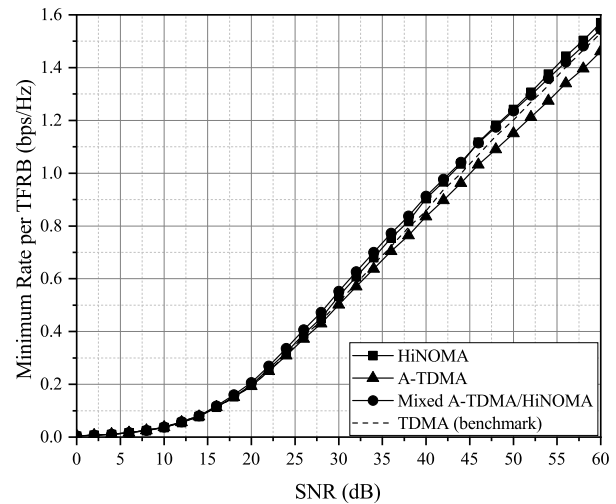


Fig. 8. Maximized R_{\min} per TFRB vs SNR.

sum rate compared to the A-TDMA in this case, reaching a higher ceiling than before. The sum rate of HiNOMA is increasing as more UEs join the system, but it offers much lower sum rate compared to the rest of the protocols, since two hops are not simultaneously operational, limiting its efficiency.

Finally, in Figs. 8 and 9, the maximized minimum rate and the sum rate when proportional fairness is optimized are presented, respectively, normalized by the utilized TFRB of each protocol. Despite the earlier findings, HiNOMA appears to prevail in this scenario, showcasing the increased spectral efficiency of NOMA against orthogonal schemes. Although in Fig 8, the mixed protocol offers greater minimum rate per TFRB for the lower SNR region. The reason behind these results is that A-TDMA and the mixed protocol employ more TFRBs, since both hops operate simultaneously. On the other hand, the benchmark TDMA outperforms the A-TDMA in this metric, since A-TDMA is designed to utilize more TFRBs to achieve much higher data rates in the same architecture.

HiNOMA's superiority is more clear in Fig. 9, followed by the mixed A-TDMA/HiNOMA. In this case, the performances of A-TDMA and TDMA are identical.

VI. CONCLUSIONS

We investigated the multiple access of a hierarchical network, proposing three novel protocols termed as Hierarchical Multiple Access (HiMA), namely HiNOMA, A-TDMA, and mixed A-TDMA/HiNOMA. Due to the varying characteristics of the many communication links in a hierarchical network and the common pool of available resources, optimization has been jointly performed for the entirety of the hierarchical network. Two different fairness metrics have been optimized, the max-min fairness and the proportional fairness. Simulation results have proven the effectiveness of the proposed protocols, since they easily outperform a similarly optimized benchmark OMA scheme. Among the proposed protocols, the mixed A-TDMA/HiNOMA offers the best performance for both metrics, while HiNOMA offers better data rate per

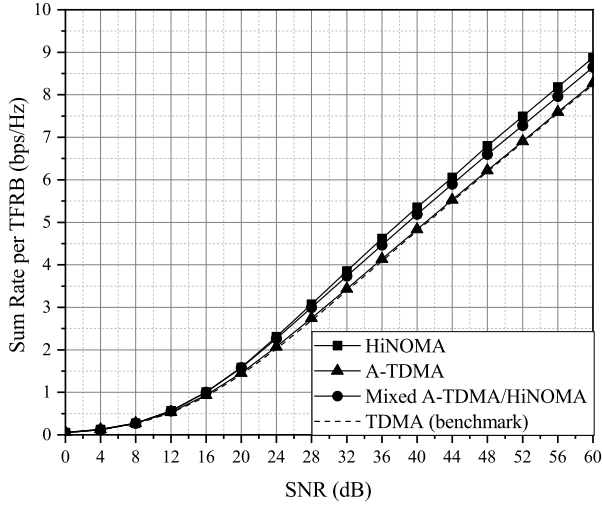


Fig. 9. Sum Rate per TFRB vs SNR when proportional fairness is maximized.

utilized TFRB. However, A-TDMA and the mixed protocol have the ability to use more TFRBs and the asynchronous protocols can benefit from it, reaching much higher data rates. Finally, the obtained results provide valuable insights about the achievable data rate and the spectral efficiency of the proposed HiMA protocols as well as their response to higher amount of connected nodes, which is crucial for optimizing RANs of the next generations of wireless access, such as C-RAN and F-RAN. We would like to note that the contribution of this work is fundamental in that it introduces the HiMA concept paving the way for more complicated schemes in the future. For example, the performance of proposed protocols can be investigated considering the use of multi-antenna APs and imperfect CSI, using the system model that was adopted in this work as benchmark.

APPENDIX A PROOF OF PROPOSITION 1

We commence by transforming the left part of constraint C_1 into a convex function. To this end, we introduce the following transformations in order to avoid products of the optimization variables from appearing in the final expressions:

$$\begin{aligned} P_n &= \exp(\tilde{P}_n), & \forall n \in \mathcal{N}, \\ q_{nm} &= \exp(\tilde{q}_{nm}), & \forall m \in \mathcal{M}_n, \forall n \in \mathcal{N}, \\ \tau_i &= \exp(\tilde{\tau}_i), & \forall i \in \{1, 2\}, \\ R_{\min} &= \exp(\tilde{R}_{\min}). \end{aligned} \quad (46)$$

The problem of (32) is formulated as

$$\begin{aligned} & \max_{\tilde{p}, \tilde{q}, \tilde{\tau}, \tilde{R}_{\min}} \exp(\tilde{R}_{\min}) \\ \text{s.t. } & C_{1a} : e^{\tilde{\tau}_2} B_n \log_2 \left(1 + \frac{|h_{nm}|^2 e^{\tilde{q}_{nm}}}{|h_{nm}|^2 \sum_{i>m}^{M_n} e^{\tilde{q}_{ni}} + \sigma^2} \right) \geq e^{\tilde{R}_{\min}}, \\ & \quad \forall m \in \mathcal{M}_n \quad \text{and} \quad \forall n \in \mathcal{N}, \\ & C_{1b} : e^{\tilde{\tau}_1} B_0 \log_2 \left(1 + \frac{|h_n|^2 e^{\tilde{P}_n}}{|h_n|^2 \sum_{i>n}^N e^{\tilde{P}_i} + \sigma^2} \right) \geq M_n e^{\tilde{R}_{\min}}, \\ & \quad \forall n \in \mathcal{N}, \\ & C_2 : e^{\tilde{\tau}_1} + e^{\tilde{\tau}_2} \leq 1, \quad C_3 : \sum_{n=1}^N e^{\tilde{P}_n} \leq P_{\text{BBU}}, \\ & C_4 : \sum_{m=1}^{M_n} e^{\tilde{q}_{nm}} \leq P_{\text{FAP}}, \quad \forall n \in \mathcal{N}, \quad C_5 : \sum_{n=1}^N e^{\tilde{\tau}_1 + \tilde{P}_n} \leq E_{\text{BBU}}, \\ & C_6 : \sum_{m=1}^{M_n} e^{\tilde{\tau}_2 + \tilde{q}_{nm}} \leq E_{\text{FAP}}, \quad \forall n \in \mathcal{N}, \\ & C_7 : \tilde{p}_n \leq \log(P_{\text{BBU}}) \quad \text{and} \quad \tilde{q}_{nm} \leq \log(P_{\text{FAP}}). \end{aligned} \quad (47)$$

Constraints $C_2, C_3, C_4, C_5, C_6,$ and C_7 are convex since they are sum of exponential functions. However, the objective function is now convex and since this is maximization problem, the objective function needs to be concave. In order to overcome this issue, the function $f' = \tilde{R}_{\min}$ is used, since the objective function is an increasing function of \tilde{R}_{\min} and its maximization is equivalent to that of f' . Finally, in their current form, C_{1a} and C_{1b} are not convex, but with some mathematical manipulations we get the following for C_{1a} :

$$\begin{aligned} & e^{\tilde{\tau}_2} B_n \log_2 \left(1 + \frac{|h_{nm}|^2 e^{\tilde{q}_{nm}}}{|h_{nm}|^2 \sum_{i=m+1}^{M_n} e^{\tilde{q}_{ni}} + \sigma^2} \right) \geq e^{\tilde{R}_{\min}}, \\ & \frac{|h_{nm}|^2 e^{\tilde{q}_{nm}}}{|h_{nm}|^2 \sum_{i=m+1}^{M_n} e^{\tilde{q}_{ni}} + \sigma^2} \geq 2^{\frac{1}{B_n} \exp(\tilde{R}_{\min} - \tilde{\tau}_2)} - 1, \\ & \log \left(\frac{|h_{nm}|^2 \sum_{i=m+1}^{M_n} e^{\tilde{q}_{ni}} + \sigma^2}{|h_{nm}|^2 e^{\tilde{q}_{nm}}} \right) + \\ & \log \left(2^{\frac{1}{B_n} \exp(\tilde{R}_{\min} - \tilde{\tau}_2)} - 1 \right) \leq 0, \\ & -\tilde{q}_{nm} - \log(|h_{nm}|^2) + \log \left(2^{\frac{1}{B_n} \exp(\tilde{R}_{\min} - \tilde{\tau}_2)} - 1 \right) \\ & + \log \left(\sigma^2 + |h_{nm}|^2 \sum_{i=m+1}^{M_n} \exp(\tilde{q}_{ni}) \right) \leq 0, \end{aligned} \quad (48)$$

$\forall m \in \mathcal{M}_n, \forall n \in \mathcal{N}$. Following the same steps for constraint C_{1b} we get By following the exact same procedure for C_1 , we get:

$$\begin{aligned} & -\tilde{P}_n - \log(|h_n|^2) + \log \left(2^{\frac{M_n}{B_0} \exp(\tilde{R}_{\min} - \tilde{\tau}_1)} - 1 \right) \\ & + \log \left(\sigma^2 + |h_n|^2 \sum_{i=n+1}^N \exp(\tilde{P}_i) \right) \leq 0, \quad \forall n \in \mathcal{N}. \end{aligned} \quad (49)$$

The first two terms of (48) are linear. The fourth term is convex as a log-sum-exp function. Finally, the third term of the left part of (48) is a function $g = \log \left(2^{\frac{1}{B_n} \exp(\tilde{R}_{\min} - \tilde{\tau}_2)} - 1 \right)$

we need to examine for convexity. By considering its Hessian matrix, which is given as

$$\mathbf{H} = \begin{bmatrix} \frac{\partial^2 g}{\partial \tilde{R}_{\min}^2} & \frac{\partial^2 g}{\partial \tilde{R}_{\min} \partial \tilde{\tau}_2} \\ \frac{\partial^2 g}{\partial \tilde{R}_{\min} \partial \tilde{\tau}_2} & \frac{\partial^2 g}{\partial \tilde{\tau}_2^2} \end{bmatrix}, \quad (50)$$

and after some algebraic manipulations it is given as

$$\mathbf{H} = \begin{bmatrix} q & -q \\ -q & q \end{bmatrix}. \quad (51)$$

It can easily be shown that \mathbf{H} has a non-zero eigenvalue that is expressed as

$$u_1 = 2q = \frac{2^z z \log(2)(2^z - z \log(2) - 1)}{(2^z - 1)^2}, \quad (52)$$

where z is defined by $z = \frac{1}{B_n} \exp(\tilde{R}_{\min} - \tilde{\tau}_2)$. Considering also that $y = 2^z - z \log(2) - 1$ is an increasing function with respect to z and when $z \rightarrow 0$, $y \rightarrow 0$, it is shown that $u_1 \geq 0$. Then, it becomes evident that the Hessian matrix of g is positive semi-definite, due to the fact that the eigenvalues of the matrix are non-negative. As a result, constraint C_{1a} is proven to be convex. Following the exact same procedure, it can easily be shown that constraint C_{1b} is convex as well.

Therefore, the non-convex problem of (32) can be transformed to an equivalent convex problem and the proof is completed.

APPENDIX B

DUAL DECOMPOSITION AND SOLUTION OF (33)

First, we obtain the Lagrangian of the problem, which is given by (53) at the top of the next page, with λ_i being the Lagrange multipliers (LMs). Using the Karush-Kuhn-Tucker (KKT) conditions for a fixed set of LMs the subproblems are solved in each iteration. Then, the LMs are updated via a subgradient method [33], [34] that offers a theoretical complexity of $O(1/\epsilon^2)$ iterations to find the ϵ -suboptimal point. The KKT conditions are given in (54), (55), (56), (57), (58),

$$\begin{aligned} \frac{dL}{d\tilde{R}_{\min}} = 0 \Leftrightarrow & 1 - \sum_{n=1}^N \sum_{m=1}^{M_n} \lambda_{nm} \times \\ & \frac{2^{\frac{\exp(\tilde{R}_{\min} - \tilde{\tau}_2)}{B_n}} \exp(\tilde{R}_{\min} - \tilde{\tau}_2) \log(2)}{2^{\frac{\exp(\tilde{R}_{\min} - \tilde{\tau}_2)}{B_n}} - 1} \frac{1}{B_n} - \sum_{n=1}^N \lambda_{M+n} \times \\ & \frac{2^{\frac{M_n \exp(\tilde{R}_{\min} - \tilde{\tau}_1)}{B_0}} \exp(\tilde{R}_{\min} - \tilde{\tau}_1) \frac{M_n \log(2)}{B_0}}{2^{\frac{M_n \exp(\tilde{R}_{\min} - \tilde{\tau}_1)}{B_0}} - 1} = 0, \quad (54) \end{aligned}$$

$$\begin{aligned} \frac{dL}{d\tilde{p}_n} = 0 \Leftrightarrow & \lambda_{M+n} - \sum_{j=1}^{n-1} \lambda_{M+j} \frac{|h_j|^2 \exp(\tilde{p}_n)}{\sigma^2 + |h_j|^2 \sum_{i=j+1}^N \exp(\tilde{p}_i)} \\ & - \lambda_{M+N+2} \exp(\tilde{p}_n) - \lambda_{M+2N+3} \exp(\tilde{p}_n + \tilde{\tau}_1) = 0, \quad (55) \end{aligned}$$

$$\begin{aligned} \frac{dL}{d\tilde{q}_{nm}} = 0 \Leftrightarrow & \lambda_{nm} - \sum_{j=1}^{m-1} \lambda_{nj} \frac{|h_{nj}|^2 \exp(\tilde{q}_{nm})}{\sigma^2 + |h_{nj}|^2 \sum_{i=j+1}^{M_n} \exp(\tilde{q}_{ni})} \\ & - \lambda_{M+N+n+2} \exp(\tilde{q}_{nm}) - \lambda_{M+2N+3+n} \exp(\tilde{q}_{nm} + \tilde{\tau}_2) = 0, \quad (56) \end{aligned}$$

$$\begin{aligned} \frac{dL}{d\tilde{\tau}_1} = 0 \Leftrightarrow & \\ & - \sum_{n=1}^N \lambda_{M+n} \frac{2^{\frac{M_n \exp(\tilde{R}_{\min} - \tilde{\tau}_1)}{B_0}} \exp(\tilde{R}_{\min} - \tilde{\tau}_1) \frac{M_n \log(2)}{B_0}}{1 - 2^{\frac{M_n \exp(\tilde{R}_{\min} - \tilde{\tau}_1)}{B_0}}} \\ & - \lambda_{M+2N+3} \sum_{n=1}^N \exp(\tilde{p}_n + \tilde{\tau}_1) - \lambda_{M+N+1} \exp(\tilde{\tau}_1) = 0, \quad (57) \end{aligned}$$

$$\begin{aligned} \frac{dL}{d\tilde{\tau}_2} = 0 \Leftrightarrow & \\ & - \sum_{n=1}^N \frac{2^{\frac{\exp(\tilde{R}_{\min} - \tilde{\tau}_2)}{B_n}} \exp(\tilde{R}_{\min} - \tilde{\tau}_2) \log(2)}{1 - 2^{\frac{\exp(\tilde{R}_{\min} - \tilde{\tau}_2)}{B_n}}} \frac{1}{B_n} \sum_{m=1}^{M_n} \lambda_{nm} \\ & - \sum_{n=1}^N \lambda_{M+2N+3+n} \sum_{m=1}^{M_n} e^{\tilde{q}_{nm} + \tilde{\tau}_2} - \lambda_{M+N+1} \exp(\tilde{\tau}_2) = 0. \quad (58) \end{aligned}$$

APPENDIX C PROOF OF PROPOSITION 2

The procedure is similar to that of the proof of Proposition 1; we commence by transforming the objective function into a concave function, due to it being a maximization problem. To this end, we introduce two auxiliary variables

$$r_{nm} \leq R_{nm} \quad \text{and} \quad r_n \leq R_n. \quad (59)$$

The problem of (34) is formulated as

$$\begin{aligned} \max_{\tau, \mathbf{p}, \mathbf{q}, \mathbf{r}} & \sum_{n=1}^N \sum_{m=1}^{M_n} \log(r_{nm}) \\ \text{s.t.} & C_1: \sum_{m=1}^{M_n} r_{nm} \leq r_n, \quad \forall n \in \mathcal{N}, \quad (60) \\ & (28).C_2, (28).C_3, (28).C_4, (28).C_5, (28).C_6, (28).C_7, \\ & C_8: r_{nm} \leq R_{nm}, \quad \forall m \in \mathcal{M}_n, \forall n \in \mathcal{N}, \\ & C_9: r_n \leq R_n, \quad \forall n \in \mathcal{N}. \end{aligned}$$

There are two new constraints that need to be satisfied due to the use of (59). The problem of (60) is still non-convex. In order to continue our proof, we introduce the following transformations:

$$\begin{aligned} P_n &= \exp(\tilde{P}_n), & \forall n \in \mathcal{N}, \\ q_{nm} &= \exp(\tilde{q}_{nm}), & \forall m \in \mathcal{M}_n, \forall n \in \mathcal{N}, \\ \tau_i &= \exp(\tilde{\tau}_i), & \forall i \in \{1, 2\}, \\ r_{nm} &= \exp(\tilde{r}_{nm}), & \forall m \in \mathcal{M}_n, \forall n \in \mathcal{N}, \\ r_n &= \exp(\tilde{r}_n), & \forall n \in \mathcal{N}. \quad (61) \end{aligned}$$

The first 3 transformations are needed to transform constraints C_5 and C_6 into convex functions. However, due to the new constraints introduced due to (59), i.e., C_8 and C_9 , the introduction of \tilde{r}_{nm} and \tilde{r}_n is also needed. Constraints $C_1 - C_7$ are convex after (61). C_8 and C_9 are then expressed in the same manner as (48) and it can easily be shown for both that they are convex functions. Therefore, the proof is completed.

$$\begin{aligned}
 L = & \tilde{R}_{\min} - \sum_{n=1}^N \sum_{m=1}^{M_n} \lambda_{nm} \left[-\tilde{q}_{nm} - 2 \log(|h_{nm}|) + \log \left(2^{\frac{\exp(\tilde{R}_{\min} - \tilde{\tau}_2)}{B_n}} - 1 \right) + \log \left(\sigma^2 + |h_{nm}|^2 \sum_{i>m}^{M_n} \exp(\tilde{q}_{ni}) \right) \right] \\
 & - \sum_{n=1}^N \lambda_{M+n} \left[-\tilde{p}_n - 2 \log(|h_n|) + \log \left(2^{\frac{M_n \exp(\tilde{R}_{\min} - \tilde{\tau}_1)}{B_0}} - 1 \right) + \log \left(\sigma^2 + |h_n|^2 \sum_{i>n}^N \exp(\tilde{p}_i) \right) \right] \\
 & - \lambda_{M+N+1} (\exp(\tilde{\tau}_1) + \exp(\tilde{\tau}_2) - 1) - \lambda_{M+N+2} \sum_{n=1}^N (\exp(\tilde{p}_n) - P_{\text{BBU}}) - \sum_{n=1}^N \lambda_{M+N+2+n} \sum_{m=1}^{M_n} (\exp(\tilde{q}_{nm}) - P_{\text{FAP}}) \\
 & - \lambda_{M+2N+3} \sum_{n=1}^N (\exp(\tilde{p}_n + \tilde{\tau}_1) - E_{\text{BBU}}) - \sum_{n=1}^N \lambda_{M+2N+3+n} \sum_{m=1}^{M_n} (\exp(\tilde{q}_{nm} + \tilde{\tau}_2) - E_{\text{FAP}})
 \end{aligned} \tag{53}$$

REFERENCES

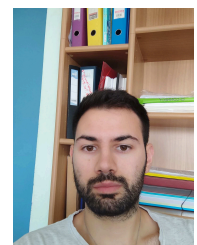
- [1] Z. Ding, X. Lei, G. K. Karagiannidis, R. Schober, J. Yuan, and V. K. Bhargava, "A Survey on Non-Orthogonal Multiple Access for 5G Networks: Research Challenges and Future Trends," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 10, pp. 2181–2195, Oct. 2017.
- [2] V. K. Papanikolaou, P. D. Diamantoulakis, and G. K. Karagiannidis, "User Grouping for Hybrid VLC/RF Networks With NOMA: A Coalitional Game Approach," *IEEE Access*, vol. 7, pp. 103 299–103 309, 2019.
- [3] L. Wei, R. Q. Hu, Y. Qian, and G. Wu, "Key elements to enable millimeter wave communications for 5G wireless systems," *IEEE Wireless Commun.*, vol. 21, no. 6, pp. 136–143, Dec. 2014.
- [4] Y. Saito, A. Benjebbour, Y. Kishiyama, and T. Nakamura, "System-level performance evaluation of downlink non-orthogonal multiple access (NOMA)," in *2013 IEEE 24th Annu. Int. Symp. Pers. Indoor, and Mobile Radio Commun. (PIMRC)*, London, UK, Sep. 2013, pp. 611–615.
- [5] S. Timotheou and I. Krikididis, "Fairness for Non-Orthogonal Multiple Access in 5G Systems," *IEEE Signal Process. Lett.*, vol. 22, no. 10, pp. 1647–1651, Oct. 2015.
- [6] X. Gu, X. Ji, Z. Ding, W. Wu, and M. Peng, "Outage Probability Analysis of Non-Orthogonal Multiple Access in Cloud Radio Access Networks," *IEEE Commun. Lett.*, vol. 22, no. 1, pp. 149–152, Jan. 2018.
- [7] K. N. Pappi, P. D. Diamantoulakis, and G. K. Karagiannidis, "Distributed Uplink-NOMA for Cloud Radio Access Networks," *IEEE Commun. Lett.*, vol. 21, no. 10, pp. 2274–2277, Oct. 2017.
- [8] P. D. Diamantoulakis and G. K. Karagiannidis, "Performance Analysis of Distributed Uplink NOMA," *IEEE Commun. Lett.*, vol. 25, no. 3, pp. 788–792, Mar. 2021.
- [9] I. Randrianantenaina, M. Kaneko, H. Dahrouj, H. ElSawy, and M. Alouini, "Interference Management in NOMA-Based Fog-Radio Access Networks via Scheduling and Power Allocation," *IEEE Trans. Commun.*, vol. 68, no. 8, pp. 5056–5071, Aug. 2020.
- [10] R. Rai, H. Zhu, and J. Wang, "Resource scheduling in non-orthogonal multiple access (NOMA) based cloud-RAN systems," in *2017 IEEE 8th Annu. Ubiquitous Computing Electron. Mobile Commun. Conf. (UEMCON)*, New York, NY, USA, Oct. 2017, pp. 418–422.
- [11] D. Boviz, C. S. Chen, and S. Yang, "Effective Design of Multi-User Reception and Fronthaul Rate Allocation in 5G Cloud RAN," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 8, pp. 1825–1836, Aug. 2017.
- [12] R. Singh, H. Zhu, and J. Wang, "Performance of Non-Orthogonal Multiple Access (NOMA) in a C-RAN System," in *2017 IEEE 86th Veh. Technol. Conf. (VTC-Fall)*, Toronto, ON, Canada, Sep. 2017, pp. 1–5.
- [13] H. Q. Tran, P. Q. Truong, C. V. Phan, and Q. Vien, "On the Energy Efficiency of NOMA for Wireless Backhaul in Multi-Tier Heterogeneous CRAN," in *2017 Int. Conf. Recent Advances Signal Process. Telecommun. Comput. (SigTelCom)*, Da Nang, Vietnam, Jan. 2017, pp. 229–234.
- [14] F. Zhou, Y. Wu, R. Q. Hu, Y. Wang, and K. K. Wong, "Energy-Efficient NOMA Enabled Heterogeneous Cloud Radio Access Networks," *IEEE Netw.*, vol. 32, no. 2, pp. 152–160, Mar. 2018.
- [15] W. Hao, Z. Chu, F. Zhou, S. Yang, G. Sun, and K. Wong, "Green Communication for NOMA-Based CRAN," *IEEE Internet Things J.*, vol. 6, no. 1, pp. 666–678, Feb. 2019.
- [16] S. Lee, S. Park, and I. Lee, "NOMA Systems With Content-Centric Multicast Transmission for C-RAN," *IEEE Wireless Commun. Lett.*, vol. 7, no. 5, pp. 828–831, Oct. 2018.
- [17] X. Cao, M. Peng, and Z. Ding, "A Game-Theoretic Approach of Resource Allocation in NOMA-Based Fog Radio Access Networks," in *2019 IEEE 90th Veh. Technol. Conf. (VTC2019-Fall)*, Honolulu, HI, USA, Sep. 2019, pp. 1–5.
- [18] B. Liu and M. Peng, "Joint Resource Block-Power Allocation for NOMA-Enabled Fog Radio Access Networks," in *ICC 2019 - 2019 IEEE Int. Conf. Commun. (ICC)*, Shanghai, China, May 2019, pp. 1–6.
- [19] R. Rai, H. Zhu, and J. Wang, "Performance Analysis of NOMA Enabled Fog Radio Access Networks," *IEEE Trans. Commun.*, vol. 69, no. 1, pp. 382–397, Jan. 2021.
- [20] X. Wen, H. Zhang, H. Zhang, and F. Fang, "Interference Pricing Resource Allocation and User-Subchannel Matching for NOMA Hierarchy Fog Networks," *IEEE J. Sel. Topics Signal Process.*, vol. 13, no. 3, pp. 467–479, Jun. 2019.
- [21] G. Liu, X. Chen, Z. Ding, Z. Ma, and F. R. Yu, "Hybrid Half-Duplex/Full-Duplex Cooperative Non-Orthogonal Multiple Access With Transmit Power Adaptation," *IEEE Trans. Wireless Commun.*, vol. 17, no. 1, pp. 506–519, Jan. 2018.
- [22] Z. Ding, M. Peng, and H. V. Poor, "Cooperative non-orthogonal multiple access in 5G systems," *IEEE Commun. Lett.*, vol. 19, no. 8, pp. 1462–1465, Aug. 2015.
- [23] L. Lv, J. Chen, and Q. Ni, "Cooperative Non-Orthogonal Multiple Access in Cognitive Radio," *IEEE Commun. Lett.*, vol. 20, no. 10, pp. 2059–2062, Oct. 2016.
- [24] X. Yue, Y. Liu, S. Kang, A. Nallanathan, and Z. Ding, "Exploiting Full/Half-Duplex User Relaying in NOMA Systems," *IEEE Trans. Commun.*, vol. 66, no. 2, pp. 560–575, Feb. 2018.
- [25] M. B. Shahab and S. Y. Shin, "Time Shared Half/Full-Duplex Cooperative NOMA with Clustered Cell Edge Users," *IEEE Commun. Lett.*, vol. 22, no. 9, pp. 1794–1797, Sep. 2018.
- [26] Q. Y. Liau and C. Y. Leow, "Successive User Relaying in Cooperative NOMA System," *IEEE Wireless Commun. Lett.*, vol. 8, no. 3, pp. 921–924, Jun. 2019.
- [27] C. Li, D. Guo, K. Guo, Y. Qin, and R. Xu, "Outage Performance of Partial Relay Selection in Underlay CR-NOMA Networks," in *2019 28th Wireless Opt. Commun. Conf. (WOCC)*. Beijing, China: IEEE, May 2019, pp. 1–5.
- [28] J. Gora and S. Redana, "In-band and out-band relaying configurations for dual-carrier LTE-advanced system," in *2011 IEEE 22nd Int. Symp. Personal, Indoor, and Mobile Radio Communications*, Toronto, ON, Canada, Sep. 2011, pp. 1820–1824.
- [29] T. M. de Moraes, M. D. Nisar, A. A. Gonzalez, and E. Seidel, "Resource allocation in relay enhanced LTE-Advanced networks," *EURASIP J. Wireless Commun. Netw.*, vol. 2012, no. 1, p. 364, Dec. 2012.
- [30] M. Vaezi, R. Schober, Z. Ding, and H. V. Poor, "Non-Orthogonal Multiple Access: Common Myths and Critical Questions," *IEEE Wireless Commun.*, vol. 26, pp. 174–180, Oct. 2019.
- [31] H. Tabassum, E. Hossain, and J. Hossain, "Modeling and analysis of uplink non-orthogonal multiple access in large-scale cellular networks using poisson cluster processes," *IEEE Trans. Commun.*, vol. 65, no. 8, pp. 3555–3570, Aug. 2017.
- [32] V. K. Papanikolaou, P. D. Diamantoulakis, P. C. Sofotasios, S. Muhaidat, and G. K. Karagiannidis, "On Optimal Resource Allocation for Hybrid VLC/RF Networks With Common Backhaul," *IEEE Trans. on Cogn. Commun. Netw.*, vol. 6, no. 1, pp. 352–365, Mar. 2020.
- [33] S. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.

- [34] S. Boyd, "Subgradient methods," *Lecture notes of EE364b, Stanford University*, 2003.



Vasilis K. Papanikolaou was born in Kavala, Greece in 1995. He received the Diploma Degree (5 years) in Electrical and Computer Engineering from the Aristotle University of Thessaloniki (AUTH), Greece, in 2018, where is currently pursuing his PhD with the Department of Electrical and Computer Engineering. He was a visitor researcher at Lancaster University, UK and at Khalifa University, Abu Dhabi, UAE. In 2018, he received the IEEE Student Travel Grant Award for IEEE WCNC 2018.

His research interests include visible light communications (VLC), non-orthogonal multiple access (NOMA), optimization theory, and game theory. He has served as a reviewer in various IEEE journals and conferences. He was also an Exemplary Reviewer of IEEE Wireless Communications Letters in 2019 (top 3% of reviewers).



Nikos A. Mitsiou was born in Achinos, Phthiotis, Greece. He received the Diploma Degree (5 years) in Electrical and Computer Engineering from the Aristotle University of Thessaloniki (AUTH), Greece, in 2021, where he is currently pursuing his PhD with the Department of Electrical and Computer Engineering. He is a member of the Wireless and Communications & Information Processing (WCIP) Group. His research interests include non-orthogonal multiple access (NOMA), optimization theory and resource allocation.



Panagiotis D. Diamantoulakis (SM IEEE) received the Diploma (five years) and PhD from the Department of Electrical and Computer Engineering, Aristotle University of Thessaloniki (AUTH), Greece, in 2012 and 2017, respectively. Since 2017, he works as a Post-doctoral Fellow in Wireless Communications & Information Processing (WCIP) Group at AUTH and, since 2021, he is also a visiting Assistant Professor in the Key Lab of Information Coding and Transmission at Southwest Jiaotong University (SWJTU), China. From 2018 to 2020, he also

worked as visiting Post-doctoral Researcher in the Key Lab of Information Coding and Transmission at SWJTU and in the Institute for Digital Communications (IDC) of the Telecommunications Laboratory (LNT) at Friedrich-Alexander-Universität Erlangen-Nürnberg (FAU), Germany. His current research interests include resource allocation in wireless communications, optimization theory and applications in wireless networks and smart grids, game theory, wireless power transfer and optical wireless communications. He is a Working Group Member in the Newfocus COST Action "European Network on Future Generation Optical Wireless Communication Technologies". He serves as an Associate Editor for IEEE Wireless Communications Letters, IEEE Open Journal of the Communications Society, Physical Communications (Elsevier), and Frontiers in Communications and Networks. He was also an Exemplary Reviewer of IEEE Communications Letters in 2014 and IEEE Transactions on Communications in 2017 and 2019 (top 3% of reviewers).



Zhiguo Ding (Fellow, IEEE) received the B.Eng. degree in electrical engineering from the Beijing University of Posts and Telecommunications, Beijing, China, in 2000, and the Ph.D. degree in electrical engineering from Imperial College London, London, U.K., in 2005. From 2005 to 2018, he was with Queen's University Belfast, Belfast, U.K., Imperial College, Newcastle University, Newcastle upon Tyne, U.K., and Lancaster University, Lancashire, U.K. Since 2018, he has been a Professor of communications with the University of Manchester,

Manchester, U.K. From 2012 to 2020, he was an Academic Visitor with Princeton University, Princeton, NJ, USA. His research interests include 5G networks, game theory, cooperative and energy harvesting networks, and statistical signal processing. He is the Area Editor of the IEEE Open Journal of the Communications Society, the Editor of IEEE Transactions on Communications, IEEE Transactions on Vehicular Technology, and Journal of Wireless Communications and Mobile Computing, and from 2013 to 2016, he was the Editor of the IEEE Wireless Communication Letters, IEEE Communication Letters. He was the recipient of the Best Paper Award in IET ICWMC-2009 and IEEE WCSP-2014, EU Marie Curie Fellowship 2012–2014, Top IEEE TVT Editor 2017, IEEE Heinrich Hertz Award 2018, IEEE Jack Neubauer Memorial Award 2018, IEEE Best Signal Processing Letter Award 2018, and the Web of Science Highly Cited Researcher 2019.



George K. Karagiannidis (M'96-SM'03-F'14) was born in Pithagorion, Samos Island, Greece. He received the University Diploma (5 years) and PhD degree, both in electrical and computer engineering from the University of Patras, in 1987 and 1999, respectively. From 2000 to 2004, he was a Senior Researcher at the Institute for Space Applications and Remote Sensing, National Observatory of Athens, Greece. In June 2004, he joined the faculty of Aristotle University of Thessaloniki, Greece where he is currently Professor in the Electrical

& Computer Engineering Dept. and Head of Wireless Communications & Information Processing (WCIP) Group. He is also Honorary Professor at South West Jiaotong University, Chengdu, China. His research interests are in the broad area of Digital Communications Systems and Signal processing, with emphasis on Wireless Communications, Optical Wireless Communications, Wireless Power Transfer and Applications and Communications & Signal Processing for Biomedical Engineering. Dr. Karagiannidis has been involved as General Chair, Technical Program Chair and member of Technical Program Committees in several IEEE and non-IEEE conferences. In the past, he was Editor in several IEEE journals and from 2012 to 2015 he was the Editor-in Chief of IEEE Communications Letters. Currently, he serves as Associate Editor-in Chief of IEEE Open Journal of Communications Society. Dr. Karagiannidis is one of the highly-cited authors across all areas of Electrical Engineering, recognized from Clarivate Analytics as Web-of-Science Highly-Cited Researcher in the six consecutive years 2015-2020.