

# Energy Efficiency Optimization for PSOAM Mode-Groups based MIMO-NOMA Systems

Yan Song, Jie Tang, *Senior Member, IEEE*, Chuting Lin, Wanmei Feng, *Student Member, IEEE*, Zhen Chen, *Member, IEEE*, Daniel Ka Chun So, *Senior Member, IEEE*, and Kai-Kit Wong, *Fellow, IEEE*

**Abstract**—Plane spiral orbital angular momentum (PSOAM) mode-groups (MGs) and multiple-input multiple-output non-orthogonal multiple access (MIMO-NOMA) serve as two emerging techniques for achieving high spectral efficiency (SE) in the next-generation networks. In this paper, a PSOAM MGs based multi-user MIMO-NOMA system is studied, where the base station transmits data to users by utilizing the generated PSOAM beams. For such scenario, the interference between users in different PSOAM mode groups can be avoided, which leads to a significant performance enhancement. We aim to maximize the energy efficiency (EE) of the system subject to the constraints of the total transmission power and the minimum data rate. This designed optimization problem is non-convex owing to the interference among users, and hence is quite difficult to tackle directly. To solve this issue, we develop a dual layer resource allocation algorithm where the bisection method is exploited in the outer layer to obtain the optimal EE and a resource distributed iterative algorithm is exploited in the inner layer to optimize the transmit power. Besides, an alternative resource allocation algorithm with Deep Belief Networks (DBN) is proposed to cope with the requirement for low computational complexity. Simulation results verify the theoretical findings and demonstrate the proposed algorithms on the PSOAM MGs based MIMO-NOMA system can obtain a better performance comparing to the conventional MIMO-NOMA system in terms of EE.

**Index Terms**—Energy efficiency (EE), plane spiral orbital angular momentum (PSOAM), non-orthogonal multiple access (NOMA), mode group, resource allocation.

This paper was presented in part at the EAI 3rd Artif. Intell. Commun. Netw. (AICON 2021), Xining, 2021. This work has been supported in part by Key Research and Development Project of Guangdong Province under Grant 2019B010156003, in part by National Key Research and Development Project under Grant 2019YFB1804100, in part by the National Natural Science Foundation of China under Grant 61971194, in part by the National Science Foundation of Guangdong Province under Grant 2019A1515011607, in part by the Research Fund Program of Guangdong Key Laboratory of Aerospace Communication and Networking Technology under Grant 2018B030322004, in part by the Special Project for Guangxi Science and Technology Bases and Talents under Grant AC22035089, AD21075054 and in part by the Open Research Fund of National Mobile Communications Research Laboratory, Southeast University under Grant 2019D06. (*Corresponding author: Jie Tang.*)

Y. Song, C. Lin, W. Feng and Z. Chen are with the school of Electronic and Information Engineering, South China University of Technology, Guangzhou 510641, China (e-mail: eesongyan1222@mail.scut.edu.cn, 13929316612@163.com, eewmfeng@mail.scut.edu.cn, chenz@scut.edu.cn).

J. Tang is with the School of Electronic and Information Engineering, South China University of Technology, Guangzhou, China, and also with the National Mobile Communications Research Laboratory, Southeast University, Nanjing 210096, China (e-mail: eejtang@scut.edu.cn).

D. K. C. So is with the School of Electrical and Electronic Engineering, University of Manchester, Manchester M13 9PL, U.K. (e-mail: d.so@manchester.ac.uk).

K.-K. Wong is with the Department of Electronic and Electrical Engineering, University College London, London WC1E 6BT, U.K. (e-mail: kai-kit.wong@ucl.ac.uk).

## I. INTRODUCTION

THE rapid development of Internet-of-Things (IoTs) applications has caused the exponential growth of wireless devices. Consequently, the sixth generation (6G) wireless networks face particular challenges to meet the further requirements in terms of reliable data connectivity and ultra-high data-rate. In addition, the data rates of devices are severely limited by the insufficient spectrum resources. These trends make spectral efficiency (SE) the main performance indicator of mobile communication networks. On the other hand, massive number of connected devices also leads to enormous energy consumption, and thus energy efficiency (EE) has become an important and global topic from both environmental and economic reasons [1].

Recently, the orbital angular momentum (OAM) technique is regarded as one of the key techniques in the next-generation networks [2]. Apart from the traditional multiplexing schemes, mode division multiplexing using OAM exploits a new degree of freedom for improving the SE owing to its orthogonality, thus it is capable of meeting the requirements of high data rate with a reliable bit error rate (BER) [3]. OAM was discovered by Allen *et al.* from Laguerre-Gaussian (LG) light waves in 1992. As the wireless communication develops rapidly, the OAM technology was first implemented in the electromagnetic (EM) waves, which could solve the radio-band congestion problem by encoding the information in an OAM state [4]. To further explore the unique properties of OAM waves under the multipath transmission conditions, the authors in [5] proposed a hybrid orthogonal division multiplexing (HODM) scheme, which could achieve high system throughput and mitigate the multipath interference when used in conjunction with orthogonal frequency division multiplexing. In [6], by applying the multiple input multiple output (MIMO) technique, a point-to-point OAM-MIMO system is proposed to address the problem of signal-to-noise ratio (SNR) degradation and the limitation of high OAM modes. It has been proven that this system can achieve over 200 Gbit/s data transmission with 21 streams at a distance of 10 meters. By considering the co-mode interference, the authors in [7] developed the concentric uniform circular arrays (UCAs) based radio vortex wireless communication model to increase the system capacity by using the optimal resource allocation schemes. To further improve the performance of the OAM-based communication systems in long-distance transmission, an OAM index modulation transmission scheme for long distance had been established in [8], where different OAM modes were mapped to the

index domain. In addition, the authors in [9] proposed an OAM spatial modulation (OAM-SM) system and analysed the capacity, EE, average bit error probability and the robust to path-loss attenuations, which outperformed the OAM-MIMO system at long distance transmission.

However, the main practical challenge for applying such technology into the electromagnetic (EM) field is the beam divergence and phase singularity caused by the OAM modes and long-distance transmission. With the increase of OAM modes number and the transmission distance, the dark zone in the center of the OAM wave becomes larger. Consequently, the receiver size will become extremely large using the conventional whole circle receiving scheme, which leads to large receiver form factor and high hardware complexity. To solve the problem, the authors in [10] proposed the partial arc sampling receiving (PASR), which could improve the SE with a low computational complexity and compact size. Moreover, S. Zheng *et al.* proposed a new form of OAM waves called plane spiral orbital angular momentum (PSOAM), which propagated along the transverse plane intelligently, and thereby avoiding the aforementioned issues of phase singularity and the diversity [11]. Subsequently, the authors proposed a PSOAM waves-based MIMO system to enhance the capacity gain and reduce the spatial correlations for the LOS channel by utilizing the diversity of different OAM waves [12]. Meanwhile, the authors in [13] had studied and proven that the PSOAM based MIMO system achieved better system performance in the NLOS channel than in the LOS channel. In addition, the authors further analysed the PSOAM beams and put forward the concept of PSOAM mode-groups (MGs), which had the promising prospect in SM-MIMO, radar detection and smart antenna. More importantly, the PSOAM MGs wave is a structural radio beam, which can manipulate its phase distribution and intensity, and thus is a feasible way for the PSOAM MGs beamforming method thanks to its directional gain and vorticity [14]. They also proved that the beamforming based PSOAM MGs was superior to the MIMO beamforming method due to its simplicity. In contrast with the scheme of OAM shift keying, PSOAM MGs can be applied in physical encryption, and thus the authors in [15] proposed a low probability of intercept system by mapping symbols to different PSOAM MGs beams, which achieved secure communication within the practical SNR ranges. To demonstrate the performance of PSOAM MGs, the authors in [16] applied PSOAM MGs based MIMO techniques into a single-user system, where the PASR method was adopted to demultiplex the PSOAM MGs data streams. Since the PSOAM MGs can enhance the SNR and decrease the spatial correction of sub-channels, the BER and SE performance of the proposed system outperforms the existing MIMO systems. Moreover, a partial arc transmitting scheme is developed in [17] to generate an OAM MG with high equivalent OAM order. This work was extended to realize a high-purity eight-mode PSOAM antenna using a rotating parabolic reflector and a coaxial resonator group [18]. The authors in [19] have further studied the peculiarity of PSOAM MGs and demonstrated the intensity patterns and the phase variation, which offered a potential application prospect for beamforming and radar imaging.

On the other hand, non-orthogonal multiple access (NOMA) is regarded as a key technique to enhance the SE in the fifth generation (5G) and beyond 5G (B5G) mobile networks [20]–[22]. It can simultaneously serve a large amount of users with the same physical resource via superposition coding (SC), where different users can be distinguished with different power levels and successive interference cancellation (SIC) is used to cancel the interferences among multiple users [23]. In particular, NOMA has been widely used in MISO and MIMO systems since it can obtain superior performance in terms of EE and SE in comparison with the orthogonal multiple access (OMA) system [24]–[26]. In [24], an EE maximization problem was investigated in a NOMA network, where the resource allocation and the time switching factors were jointly considered and a dual-layer algorithm using Dinkelbach method was proposed to tackle the non-convex problem. In [25], a power minimization problem was studied in a multi-user MIMO-NOMA network, where the beamforming vectors and power allocation were optimized, and a closed-form solution for beamforming vectors were obtained. Moreover, to overcome the limitation of beamspace MIMO, the authors in [26] proposed a beamspace MIMO-NOMA system that the supported users number was larger than the RF chains number. Simulation results demonstrated that it achieved better EE performance compared with the existing beamspace MIMO systems. In [27], by integrating different OAM modes and NOMA techniques, an efficient downlink NOMA-OAM-MIMO scenario was presented to further enhance the system capacity.

#### A. Motivation and Contributions

NOMA can simultaneously serve a large amount of users with the same physical resource via superposition coding, and thus can significantly improve the system performance. Due to the increasing demand of high data-rate communication, utilizing only NOMA technique does not meet the ever-growing needs of spectrum resources in the next-generation wireless communication systems. PSOAM mode-groups technique takes a different approach to improve the spectral efficiency by introducing the orthogonality in azimuthal domain. Compared with the conventional OAM technique, PSOAM can avoid the issues of phase singularity, and thus the coaxial propagation for different OAM modes can be achieved easily [11]. Therefore, the two emerging techniques can reinforce each other for achieving high spectral efficiency. Motivated by this, this work integrates PSOAM technique into MIMO-NOMA to further enhance the system performance. Specifically, in contrast to previous literature which analysed SE or EE in single-user OAM systems [9], single-user PSOAM MGs systems [16] and multi-user OAM-MIMO-NOMA systems [27], in this paper, we investigate the EE optimization problem for the PSOAM MGs based multi-user MIMO-NOMA system, where the transmit power is optimized through the proposed dual layer resource allocation. To the best of our knowledge, this is the first work that explore the EE performance of the PSOAM MGs based multi-user MIMO-NOMA networks. Our main contributions can be summarized as follows:

- We design a novel resource allocation framework for the downlink PSOAM MGs based MIMO-NOMA system, where a base station (BS) transmits data to multiple users simultaneously through the generated PSOAM beams. The goal of this work is to investigate the EE maximization problem of the proposed framework while satisfying the transmit power and minimum user rate constraints. As a result, two schemes are developed in this paper to cope with this optimization problem.
- Firstly, we present a dual layer resource allocation algorithm to solve this problem. The resultant optimization problem is non-convex and NP-hard involving the optimization of the transmit power, which cannot be solved directly. By applying the generalized fractional programming technique, we transform the original EE maximization into a parametric programming problem, which is still difficult to solve directly. Then, we adopt the first order Taylor approximation to transform the problem into a convex problem, and present a dual layer resource allocation algorithm to solve this problem. The outer layer exploits a bisection method to optimize the EE. In each iteration, the inner layer uses a resource distributed iterative algorithm to optimize the transmit power with a fixed EE.
- The complexity of the proposed dual-layer iterative resource allocation algorithm grows considerably as the number of users increases, thus it is hard to cope with the requirement for low computational complexity. To overcome this challenge, we propose a resource allocation scheme with Deep Belief Networks (DBN) to achieve the goal of EE maximization. The learning process includes training-data generation, model optimization, and model testing.
- Simulation results validate that the proposed dual layer resource allocation algorithm and the proposed resource allocation algorithm with DBN are capable of achieving the optimal EE of our proposed framework. In addition, numerical results also demonstrate that the proposed PSOAM MGs based multi-user MIMO-NOMA system can achieve superior EE performance in comparison with the conventional multi-user MIMO-NOMA systems.

## B. Organization and Notation

The remainder of the paper is organized as follows. The proposed PSOAM MGs based MIMO-NOMA system model and the problem formulation of optimizing EE are described in Section II. In Section III, we propose a dual-layer resource allocation algorithm to solve the non-convex EE optimization problem. In Section IV, another alternative scheme with the help of deep learning is investigated to acquire the solution that maximizes the EE of the proposed system. In Section V, simulation results are presented to validate the effectiveness of the theoretical findings. Finally, we draw the conclusions for this paper in Section VI.

The notations used throughout the full paper are illustrated as follows. Bold case letters represent the vector and non-bold case letters represent the scalar.  $\mathbb{E}[|a|^2]$  denotes the power of  $a$ .

The gradient of  $f(m)$  at the point  $m_0$  is denoted by  $\nabla f(m_0)$ .  $\mathbf{h}^T$  denotes the transpose of the vector  $\mathbf{h}$ , and  $\|\mathbf{a}\|^2$  denotes the  $l_2$ -norm of  $\mathbf{a}$ . We use  $[s]^+$  to represent  $\max(s, 0)$ .

## II. SYSTEM MODEL AND PROBLEM FORMULATION

### A. System model

As shown in Fig. 1, we consider a downlink PSOAM MGs based MIMO-NOMA system, which includes one BS with  $N_t$  antennas. We consider  $K$  ( $K=4$ ) randomly distributed users in a fan-shaped area. The considered system is composed of uniform linear arrays (ULAs) and the antenna spacing is  $\zeta$ . At the transmitting side, each antenna only sends data streams to a corresponding user, in which two superposed PSOAM MGs waves are radiated into the free space using the structure of the circular travelling-wave antenna with a ring horn [28]. The total number of PSOAM modes superposed in one PSOAM mode group is marked as  $G$ , which is sorted in the ascending order as  $\{l^0, l^0 + \Delta l, l^0 + 2\Delta l, \dots, l^0 + (G - 1)\Delta l\}$ , and the mode interval is marked as  $\Delta l$ . As stated in [16], the equivalent PSOAM MGs phase slope can be calculated by the first and the last modes in a MG, and we assume that the signals of the  $n^{\text{th}}$  transmitting antenna is sent to the  $n^{\text{th}}$  user. At the receiver, each user is equipped with  $N_r = 2$  receiving antennas, which are placed within the main lobe of the superposed PSOAM MGs waves. Due to the orthogonality of different PSOAM mode groups in the main lobe zone [19], there is no interference between each different PSOAM mode group. However, since all users share the same bandwidth, the interferences among the users can not be ignored when decoding the signals for the same PSOAM mode group. The total transmit power is restricted to  $P_{max}$  and the signal transmitted to user  $k$  in the  $mg^{\text{th}}$  mode group is written as

$$X_{k,mg} = \sqrt{p_{k,mg}} \cdot x_{k,mg}, \quad (1)$$

where  $p_{k,mg}$  denotes the power allocation of the  $mg^{\text{th}}$  PSOAM mode group at the  $k^{\text{th}}$  user and  $x_{k,mg}$  is the transmit data symbol for the  $mg^{\text{th}}$  mode group of user  $k$ , i.e.,  $\mathbb{E}[|x_{k,mg}|^2] = 1$ .

Let  $\varphi_0^k$  be the initial azimuthal angle of the PSOAM MG pattern of user  $k$  for the  $mg^{\text{th}}$  mode group, the link of channel gain  $h_{k,n_r,n_t,mg}$  between the  $n_t^{\text{th}}$  transmitting antenna and the  $n_r^{\text{th}}$  receiving antenna is formulated as

$$\begin{aligned} h_{k,n_r,n_t,mg} &= \Gamma_1 \frac{1}{\sqrt{G^{mg}}} \sum_{g^{mg}=1}^{G^{mg}} e^{-j(l_0^{mg} + g^{mg} \Delta l) \varphi_{k,n_r,n_t}} \\ &= \Gamma_1 \frac{1}{\sqrt{G^{mg}}} e^{-j l_{eq}^{mg} \varphi_{k,n_r,n_t}}, \end{aligned} \quad (2)$$

where  $\Gamma_1 = \beta_{k,n_t} \frac{\lambda}{4\pi d_{k,n_r,n_t}} e^{-j \frac{2\pi}{\lambda} d_{k,n_r,n_t}}$ ,  $l_0^{mg}$  and  $l_{eq}^{mg}$  represent the first mode among mode group  $mg$  and the equivalent OAM order of mode group  $mg$ , respectively,  $d_{k,n_r,n_t}$  represents the distance between the  $n_t^{\text{th}}$  transmitting antenna and the  $n_r^{\text{th}}$  receiving antenna of user  $k$ ,  $\lambda$  is the wave length, and  $\beta_{k,n_t} = \sqrt{G_t G_r}$  is a constant that related to the antenna gain of the  $n_t^{\text{th}}$  transmitting antenna and the  $n_r^{\text{th}}$  receiving antenna. Specifically, the interference of the  $n_t^{\text{th}}$  transmitting antenna to the other  $K - 1$  non-intended users can be calculated by the

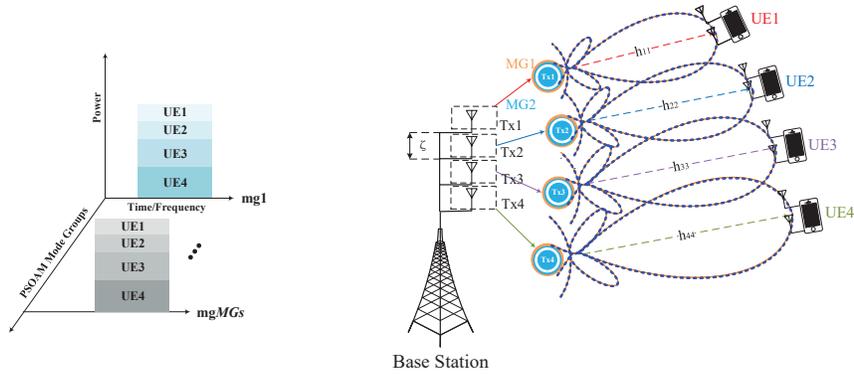


Fig. 1: Illustration of a PSOAM MGs based MIMO-NOMA system model.

antenna gain according to the amplitude and phase distribution of the PSOAM MGs wave [19].  $G_t$  can be determined by the distribution of the users and the interference of the minor lobe. In addition,  $\varphi_{k,n_r,n_t}$  represents the azimuthal angle between the  $n_t^{th}$  transmitting antenna and the  $n_r^{th}$  receiving antenna of user  $k$ , which is described in Fig. 2. To calculate the phase  $\varphi_{k,n_r,n_t}$ , there are two conditions that should be considered. One is when the initial azimuthal angle of the PSOAM MG pattern  $\varphi_0^k > 0$  and the other is when  $\varphi_0^k < 0$ . For  $\varphi_0^k > 0$ , it can be observed in Fig. 2 that it can be further divided into three cases, which are discussed as follows:

1) *Case 1*  $k = n_t$ : We define the vertical distance between the center of the two receiving antennas of user  $k$  and the corresponding transmitting antenna as  $d_k$ . The radius of the receiving aperture is marked as  $R$ . The distance between the  $n_r^{th}$  receiving antenna and the  $n_t^{th}$  transmitting antenna can be calculated by

$$d_{k,n_r,n_t,k=n_t} = \sqrt{d_k^2 + R^2}, \quad (3)$$

where  $R = d_k \tan\left(\frac{\pi}{2|l_{eq}^{mg1} - l_{eq}^{mg2}|}\right)$  and  $d_k$  is the relative distance [16]. The phase  $\varphi_k$  between  $d_k$  and  $d_{k,n_r,n_t,k=n_t}$  is given by

$$\varphi_k = (-1)^{n_r} \cdot \arctan\left(\frac{R}{d_k}\right). \quad (4)$$

Hence, the phase  $\varphi_{k,n_r,n_t,k=n_t}$  is calculated by

$$\varphi_{k,n_r,n_t,k=n_t} = \varphi_k + \varphi_0^k. \quad (5)$$

2) *Case 2*  $k < n_t$ : According to the cosine theorem, we can calculate the distance between the  $n_r^{th}$  receiving antenna of the non-intended user  $k$  and the  $n_t^{th}$  transmitting antenna as follows in (6). In addition, the azimuthal between the  $n_r^{th}$  receiving antenna and the  $n_t^{th}$  transmitting antenna of user  $k$  is defined as

$$\varphi_{k,n_r,n_t} = \frac{\pi}{2} - \omega, \quad (7)$$

$$\omega = \arccos\left(\frac{((n_t - k) \cdot \zeta)^2 + d_{k,n_r,n_t}^2 - d_{k,n_r,n_t,k=n_t}^2}{2 \cdot d_{k,n_r,n_t} \cdot (|n_t - k| \cdot \zeta)}\right). \quad (8)$$

3) *Case 3*  $k > n_t$ : Similarly, the distance between the  $n_r^{th}$  receiving antenna of user  $k$  and the  $n_t^{th}$  transmitting antenna is formulated as follows in (9). The azimuthal angle of the  $k^{th}$  user from  $n_r^{th}$  antenna to the  $n_t^{th}$  transmitting antenna is

$$\varphi_{k,n_r,n_t} = \omega - \frac{\pi}{2}, \quad (10)$$

where  $\omega$  can be calculated by (8).

For the condition  $\varphi_0^k < 0$ , the related angle and distance can be calculated using the same approach as the condition  $\varphi_0^k > 0$ . Due to  $d_{k,n_r,n_t} \gg R$  and  $d_{k,n_r,n_t} \gg \zeta$ , we can make an approximation that  $d_{k,n_r,n_t} \approx d_k$ . Hence, the channel gain is expressed as

$$h_{k,n_r,n_t,mg} = \beta_{k,n_t} \frac{\lambda}{4\pi d_k} e^{-j\frac{2\pi}{\lambda} d_k} \frac{1}{\sqrt{G^{mg}}} e^{-j l_{eq}^{mg} \varphi_{k,n_r,n_t}}. \quad (11)$$

Inspired by the PASR method, we use 1/2 circular arc to receive the PSOAM MGs waves, where  $N_r$  receiving antennas are evenly distributed at the  $\pi$  angular arc. Each antenna only sends data streams to a corresponding user at the transmitting side for the case of PSOAM MGs-based NOMA. In particular, the channel from the  $l^{th}$  transmitting antenna to the  $k^{th}$  user in mode group  $mg$  can be denoted as  $h_{k,l,mg}$ . We assume the channel state information is known and SIC will be adopted to cancel the interferences among multiple users in the NOMA system. Particularly, within the same mode group, user  $k$  can detect the signal of user  $l$  ( $l < k$ ) and remove the signal from its observation, while the signal to user  $l$  ( $l > k$ ) is treated as the interference at user  $k$ . As a result, the total rate of all  $K$  users can be formulated as

$$R_{total} = \sum_{k=1}^{K-1} \sum_{mg=1}^{MGs} B \log_2 \left( 1 + \frac{p_{k,mg} \lambda_{k,k,mg}^2}{\sum_{l=k+1}^K p_{l,mg} \lambda_{k,l,mg}^2 + \sigma^2} \right) + \sum_{mg=1}^{MGs} B \log_2 \left( 1 + \frac{p_{K,mg} \lambda_{K,K,mg}^2}{\sigma^2} \right), \quad (12)$$

where  $\lambda_{k,l,mg}$  is the singular values of the channel matrix that satisfies the condition:  $\lambda_{1,1,mg} \leq \lambda_{2,2,mg} \leq \dots \leq \lambda_{K,K,mg}$ .

### B. Power consumption model

In general, the power consumption of the PSOAM MGs based MIMO-NOMA system consists of transmit power and

$$d_{k,n_r,n_t} = \sqrt{((n_t - k) \cdot \zeta)^2 + d_{k,n_r,n_t,k=n_t}^2 - 2d_{k,n_r,n_t,k=n_t} \cdot \zeta \cdot (n_t - k) \cdot \cos\left(\frac{\pi}{2} + \varphi_{k,n_r,n_t,k=n_t}\right)}. \quad (6)$$

$$d_{k,n_r,n_t} = \sqrt{((n_t - k) \cdot \zeta)^2 + d_{k,n_r,n_t,k=n_t}^2 - 2d_{k,n_r,n_t,k=n_t} \cdot \zeta \cdot (k - n_t) \cdot \cos\left(\frac{\pi}{2} - \varphi_{k,n_r,n_t,k=n_t}\right)}. \quad (9)$$

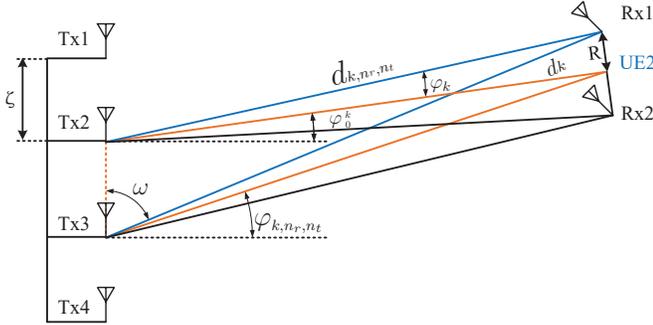


Fig. 2: The structure of the PSOAM MGs based MIMO-NOMA system.

circuit power, which is defined as follows

$$PC_{total} = \alpha \sum_{k=1}^K \sum_{mg=1}^{MGs} p_{k,mg} + P_c, \quad (13)$$

where  $\alpha$  is the power amplifier drain efficiency and  $P_c = \sum_{n_t=1}^{N_t} P_{T_x}^{(n_t)} + \sum_{k=1}^K \sum_{n_r=1}^{N_r} P_{R_x}^{(k,n_r)}$  is the circuit power consumption of system hardware. The circuit power consumed on the RF chain attached to the  $n_t^{th}$  transmit antenna and the RF chain attached to the  $n_r^{th}$  receiving antenna at the  $k^{th}$  user is denoted as  $P_{T_x}^{(n_t)}$  and  $P_{R_x}^{(k,n_r)}$ , which account for mixer, filter, intermediate frequency amplifier and so on [29]–[31].

### C. Problem Formulation

The resultant EE of the PSOAM MGs based MIMO-NOMA system can be given by

$$\gamma_{EE} \triangleq \frac{R_{total}}{PC_{total}} = \frac{\sum_{k=1}^K \sum_{mg=1}^{MGs} B \log_2 \left( 1 + \frac{p_{k,mg} \lambda_{k,k,mg}^2}{\sum_{l=k+1}^K \sum_{mg=1}^{MGs} p_{l,mg} \lambda_{k,l,mg}^2 + \sigma^2} \right)}{\alpha \sum_{k=1}^K \sum_{mg=1}^{MGs} p_{k,mg} + P_c}. \quad (14)$$

The purpose of this work is to maximize the EE of the PSOAM MGs based MIMO-NOMA system with the minimum user rate and the transmit power constraints. Therefore, the EE optimization problem is expressed in (15)–(18), where  $\mathcal{K} = \{1, 2, \dots, K\}$  represents the set of users,  $\mathcal{MG} = \{mg1, mg2, \dots, MGs\}$  denotes the set of all PSOAM mode groups. C1 guarantees the minimum user rate constraint. C2 ensures that the total transmit power is limited to  $P_{max}$ . In C3,  $p_{k,mg}$  should be a positive integer and  $mg \in \mathcal{MG}$  for

any  $k \in \mathcal{K}$  is requested. The considered EE maximization problem is non-convex and thus the solution cannot be obtained directly. In order to solve this problem, we propose an effective dual layer resource allocation algorithm that can optimize  $\gamma_{EE}$  in the outer-layer and then turn to deal with  $\mathbf{P}$  in the inner-layer with a fixed  $\gamma_{EE}$  iteratively.

### III. PROPOSED DUAL-LAYER ITERATIVE RESOURCE ALLOCATION ALGORITHM

In this section, we consider the resource allocation algorithm for the proposed system. The optimization problem in (15)–(18) belongs to a non-convex and non-linear fractional programming problem and thus very difficult to acquire the solution directly. As a result, we develop a dual-layer iterative method where the optimal power allocation  $\mathbf{P}$  is optimized in the inner-layer and the optimal EE is obtained in the outer-layer. To tackle this problem, we transform the fractional objective function into a subtractive form of numerator and denominator according to the following proposition.

**Proposition 1:** For  $R_{k,mg}(\mathbf{P}) > 0$  and  $PC_{total}(\mathbf{P}) > 0$ , the optimal solution  $\mathbf{P}^*$  of (15) is achieved when the following equation is consistent:

$$\max_{\mathbf{P} \in \{C1, C2, C3\}} [R_{total}(\mathbf{P}) - \gamma_{EE}^* PC_{total}(\mathbf{P})] = [R_{total}(\mathbf{P}^*) - \gamma_{EE}^* PC_{total}(\mathbf{P}^*)] = 0, \quad (19)$$

$$\gamma_{EE}^* = \max_{\mathbf{P} \in \{C1, C2, C3\}} \frac{R_{total}(\mathbf{P})}{PC_{total}(\mathbf{P})} = \frac{R_{total}(\mathbf{P}^*)}{PC_{total}(\mathbf{P}^*)}, \quad (20)$$

where  $\gamma_{EE}^*$  can be the optimal EE and  $\mathbf{P}^*$  can be the optimal resource allocation of the optimization problem.

*Proof:* See Appendix A. ■

Proposition 1 gives the condition for developing the optimal resource allocation method from the perspective of necessary and sufficient. We define  $\Upsilon(\gamma_{EE})$  as follows

$$\Upsilon(\gamma_{EE}) = \max_{\mathbf{P}} [R_{total}(\mathbf{P}) - \gamma_{EE} PC_{total}(\mathbf{P})], \quad (21)$$

where  $\Upsilon(\gamma_{EE})$  is a function with independent variable  $\gamma_{EE}$ . To solve the equivalent problem in (19), we first obtain the  $\gamma_{EE}^*$ . Then, we have the following proposition.

**Proposition 2:** Problem (21) is monotonically decreasing with respect to  $\gamma_{EE}$ .

*Proof:* See Appendix B. ■

By applying proposition 2, the bisection method can be used to find  $\gamma_{EE}^*$ . We can adjust the upper bound of  $\gamma_{EE}$  when  $\Upsilon(\gamma_{EE}) > 0$ , whilst the lower bound when  $\Upsilon(\gamma_{EE}) < 0$  until the optimal EE is found.

For a given  $\gamma_{EE}^i$  at the  $i^{th}$  iteration, the optimization problem is now transformed as

$$\max_{P_{k,mg}} \frac{\sum_{k=1}^K \sum_{mg=1}^{MGs} B \log_2 \left( 1 + \frac{p_{k,mg} \cdot \lambda_{k,k,mg}^2}{\sum_{l=k+1}^K p_{l,mg} \lambda_{k,l,mg}^2 + \sigma^2} \right)}{\alpha \sum_{k=1}^K \sum_{mg=1}^{MGs} p_{k,mg} + P_c} \quad (15)$$

$$\text{s.t. } C1 : B \log_2 \left( 1 + \frac{p_{k,mg} \cdot \lambda_{k,k,mg}^2}{\sum_{l=k+1}^K p_{l,mg} \lambda_{k,l,mg}^2 + \sigma^2} \right) \geq R_{req}, \forall k \in \mathcal{K}, \forall mg \in \mathcal{MG}, \quad (16)$$

$$C2 : \sum_{k=1}^K \sum_{mg=1}^{MGs} p_{k,mg} \leq P_{max}, \quad (17)$$

$$C3 : p_{k,mg} \geq 0, \forall k \in \mathcal{K}, \forall mg \in \mathcal{MG}, \quad (18)$$

$$\begin{aligned} \max_{\mathbf{P}} \quad & R_{total}(\mathbf{P}) - \gamma_{EE}^i PC_{total}(\mathbf{P}) \\ \text{s.t.} \quad & C1, C2, C3. \end{aligned} \quad (22)$$

To solve problem (22), we define the function  $\Upsilon(\gamma_{EE}^i)$  as follows

$$\Upsilon(\gamma_{EE}^i) = \max_{\mathbf{P}} R_{total}(\mathbf{P}) - \gamma_{EE}^i PC_{total}(\mathbf{P}). \quad (23)$$

Since different PSOAM MGs are orthogonal to each other,  $\Upsilon(\gamma_{EE}^i)$  can be further simplified. Hence, we can split  $R_{total}$  into  $MGs$  parts according to different PSOAM mode groups, and the total transmit power can be split for the same reason. Therefore, (23) can be transformed by

$$R_{total}(\mathbf{P}) - \gamma_{EE}^i PC_{total}(\mathbf{P}) = F(\mathbf{P}) - H(\mathbf{P}), \quad (24)$$

$$F(\mathbf{P}) = f_{mg1}(\mathbf{P}_{mg1}) + \dots + f_{mgMGs}(\mathbf{P}_{mgMGs}), \quad (25)$$

$$H(\mathbf{P}) = h_{mg1}(\mathbf{P}_{mg1}) + \dots + h_{mgMGs}(\mathbf{P}_{mgMGs}). \quad (26)$$

For each PSOAM MG, by using a logarithmic transformation, we can obtain the expression of the function  $f_{mg}(\mathbf{P}_{mg})$  and  $h_{mg}(\mathbf{P}_{mg})$  as follows

$$\begin{aligned} f_{mg}(\mathbf{P}_{mg}) &= \sum_{k=1}^K B \log_2 \left( \sum_{l=k}^K p_{l,mg} \lambda_{k,l,mg}^2 + \sigma^2 \right) \\ &- \gamma_{EE}^i \left( \alpha \sum_{l=k}^K p_{k,mg} + \frac{P_c}{MGs} \right), \end{aligned} \quad (27)$$

$$h_{mg}(\mathbf{P}_{mg}) = \sum_{k=1}^K B \log_2 \left( \sum_{l=k+1}^K p_{l,mg} \lambda_{k,l,mg}^2 + \sigma^2 \right), \quad (28)$$

where  $p_{l,mg}$  is an element in the vector  $\mathbf{P}$  that is expressed as

$$p_{l,mg} = \mathbf{P}(1, (mg-1) \cdot K + l). \quad (29)$$

Besides, the non-convex constraint  $C1$  in problem (22) is mathematically transformed into an equivalent convex linear form, which is formulated by

$$\begin{aligned} C1' : (1 - 2^{-\frac{R_{req}}{B}}) &\left( \sum_{l=k+1}^K p_{l,mg} \lambda_{k,l,mg}^2 + \sigma^2 \right) \\ &+ p_{k,mg} \lambda_{k,k,mg}^2 \geq 0, \forall k, \forall mg. \end{aligned} \quad (30)$$

Now, the optimization problem (22) is equivalent to

$$\begin{aligned} \max_{\mathbf{P}} \quad & F(\mathbf{P}) - H(\mathbf{P}) \\ \text{s.t.} \quad & C1', C2, C3. \end{aligned} \quad (31)$$

**Proposition 3:** Although the constraints of the optimization problem (31) are convex sets, (15), (22) and (31) remain to be NP-hard problems.

*Proof:* See [32] for the proof of Proposition 3.  $\blacksquare$

From (24), the optimization problem can be regarded as the sum of the expression pairs of PSOAM MGs. We define the expression pairs as  $f$  minus  $h$  and each expression pair is regarded as two concave functions. Therefore, the corresponding resource allocation problem is still non-convex, which is hard to solve directly. Fortunately, we can approximate the function  $h$  to an affine function by applying the first order Taylor approximation to formulate the function. Hence, the corresponding objective function can be written as a concave function minus an affine function, and we can transform the problem into a convex optimization problem.

To solve this issue, we can obtain  $\mathbf{P}_{mg}^q$  through an iterative resource allocation algorithm at the  $q^{th}$  iteration. In particular, the first-order Taylor expansion at  $\mathbf{P}^q$  is formulated as

$$h_{mg}(\mathbf{P}_{mg}^q) + \nabla h_{mg}^T(\mathbf{P}_{mg}^q)(\mathbf{P}_{mg} - \mathbf{P}_{mg}^q)^T, \quad (32)$$

where  $\nabla h_{mg}(\mathbf{P}_{mg})$  represents the gradient of  $h_{mg}(\mathbf{P}_{mg})$ ,  $\mathbf{P}_{mg} = \mathbf{P}((mg-1)K+1, mg \cdot k)$  and the second part in (32) is written as

$$\begin{aligned} &\nabla h_{mg}^T(\mathbf{P}_{mg})(\mathbf{P}_{mg} - \mathbf{P}_{mg}^q)^T \\ &= \sum_{k=1}^K \frac{B}{\sum_{l=k+1}^K p_{l,mg} \lambda_{k,l,mg}^2 + \sigma^2} \mathbf{e}_{mg} \times (\mathbf{P}_{mg} - \mathbf{P}_{mg}^q)^T, \end{aligned} \quad (33)$$

where  $\mathbf{e}_{mg}$  is a  $1 \times K$  matrix. If  $k+1 < j < K$ , we have  $\mathbf{e}_{mg}(j) = \lambda_{k,l,mg}^2 / \ln 2$ . In other cases, we have  $\mathbf{e}_{mg}(j) = 0$ .

TABLE I: PROPOSED DUAL-LAYER ITERATIVE RESOURCE ALLOCATION ALGORITHM.

```

1: Initialization
   Set iteration index  $i = 0$  and stopping criterion  $\varepsilon > 0$ .
   Set  $\gamma_{EE}^{min}$  and  $\gamma_{EE}^{max}$ , let  $\gamma_{EE}^{min} \leq \gamma_{EE}^* \leq \gamma_{EE}^{max}$ .
   Set the maximum iteration number  $L_{max}$ 
2: repeat
3:    $\gamma_{EE}^i = (\gamma_{EE}^{max} + \gamma_{EE}^{min})/2$ .
4:   Solve (22) with a given  $\gamma_{EE}^i$  and get  $P^i$ :
   1): Initialization
      Set the iteration index  $q = 0$ .
      Set the stopping criterion  $\epsilon > 0$ .
      Set the initial value of the transmit power  $P^{(0)}$ .
      Calculate  $I^0 = F(P^{(0)}) - H(P^{(0)})$ .
   2): repeat
   3):   Set the iteration index  $s_1 = 0, \epsilon_1 > 0$ 
        and  $\epsilon_2 > 0$  as the stopping criterions.
        Initialize  $\mu^{(0)} \geq \mathbf{0}, \nu^{(0)} \geq \mathbf{0}$  and  $\psi^{(0)} \geq 0$ .
   4):   Solve (34) to get optimal transmit power  $P^*$ 
   5):   repeat
   6):     Set  $s_2 = 0$  and initialize  $p_{k,mg}^{(0)}$ .
   7):     repeat
   8):       for  $k=1:K$ 
   9):         for  $mg=1:MGs$ 
   10):          Update  $p_{k,mg}^{(s_2)}$  according to (38)-(39).
   11):          end for
   12):        end for
   13):         $s_2 = s_2 + 1$ .
   14):        until  $\|\nabla_{P_{k,mg}} \mathcal{L}\|^2 \leq \epsilon_1$ .
   15):         $s_1 = s_1 + 1$ .
   16):        Update  $\mu^{(s_1)}, \nu^{(s_1)}$  and  $\psi^{(s_1)}$ .
   17):        until  $\mu, \nu$  and  $\psi$  converge,
        i.e.  $\|\mu^{(s_1)} - \mu^{(s_1-1)}\|^2 \leq \epsilon_2$ ,
         $\|\nu^{(s_1)} - \nu^{(s_1-1)}\|^2 \leq \epsilon_2$ , and
         $|\psi^{(s_1)} - \psi^{(s_1-1)}| \leq \epsilon_2$ .
   18):   Set  $q = q + 1$ , and  $P^q = P^*$ .
   19):   Calculate  $I^q = F(P^q) - H(P^q)$ .
   20):   until  $|I^q - I^{q-1}| = |\Xi| \leq \epsilon$ .
5:   if  $|\Upsilon(\gamma_{EE}^i)| = |R_{total}(P^i) - \gamma_{EE}^i PC_{total}(P^i)| \leq \varepsilon$ 
   then  $P^* = P^i$  and  $\gamma_{EE}^* = R_{total}(P^i)/PC_{total}(P^i)$ 
6:   break.
7:   else
8:     if  $\Upsilon(\gamma_{EE}^i) < 0$ , then
9:        $\gamma_{EE}^{max} = \gamma_{EE}^i$ .
10:    else
11:       $\gamma_{EE}^{min} = \gamma_{EE}^i$ .
12:    end if
13:  end if
14:   $i = i + 1$ .
15: until  $i > L_{max}$ .

```

Combining (32) with (33), the optimization problem (31) can be further transformed as follows in (34). Fortunately, (34) is a standard convex optimization problem, which is tackled using the Lagrange duality algorithm. To solve the

problem, we define  $\mu, \nu$  and  $\psi$  as the Lagrangian multipliers of  $C1', C2$  and  $C3$  respectively, and the Lagrangian function of the optimization problem is given as follows in (35). Hence, the dual optimization problem is denoted as

$$g(\mu, \nu, \psi) = \max_{\mathbf{P}} \mathcal{L}(\mathbf{P}, \mu, \nu, \psi), \quad (36)$$

and the dual optimization problem is given by

$$\begin{aligned} \min_{\mu, \nu, \psi} \quad & g(\mu, \nu, \psi) \\ \text{s.t.} \quad & \mu \geq \mathbf{0}, \nu \geq \mathbf{0} \text{ and } \psi \geq 0. \end{aligned} \quad (37)$$

We propose an iterative algorithm based on gradient descent algorithm first to achieve the optimal  $\mathbf{P}$  for the dual optimization problem (37). In particular,  $p_{k,mg}$  is updated successively according to the gradient of the Lagrangian function (35) in terms of  $p_{k,mg}$ , which is expressed as

$$\begin{aligned} \nabla_{p_{k,mg}} \mathcal{L} := & \frac{B}{\ln 2} \sum_{n=1}^k \frac{\lambda_{n,k,mg}^2}{\sum_{l=n}^K p_{l,mg} \lambda_{n,l,mg}^2 + \sigma^2} \\ & - \frac{B}{\ln 2} \sum_{k=1}^K \frac{(e_{mg})_k}{\sum_{l=k+1}^K p_{l,mg}^q \lambda_{(k+1),l,mg}^2 + \sigma^2} \\ & + \sum_{m=1}^{k-1} \mu_{m,mg} \left(1 - 2^{\frac{R_{req}}{B}}\right) (\lambda_{m,k,mg})^2 \\ & + \mu_{k,mg} (\lambda_{k,k,mg})^2 + \nu_{k,mg} - \psi, \\ p_{k,mg}^{(s_2)} = & \left[ p_{k,mg}^{(s_2-1)} + t \nabla_{p_{k,mg}}^{(s_2-1)} \mathcal{L} \right]^+, \end{aligned} \quad (38)$$

where  $(e_{mg})_k$  represents the  $k^{th}$  element of  $e_{mg}$ ,  $t$  represents the iteration step size of  $p_{k,mg}$ , while  $p_{k,mg}^{(s_2)}$  and  $p_{k,mg}^{(s_2-1)}$  indicate the power allocation of the  $s_2^{th}$  and  $(s_2-1)^{th}$  iteration.

We then address the optimization problem (37) to obtain the optimal Lagrangian multipliers. It is worth noting that the objective function and the constraints are linear with reference to Lagrangian multipliers. Hence, the dual problem is convex over the dual variables  $\mu, \nu, \psi$ , which is optimized according to the one dimensional searching algorithm. However, the gradient algorithm is not necessarily feasible since the dual function (36) is not guaranteed to be differentiable. Hence, we update  $\mu, \nu, \psi$  through the sub-gradient algorithm, which is presented in Lemma as follows.

**Lemma:** The sub-gradient of the Lagrange multipliers in  $g(\mu, \nu, \psi)$  is denoted as

$$\nabla_{\mu_{k,mg}} g := \left(1 - 2^{\frac{R_{req}}{B}}\right) \left( \sum_{l=k+1}^K p_{l,mg} \lambda_{k,l,mg}^2 + \sigma^2 \right) \quad (40)$$

$$\begin{aligned} & + p_{k,mg} \lambda_{k,k,mg}^2, \forall k \in \mathcal{K}, \forall mg \in \mathcal{MG}, \\ \nabla_{\nu_{k,mg}} g := & p_{k,mg}, \forall k \in \mathcal{K}, \forall mg \in \mathcal{MG}, \end{aligned} \quad (41)$$

$$\nabla_{\psi} g := P_{max} - \sum_{mg=1}^{MGs} \sum_{k=1}^K p_{k,mg}. \quad (42)$$

*Proof:* See the Lemma 1 in [33] for detailed proof.  $\blacksquare$  Therefore, with the sub-gradient of the Lagrange multipli-

$$\begin{aligned} \max_{\mathbf{P}} \quad & \sum_{mg=1}^{MGs} (f_{mg}(\mathbf{P}_{mg}) - [h_{mg}(\mathbf{P}_{mg}^q) + \nabla h_{mg}^T(\mathbf{P}_{mg}^q)(\mathbf{P}_{mg} - \mathbf{P}_{mg}^q)^T]) \\ \text{s.t.} \quad & C1', C2, C3. \end{aligned} \quad (34)$$

$$\begin{aligned} \mathcal{L}(\mathbf{P}, \boldsymbol{\mu}, \boldsymbol{\nu}, \psi) = & \sum_{mg=1}^{MGs} (f_{mg}(\mathbf{P}_{mg}) - [h_{mg}(\mathbf{P}_{mg}^q) + \nabla h_{mg}^T(\mathbf{P}_{mg}^q)(\mathbf{P}_{mg} - \mathbf{P}_{mg}^q)^T]) \\ & + \sum_{mg=1}^{MGs} \sum_{k=1}^K \mu_{k,mg} \left( (1 - 2^{-\frac{R_{req}}{B}}) \left( \sum_{l=k+1}^K p_{l,mg} \lambda_{k,l,mg}^2 + \sigma^2 \right) + p_{k,mg} \lambda_{k,k,mg}^2 \right) \\ & + \sum_{mg=1}^{MGs} \sum_{k=1}^K \nu_{k,mg} \cdot p_{k,mg} + \psi \left( P_{max} - \sum_{mg=1}^{MGs} \sum_{k=1}^K p_{k,mg} \right). \end{aligned} \quad (35)$$

ers, we can update  $\boldsymbol{\mu}$ ,  $\boldsymbol{\nu}$  and  $\psi$  as follows

$$\begin{aligned} \mu_{k,mg}^{(s_1)} = & \left[ \mu_{k,mg}^{(s_1-1)} + \varpi^{(s_1-1)} \nabla_{\mu_{k,mg}} g \right]^+, \\ \forall k \in \mathcal{K}, \forall mg \in \mathcal{MG}, \end{aligned} \quad (43)$$

$$\begin{aligned} \nu_{k,mg}^{(s_1)} = & \left[ \nu_{k,mg}^{(s_1-1)} + \varpi^{(s_1-1)} \nabla_{\nu_{k,mg}} g \right]^+, \\ \forall k \in \mathcal{K}, \forall mg \in \mathcal{MG}, \end{aligned} \quad (44)$$

$$\psi_{k,mg}^{(s_1)} = \left[ \psi_{k,mg}^{(s_1-1)} + \varpi^{(s_1-1)} \nabla_{\psi_{k,mg}} g \right]^+, \quad (45)$$

where  $\varpi$  represents the iteration step size of Lagrange multipliers.

Now we can conclude the approach to deal with the optimization problem in (15)-(18) under a given  $\gamma_{EE}$ . The detail information of the proposed dual-layer resource allocation algorithm is presented in TABLE I.

We denote the computational complexity of the proposed algorithm by  $O(N_1 N_2 N_3)$ , which consists of three parts: the outer layer part  $N_1$ , the inner layer part  $N_2$  and the SIC part  $N_3$  [24]. The number of iterations of the outer layer is bounded by  $N_1 = L_{max}$ , where  $L_{max}$  is the maximum iteration number [32]. The complexity of the inner layer is approximated as  $N_2 = O(\lceil I^{max} - I^0 - \frac{\epsilon}{\Delta I} \rceil \cdot K \cdot MGs \cdot \log \frac{1}{\epsilon})$ , where  $I^{max} = \max_{\mathbf{P}} [F(\mathbf{P}) - H(\mathbf{P})]$ ,  $\Delta I = \min_q [I^q - I^{q-1}]$  [34]. The interference of each terminal is successfully cancelled in SIC process and thus the computational complexity depends on the number of users is approximately of order  $O(K^{2.376})$  [35]. In general, the total complexity of the proposed resource allocation iterative algorithm is approximately  $O(L_{max} \cdot \lceil I^{max} - I^0 - \frac{\epsilon}{\Delta I} \rceil \cdot K \cdot MGs \cdot \log \frac{1}{\epsilon} \cdot K^{2.376})$ .

#### IV. PROPOSED RESOURCE ALLOCATION ALGORITHM WITH DBN

As the number of users increases, the complexity of the proposed dual-layer iterative resource allocation algorithm grows considerably, which is hard to cope with the requirement for low latency applications [36]. Deep learning is capable of learning complex algorithms from only observed data and has

achieved remarkable success in communication systems. To overcome this challenge, we investigate a resource allocation algorithm with deep learning framework to acquire the solution that maximizes the EE of the proposed system. Specifically, we adopt a deep belief network (DBN) as the backbone network since its robustness and flexibility in solving the resource allocation optimization problem. Before applying the DBN to the optimization problem of this paper, we briefly review its advantages and the specific implementation details. Then, the involved training data generation, model optimization, and model prediction are presented as follows.

##### A. Preliminaries

A DBN consists of multiple Restricted Boltzmann Machines (RBMs), which include an input layer and a hidden layer. Input data is added to the input layer of the DBN and then passed to the hidden layer, where this information is processed by multiple neurons. The optimization process of the DBN includes pre-training and fine-tuning. Firstly, the DBN is pre-trained layer by layer using the energy function to obtain the initial parameters. Secondly, based on the back-propagation algorithm, the DBN can be fine-tuned by supervision of training data.

##### B. Training Data Generation

The success of deep neural networks is inseparable from the support of a large amount of training data. Hence, the generation of training data is necessary. The vector of the randomly generated channel gains serve as the input  $\mathbf{X}$  of the DBN. Besides, the output vector  $\mathbf{Y}$  of the DBN is provided by the corresponding optimal solution  $\{\mathbf{p}^*\}$ , which is obtained through the complete algorithm developed in Section III. The BS with  $N_t$  antennas radiating  $MGs$  mode groups communicate with  $K$  users, the set of optimal solution  $\{\mathbf{p}^*\}$  is composed of  $MGs \cdot K$  components. Hence, the number of DBN is set to  $MGs \cdot K$  in this paper. In particular,  $\hat{y}_l$ , the resource allocated to user  $k$  under the  $mg^{th}$  mode group, is predicted by  $DBN_l$ .

### C. Model Optimization

**Pre-training:** Following the practice of deep learning, network optimization can benefit from better initialization in pre-training. Let  $\mathbf{v}$  and  $\mathbf{h}$  represent the input and hidden layer of the RBM in  $\text{DBN}_l$ , respectively. The biases of  $\mathbf{v}$  and  $\mathbf{h}$  is accordingly marked as  $\mathbf{b}_v$  and  $\mathbf{b}_h$ . The weight between  $\mathbf{v}$  and  $\mathbf{h}$  is defined as  $\mathbf{w}$ . Let  $\Lambda = \{\mathbf{w}, \mathbf{b}_v, \mathbf{b}_h\}$ . The loss of pre-training  $\mathcal{L}_{pre}(\mathbf{v}(t))$  for set  $\Lambda$  can be formulated as follows:

$$\Lambda_{(t+1)} = \Lambda_{(t)} - \chi \frac{-\partial \log \mathcal{L}_{pre}(\mathbf{v}(t))}{\partial \Lambda_{(t)}}, \quad (46)$$

where  $t$  and  $\chi$  are the number of iterations and the learning rate, respectively.  $\mathcal{L}_{pre}(\mathbf{v}(t))$  indicates the probability distribution of  $\mathbf{v}(t)$  and can be normalized as the following function

$$\begin{aligned} \mathcal{L}_{pre}(\mathbf{v}(t)) &= \sum_{\mathbf{h}(t)} \mathcal{L}_{pre}(\mathbf{v}(t), \mathbf{h}(t)) \\ &= \sum_{\mathbf{h}(t)} \frac{\exp(-E(\mathbf{v}(t), \mathbf{h}(t)))}{\sum_{\mathbf{v}(t)} \sum_{\mathbf{h}(t)} \exp(-E(\mathbf{v}(t), \mathbf{h}(t)))}, \end{aligned} \quad (47)$$

where  $E(\mathbf{v}(t), \mathbf{h}(t))$  is the energy function and is defined as

$$E(\mathbf{v}(t), \mathbf{h}(t)) = -\mathbf{v}'_{(t)} \mathbf{w}(t) \mathbf{h}(t) - \mathbf{b}'_{v(t)} \mathbf{v}(t) - \mathbf{b}'_{h(t)} \mathbf{h}(t). \quad (48)$$

**Fine-tuning:** With the supervision of training data, the parameter set  $\Lambda$  can be fine-tuned iteratively by applying the back-propagation algorithm, which can be formulated by

$$\Lambda_{(t'+1)} = \Lambda_{(t')} - \tilde{\chi} \frac{\partial \mathcal{L}_f}{\partial \Lambda_{(t')}}, \quad (49)$$

where  $\tilde{\chi}$  represents the learning rate and  $t'$  is the number of iteration in the fine-tuning stage. Particularly, the loss function of fine-tuning stage  $\mathcal{L}_f$  is adopted as the cross entropy function to reduce the estimation error of  $\text{DBN}_l$ . The fine-tuning procedure can be calculated as

$$\mathcal{L}_f = -\frac{1}{N} \sum_{i=1}^N \left( y_l^{(i)} \log(\hat{y}_l^{(i)}) + (1 - y_l^{(i)}) \log(1 - \hat{y}_l^{(i)}) \right), \quad (50)$$

where  $N$  is the number of training instances,  $y_l^{(i)}$  is the output of corresponding instance in the  $i^{\text{th}}$  iteration and  $\hat{y}_l^{(i)}$  represent the corresponding prediction from  $\text{DBN}_l$ .

### D. Model Prediction

The channel coefficients  $h_{k,l,mg}$  are randomly generated and the input layers of all DBNs are given by  $\mathbf{X} = [h_{1,1,mg1}, \dots, h_{1,N_t,mg1}, \dots, h_{K,1,mg1}, \dots, h_{K,N_t,mg1}, \dots, h_{K,N_t,mg_{MGs}}]^T$ . With the well-trained model and the input  $\mathbf{X}$ , the estimation of the optimal output is calculated directly from model's prediction. Hence, the approximated solution  $\mathbf{p}^*$  of the proposed framework can be obtained as  $\mathbf{p}^* = [\hat{y}_{1,1}, \dots, \hat{y}_{K,1}, \dots, \hat{y}_{K,MGs}]^T$ . Details of the proposed DBN-based resource allocation algorithm is summarized in TABLE II.

## V. SIMULATION RESULTS

In this section, simulation results are presented to evaluate the performance of the bisection-based resource allocation algorithm. To evaluate the EE performance, we employ a channel with the carrier frequency operating at 10 GHz [37].

TABLE II: RESOURCE ALLOCATION ALGORITHM BASED ON DBN.

- 1): Set the stopping criterion  $\epsilon_3$  and  $\epsilon_4$ .
- 2): **Training Data Generation Part:**
- 3): generate plenty of training data samples  $\{\mathbf{X}, \mathbf{Y}\}$
- 4): **Model Optimization Part:**
- 5): **for**  $l = 1 : K \cdot MGs$
- 6): **for**  $m = 1 : M$
- 7): initialize  $\Lambda$  for the  $m^{\text{th}}$  RBM
- 8): **Pre-training phase**
- 9): **repeat**
- 10): calculate  $\Lambda$  according to equation (46)
- 11): **until**  $\|\Lambda_{(t+1)} - \Lambda_{(t-1)}\|^2 \leq \epsilon_3$
- 12): **Fine-tuning phase**
- 13): **repeat**
- 14):  $\Lambda$  is fine-tune on account of equation (49)
- 15): **until**  $\|\Lambda_{(t'+1)} - \Lambda_{(t'-1)}\|^2 \leq \epsilon_4$
- 16): **end**
- 17): **end**
- 18): **Model Prediction Part:**
- 19): randomly generate channel gains and define  $\mathbf{X}$  as  $[h_{1,1,mg1}, \dots, h_{1,N_t,mg1}, \dots, h_{K,N_t,mg_{MGs}}]^T$
- 20): load the well-trained model
- 21): predict  $\hat{y}_l (1 \leq l \leq MGs \cdot K)$  with  $\mathbf{X}$  and obtain the approximation of solution

The power amplifier drain efficiency is set to  $\alpha = 2$ , the circuit power  $P_{Tx}^{(n_i)}$  and  $P_{Rx}^{(k,n_i)}$  is set to 100 mW. In particular, we consider  $K = 4$  users randomly distributed in a fan-shaped area, which is 2 km away from the BS and all the results are averaged over various random locations of terminals. The proposed system is composed of uniform linear arrays with 4 antennas at the transmitter and the element spacing is  $\zeta = 7\lambda$ . The selected PSOAM MGs of each transmitting antenna are  $mg^1 = \{-55, -54, \dots, -46, -45\}$  and  $mg^2 = \{45, 46, \dots, 54, 55\}$ . The radius of the receiving aperture R is approximately 1.57 m when the relative distance is 100 m, and hence the EE is calculated at a R of 1.57 m [16]. The bandwidth of the system is normalized to 1 Hz. According to the algorithm we proposed, the stopping criteria are set to  $\epsilon = \epsilon_3 = 10^{-3}$ . It is worth noting that the parameters in this system are selected to prove the performance of EE as an example and can be replaced by other reasonable parameters according to the specific scenarios.

Our DBNs are implemented based on *Pytorch*. We prepare 1000 data samples  $\{\mathbf{X}, \mathbf{Y}\}$  to train the DBNs in each case of various system parameters setting. The number of neurons of each hidden layer in the DBN model is taken as 32, 64 and 32, respectively. The learning rate  $\chi$  and  $\tilde{\chi}$  are taken as  $1e^{-4}$ . Besides, the training epoches is 2000. The terminating threshold  $\epsilon_3 = 1e^{-3}$  and  $\epsilon_4 = 1e^{-3}$ .

In the first simulation, the convergence behaviour of the proposed bisection-based resource allocation algorithm is evaluated by demonstrating how the  $\Xi$  and  $\Upsilon$  behave with the number of iterations. We set  $P_{max} = 2W$ ,  $R_{req} = 1\text{bit/s/Hz}$ .

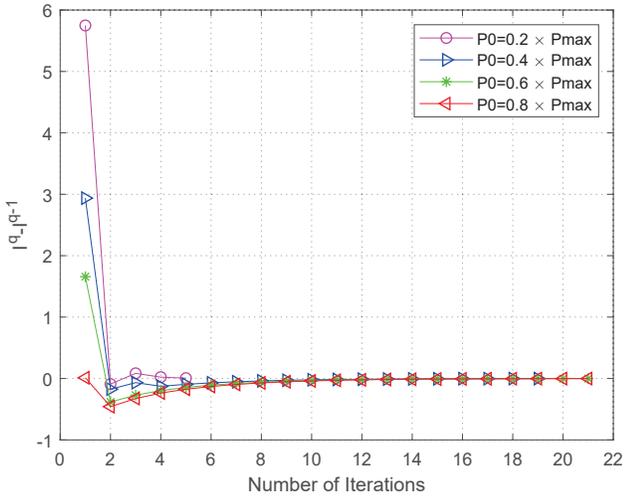


Fig. 3: The convergence evolution of the proposed resource allocation algorithm.

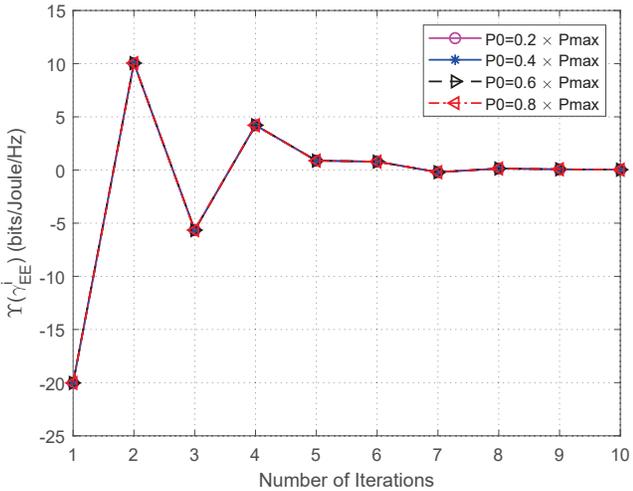


Fig. 4: The convergence evolution of the proposed bisection-based EE optimization algorithm.

As it can be seen in Fig. 3, the inner layer of the proposed algorithm converges to zero, and the initiation point  $P_0$  affects the convergence speed of the proposed algorithm. Specifically,  $\Xi$  converges to zero after eight iteration when  $P_0 = 0.2P_{max}$ ,  $0.4P_{max}$ ,  $0.6P_{max}$ ,  $0.8P_{max}$ . Moreover, as it can be seen in Fig. 4, the outer layer of the proposed algorithm can also converge to zero at approximately eight iterations, which demonstrates that the convergence of our proposed two-layer algorithm can converge to a stable value. We also investigate the  $\gamma_{EE}^i$  of the proposed EE algorithm with different constraints of  $\gamma_{EE}^{max}$ . We set  $P_0 = 0.2P_{max}$ ,  $P_{max} = 2W$ ,  $R_{req} = 1\text{bit/s/Hz}$ . As shown in Fig. 5,  $\gamma_{EE}^i$  converges to a stable value after eight iteration with  $\gamma_{EE}^{max} = 3, 5, 10$  bits/Joule/Hz, and the  $\gamma_{EE}^{max}$  affects the convergence speed of the proposed bisection method. These results demonstrate the stability and validity of the proposed algorithm.

In the next simulation, we study the EE performance of the presented algorithm under different number of users with

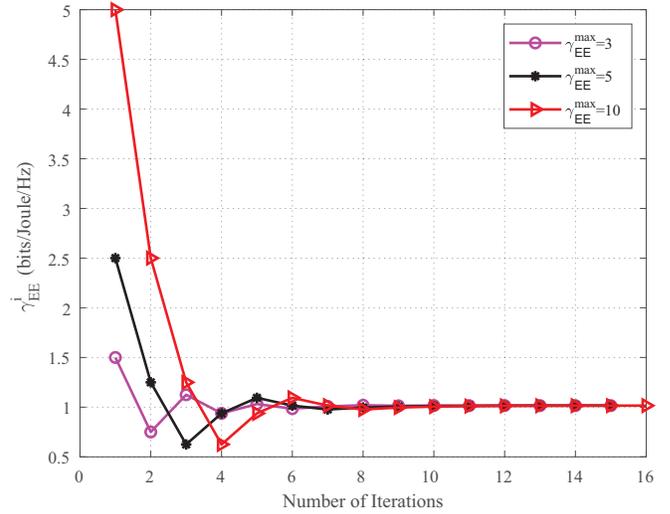


Fig. 5: An example of convergence evolution of the proposed bisection-based EE optimization algorithm in terms of  $\gamma_{EE}^{max}$ .

various circuit power  $P_c$ . The number of users is ranged from one to seven, and circuit power  $P_{T_x}^{(n_t)}$  and  $P_{R_x}^{(k,n_r)}$  are set to 100mW, 120mW and 140mW. As shown in Fig. 6, with the increase of  $P_c$ ,  $\gamma_{EE}^*$  decreases accordingly as anticipated. The reason is that  $\gamma_{EE}$  is inversely proportional to  $P_{total}$ , and hence increasing circuit power leads to increasing the total power consumption of the system, thus leads to a poor  $\gamma_{EE}^*$ . In addition, if  $P_c$  remains constant,  $\gamma_{EE}^*$  decreases as the number of users increases. This is because the interference within the same PSOAM MGs will be enhanced with the increasing users number. Therefore, higher transmit power is required to meet the minimum user rate and the hardware circuit consumption when the number of users is large. Hence, this results in a poor  $\gamma_{EE}^*$  when the users number raises. Furthermore, to show the EE gain achieved by the proposed PSOAM MGs based MIMO-NOMA system, we compare our proposed algorithm with the SE maximization for PSOAM MGs based MIMO-NOMA and the EE maximization for PSOAM MGs based OMA. As shown in Fig. 7, EE achieved by the proposed PSOAM MGs based MIMO-NOMA system outperform the PSOAM MGs based OMA scheme regardless of the number of users. This is because NOMA can enhance the EE by allowing simultaneously serving multiple users with the same physical resource. The results also show that NOMA combined with PSOAM can effectively improve the EE performance. Besides, the PSOAM MGs based MIMO-NOMA system achieves higher EE compared with the SE maximization scheme for PSOAM MGs based MIMO-NOMA, which demonstrates the effectiveness of our proposed resource allocation algorithm. Particularly, the superiority is more obvious as the number of users increase. This is because the interference increases with the growth of the number of users and the increasing sum rate of the system cannot offset the consumption of the transmit power, leading to a decrease in EE. In addition, the simulation results for the proposed resource allocation scheme with DBN are in accordance with the results obtained by the dual layer resource allocation algorithm, which prove the validity of the

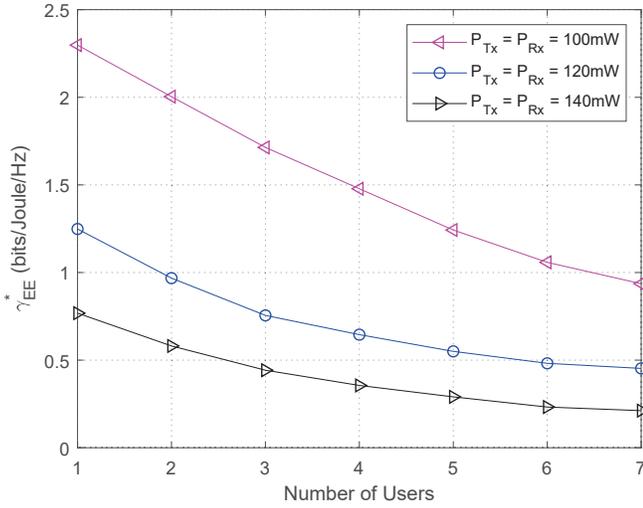


Fig. 6: The performance of the proposed algorithm with different number of users and circuit power.

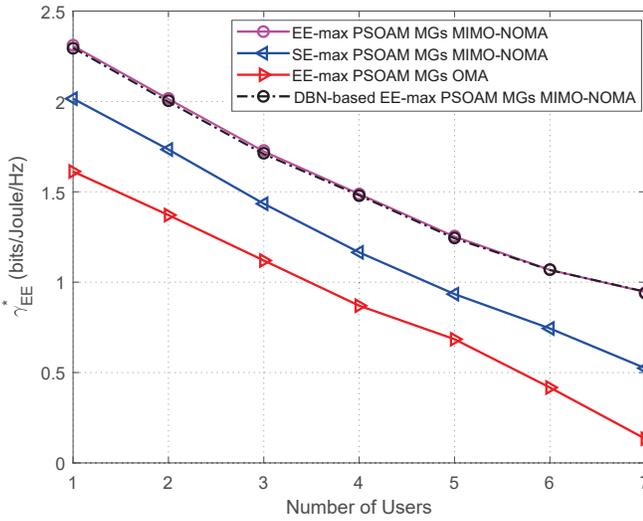


Fig. 7: Comparison of the optimal EE for different number of users in the network.

proposed DBN-based framework.

Next, we aim to study the EE performance under the constraints of transmit power and minimum required data rate. We first show the optimal EE  $\gamma_{EE}^*$  of the proposed bisection-based resource allocation algorithm with different transmit power  $P_{max}$  and minimum required data rate  $R_{req}$ . To prove the effectiveness of our proposed method, we apply the algorithms in the PSOAM MGs based MIMO-NOMA system and the traditional multi-user MIMO-NOMA system [38] for comparison. As can be seen in Fig. 8, the  $\gamma_{EE}^*$  obtained by all the algorithms are monotonically decreasing with the increase of  $R_{req}$ . For the proposed dual layer resource allocation algorithm, a significant drop occurs in our proposed system when the minimum data rate of users is larger than 4 bit/s/Hz. This is owing to the fact that the limitation of transmit power cannot satisfy the QoS requirement of each user. Similarly, the EE of the conventional MIMO system decreases slowly in the lower rate region, i.e.,  $0.5 \leq R_{req} < 2.5$  bit/s/Hz, and

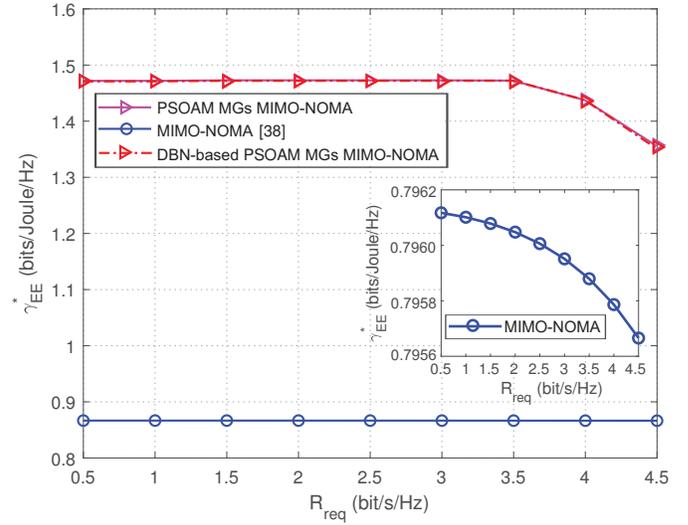


Fig. 8: The performance of the proposed algorithms with different minimum rate requirements.

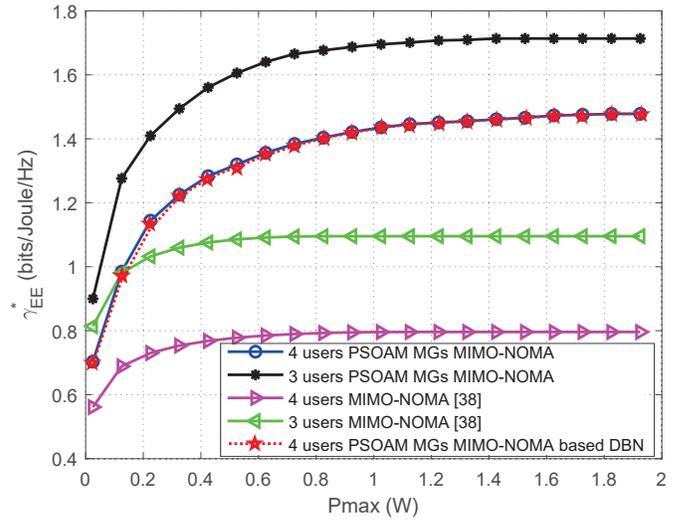


Fig. 9: The performance of the proposed algorithms with different entire transmit power constraints.

then decreases sharply when  $R_{req} > 2.5$  bit/s/Hz. This is because a lower rate constraint needs smaller transmit power and thus the transmit power is allocated to achieve the  $\gamma_{EE}^*$ . Inversely, more transmit power ought to be allocated to satisfy the minimum user rate when  $R_{req}$  is high, which leads to the rapidly decreasing curve. Compared with the conventional multi-user MIMO-NOMA system, our proposed solution is capable of achieving a significant performance gain in terms of EE owing to the degree of freedom provided by PSOAM MGs and MIMO-NOMA techniques, which is a splendid solution for the improvement of EE.

Finally, we investigate the  $\gamma_{EE}^*$  of the proposed solution with various transmit power  $P_{max}$  as well as different users number. As can be seen in Fig. 9, the  $\gamma_{EE}^*$  obtained by the two approaches first increase and then become constant at high transmit power region since the transmitter clips the transmit power once the maximum EE is achieved. Specifically, the

$\gamma_{EE}^*$  increases dramatically with a lower  $P_{max}$ , and then achieves an asymptotic value when the balance between the available rates and the energy consumption is obtained. This is because in the high  $P_{max}$  region, a portion of the maximum transmit power is used to remain the optimal EE instead of exploiting the full  $P_{max}$  region to maximize the total rate, hence the total rates will not further increase with  $P_{max}$ . Additionally, higher  $P_{max}$  is required to achieve the stable  $\gamma_{EE}^*$  when the users number increases in the system. This is because if the users number becomes larger, the total harvested power will increase and hence lead to a poor EE performance. In addition, it can be seen that our presented PSOAM MGs based multi-user MIMO-NOMA system achieves much higher EE than the conventional multi-user MIMO-NOMA system for all feasible  $P_{max}$  values. This is because the number of sub-channels in our proposed system is larger than that of the conventional MIMO-NOMA system due to mode division multiplexing using PSOAM MGs. Thus, it is more efficient to use the structure PSOAM MGs beams to achieve higher data rate or EE compared with the conventional multi-user MIMO-NOMA system.

## VI. CONCLUSION

In this paper, we explore the EE optimization problem for a PSOAM MGs based multi-user MIMO-NOMA system. We aim to maximize the EE while meeting the constraints of total transmit power and minimum user rate. The corresponding problem of maximizing EE is NP-hard and cannot be tackled directly. To solve this problem, a dual layer resource allocation algorithm is developed. Particularly, we achieve the optimal EE in the outer layer via the bisection-based solution and achieve the optimal resource allocation in the inner layer through the resource distribution iterative algorithm. We also investigate another alternative algorithm to achieve the goal of EE maximization with the help of deep learning. Numerical results validate the superiority of the proposed PSOAM MGs based MIMO-NOMA system in EE compared with the conventional multi-user MIMO-NOMA system. Future works will investigate the performance of PSOAM MGs for the channel with partial channel state information (CSI). Besides, the EE and SE trade-off problem for the PSOAM MGs based multi-user MIMO-NOMA system is also a significant task in the future.

### APPENDIX A

#### PROOF OF PROPOSITION 1

We can make an assumption that the optimal transmit resource allocation for (19) is  $\mathbf{P}^*$ . For all feasible transmit resource allocation  $\mathbf{P} \in \{C1, C2, C3\}$ , we can obtain the formulas

$$R_{total}(\mathbf{P}) - \gamma_{EE}^* PC_{total}(\mathbf{P}) \leq 0, \quad (51)$$

$$R_{total}(\mathbf{P}^*) - \gamma_{EE}^* PC_{total}(\mathbf{P}^*) = 0. \quad (52)$$

Equation (51) and (52) can be transformed as  $\frac{R_{total}(\mathbf{P})}{PC_{total}(\mathbf{P})} \leq \gamma_{EE}^*$  and  $\frac{R_{total}(\mathbf{P}^*)}{PC_{total}(\mathbf{P}^*)} = \gamma_{EE}^*$ , respectively. Therefore,  $\mathbf{P}^*$  is the optimal solution for (15) as well as (19).

Now, we start to prove the necessity of *Proposition*

I. According to (15), we can obtain  $\frac{R_{total}(\mathbf{P})}{PC_{total}(\mathbf{P})} \leq \gamma_{EE}^*$  and  $\frac{R_{total}(\mathbf{P}^*)}{PC_{total}(\mathbf{P}^*)} = \gamma_{EE}^*$  for all feasible solution  $\mathbf{P} \in \{C1, C2, C3\}$  which can be transformed into

$$R_{total}(\mathbf{P}) - \gamma_{EE}^* PC_{total}(\mathbf{P}) \leq 0, \quad (53)$$

$$R_{total}(\mathbf{P}^*) - \gamma_{EE}^* PC_{total}(\mathbf{P}^*) = 0. \quad (54)$$

$\mathbf{P}^*$  is the optimal solution for (19) as well.

### APPENDIX B

#### PROOF OF PROPOSITION 2

We assume that there are two values of EE denoted as  $\gamma_{EE}^a$  and  $\gamma_{EE}^b$ . Let  $\gamma_{EE}^a > \gamma_{EE}^b$ . The optimal transmit resource allocation corresponding to  $\gamma_{EE}^a$  and  $\gamma_{EE}^b$  is  $\mathbf{P}^a$  and  $\mathbf{P}^b$ . We have

$$\begin{aligned} \Upsilon(\gamma_{EE}^a) &= \max_{\mathbf{P}} R_{total}(\mathbf{P}) - \gamma_{EE}^a PC_{total}(\mathbf{P}) \\ &= R_{total}(\mathbf{P}^a) - \gamma_{EE}^a PC_{total}(\mathbf{P}^a) \\ &< R_{total}(\mathbf{P}^a) - \gamma_{EE}^b PC_{total}(\mathbf{P}^a) \\ &\leq R_{total}(\mathbf{P}^b) - \gamma_{EE}^b PC_{total}(\mathbf{P}^b) \\ &= \Upsilon(\gamma_{EE}^b). \end{aligned} \quad (55)$$

Hence, we have proved the strictly monotonically decreasing character of (21) in  $\gamma_{EE}$ .

### REFERENCES

- [1] F. Guo, F. R. Yu, H. Zhang, X. Li, H. Ji, and V. C. M. Leung, "Enabling massive IoT toward 6G: A comprehensive survey," *IEEE Internet Things J.*, pp. 1–1, 2021.
- [2] O. Edfors and A. J. Johansson, "Is orbital angular momentum (OAM) based radio communication an unexploited area?" *IEEE Trans. Antennas Propag.*, vol. 60, no. 2, pp. 1126–1131, 2012.
- [3] W. Zhang, S. Zheng, X. Hui, R. Dong, X. Jin, H. Chi, and X. Zhang, "Mode division multiplexing communication using microwave orbital angular momentum: An experimental study," *IEEE Trans. Wireless Commun.*, vol. 16, no. 2, pp. 1308–1318, 2017.
- [4] B. Thidé, H. Then, J. Sjöholm, K. Palmer, J. Bergman, T. Carozzi, Y. N. Istomin, N. Ibragimov, and R. Khamitova, "Utilization of photon orbital angular momentum in the low-frequency radio domain," *Phys. Rev. Lett.*, vol. 99, no. 8, p. 087701, 2007.
- [5] L. Liang, W. Cheng, W. Zhang, and H. Zhang, "Joint OAM multiplexing and OFDM in sparse multipath environments," *IEEE Trans. Veh. Technol.*, vol. 69, no. 4, pp. 3864–3878, 2020.
- [6] Y. Yagi, H. Sasaki, T. Yamada, and D. Lee, "200 Gbit/s wireless transmission using dual-polarized OAM-MIMO multiplexing with uniform circular array on 28 GHz band," *IEEE Antennas Wireless Propag. Lett.*, pp. 1–1, 2021.
- [7] H. Jing, W. Cheng, Z. Li, and H. Zhang, "Concentric UCAs based low-order OAM for high capacity in radio vortex wireless communications," *J. of Commun. Inf. Netw.*, vol. 3, no. 4, pp. 85–100, 2018.
- [8] C. Zhang and Y. Zhao, "Orbital angular momentum nondegenerate index mapping for long distance transmission," *IEEE Trans. Wireless Commun.*, vol. 18, no. 11, pp. 5027–5036, 2019.
- [9] X. Ge, R. Zi, X. Xiong, Q. Li, and L. Wang, "Millimeter wave communications with OAM-SM scheme for future mobile networks," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 9, pp. 2163–2177, 2017.
- [10] W. Zhang, S. Zheng, Y. Chen, X. Jin, H. Chi, and X. Zhang, "Orbital angular momentum-based communications with partial arc sampling receiving," *IEEE Commun. Lett.*, vol. 20, no. 7, pp. 1381–1384, 2016.
- [11] S. Zheng, X. Hui, X. Jin, H. Chi, and X. Zhang, "Transmission characteristics of a twisted radio wave based on circular traveling-wave antenna," *IEEE Trans. Antennas Propag.*, vol. 63, no. 4, pp. 1530–1536, 2015.
- [12] Z. Zhang, S. Zheng, W. Zhang, X. Jin, H. Chi, and X. Zhang, "Experimental demonstration of the capacity gain of plane spiral OAM-based MIMO system," *IEEE Microw. Wireless Compon. Lett.*, vol. 27, no. 8, pp. 757–759, 2017.

- [13] S. Zheng, R. Dong, Z. Zhang, X. Yu, X. Jin, H. Chi, Z. N. Chen, and X. Zhang, "Non-line-of-sight channel performance of plane spiral orbital angular momentum MIMO systems," *IEEE Access*, vol. 5, pp. 25 377–25 384, 2017.
- [14] Z. Zhang, S. Zheng, Jiayu Zheng, X. Jin, H. Chi, and X. Zhang, "Plane spiral orbital angular momentum wave and its applications," in *IEEE MTT-S Int. Microw. Symp. Dig.*, 2016, pp. 1–4.
- [15] J. Zhou, S. Zheng, X. Yu, X. Jin, and X. Zhang, "Low probability of intercept communication based on structured radio beams using machine learning," *IEEE Access*, vol. 7, pp. 169 946–169 952, 2019.
- [16] X. Xiong, S. Zheng, Z. Zhu, X. Yu, X. Jin, and X. Zhang, "Performance analysis of plane spiral OAM mode-group based MIMO system," *IEEE Commun. Lett.*, vol. 24, no. 7, pp. 1414–1418, 2020.
- [17] X. Xiong, S. Zheng, Z. Zhu, Z. Wang, Y. Chen, X. Yu, and X. Zhang, "Direct generation of OAM mode-group and its application in LoS-MIMO system," *IEEE Commun. Lett.*, vol. 24, no. 11, pp. 2628–2631, 2020.
- [18] Z. Zhu, S. Zheng, X. Xiong, Y. Chen, X. Jin, X. Yu, and X. Zhang, "A compact pattern reconfiguration antenna based on multimode plane spiral OAM," *IEEE Trans. Antennas Propag.*, vol. 69, no. 2, pp. 1168–1172, 2021.
- [19] Z. Wang, S. Zheng, X. Xiong, Z. Zhu, Y. Chen, X. Yu, X. Jin, and X. Zhang, "Structure radio beam construction in azimuthal domain," *IEEE Access*, vol. 8, pp. 9395–9402, 2020.
- [20] L. Song, Y. Li, Z. Ding, and H. V. Poor, "Resource management in non-orthogonal multiple access networks for 5g and beyond," *IEEE Network*, vol. 31, no. 4, pp. 8–14, 2017.
- [21] Y. Yu, H. Chen, Y. Li, Z. Ding, and B. Vucetic, "On the performance of non-orthogonal multiple access in short-packet communications," *IEEE Commun. Lett.*, vol. 22, no. 3, pp. 590–593, 2018.
- [22] Y. Zhang, H. Wang, T. Zheng, Q. Yang, "Energy-efficient transmission design in non-orthogonal multiple access," *IEEE Trans. Veh. Technol.*, vol. 66, no. 3, pp. 2852–2857, 2017.
- [23] 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 Process. Lett.*, vol. 21, no. 12, pp. 1501–1505, 2014.
- [24] J. Tang, J. Luo, M. Liu, D. K. C. So, E. Alsusa, G. Chen, K. Wong, and J. A. Chambers, "Energy efficiency optimization for NOMA with SWIPT," *IEEE J. Sel. Topics in Signal Process.*, vol. 13, no. 3, pp. 452–466, 2019.
- [25] J. Ding, J. Cai, and C. Yi, "An improved coalition game approach for MIMO-NOMA clustering integrating beamforming and power allocation," *IEEE Trans. Veh. Technol.*, vol. 68, no. 2, pp. 1672–1687, 2019.
- [26] B. Wang, L. Dai, Z. Wang, N. Ge, and S. Zhou, "Spectrum and energy-efficient beamspace MIMO-NOMA for millimeter-wave communications using lens antenna array," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 10, pp. 2370–2382, 2017.
- [27] A. A. Amin and S. Y. Shin, "Channel capacity analysis of non-orthogonal multiple access with OAM-MIMO system," *IEEE Wireless Commun. Lett.*, vol. 9, no. 9, pp. 1481–1485, 2020.
- [28] X. Zhang and S. Zheng, "Grouping plane spiral electromagnetic waves for structured RF beams," in *Proc. 6th Asia-Pacific Conf. Antennas Propag. (APCAP)*, 2017, pp. 1–3.
- [29] H. Li, L. Song, and M. Debbah, "Energy efficiency of large-scale multiple antenna systems with transmit antenna selection," *IEEE Trans. Commun.*, vol. 62, no. 2, pp. 638–647, 2014.
- [30] S. Cui, A.J. Goldsmith, and A. Bahai, "Energy-efficiency of mimo and cooperative mimo techniques in senior networks," *IEEE J. Sel. Areas Commun.*, vol. 22, no. 6, pp. 1089–1098, 2004.
- [31] Z. Liu, J. Li, and D. Sun, "Circuit power consumption-unaware energy efficiency optimization for massive mimo systems," *IEEE Wireless Commun. Lett.*, vol. 6, no. 3, pp. 370–373, 2017.
- [32] Y. Li, M. Sheng, X. Wang, Y. Zhang, and J. Wen, "Max–min energy-efficient power allocation in interference-limited wireless networks," *IEEE Trans. Veh. Technol.*, vol. 64, no. 9, pp. 4321–4326, 2015.
- [33] L. Zhang, Y. Xin, and Y. Liang, "Weighted sum rate optimization for cognitive radio MIMO broadcast channels," *IEEE Trans. Wireless Commun.*, vol. 8, pp. 2950–2959, 2009.
- [34] J. J. M. Paul and H. Calamai, "Projected gradient methods for linearly constrained problems," *Math. Program.*, vol. 39, no. 1, pp. 93–116, 1987.
- [35] X. Huang and V. Y. Pan, "Fast rectangular matrix multiplication and applications," *J. Complexity*, vol. 14, no. 2, pp. 257–299, 1998.
- [36] X. Cao, R. Ma, L. Liu, H. Shi, Y. Cheng, and C. Sun, "A machine learning-based algorithm for joint scheduling and power control in wireless networks," *IEEE Internet Things J.*, vol. 5, no. 6, pp. 4308–4318, 2018.
- [37] J. Tang, Y. Song, C. Lin, W. Feng, Z. Chen, X. Zhang, K. Wong, "Energy Efficiency Optimization for Plane Spiral OAM Mode-Group Based MIMO-NOMA Systems," *Proc. Int. Conf. Artif. Intell. Commun. Netw.*, pp. 177–188, 2021.
- [38] Y. Liu, Z. Qin, M. ElKashlan, A. Nallanathan, and J. A. McCann, "Non-orthogonal multiple access in large-scale heterogeneous networks," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 12, pp. 2667–2680, 2017.