

Effective Throughput Maximization of NOMA with Practical Modulations

Yuan Wang, *Student Member, IEEE*, Jiaheng Wang, *Senior Member, IEEE*,
Vincent W.S. Wong, *Fellow, IEEE*, and Xiaohu You, *Fellow, IEEE*

Abstract—Non-orthogonal multiple access (NOMA) has been considered as a promising technology for future wireless communications. In most of the existing NOMA schemes, the ideal information rate based on Shannon capacity is used as the performance metric, assuming perfect successive interference cancellation (SIC) and Gaussian transmit signals without considering practical modulations. The implicit assumptions and the resulting schemes may lead to suboptimal performance in practical NOMA systems. In this paper, we consider multi-user multi-channel NOMA systems using practical quadrature amplitude modulation (QAM) with imperfect SIC. We aim to maximize a more practical performance metric, namely the *effective throughput*, which takes into account the data rate and error performance. To achieve this goal, we derive both the exact and approximate expressions of the effective throughput. We also formulate a joint resource optimization problem of the power allocation, channel assignment, and modulation selection to maximize the effective throughput. We develop an efficient power allocation solution by proposing a closed-form power allocation within channels and a waterfilling-form power budget allocation among channels. We also develop efficient channel assignment and modulation selection methods with the aid of matching theory and machine learning, respectively. Consequently, we provide an efficient joint resource allocation algorithm via iterative optimization to maximize the effective throughput. Numerical results are presented to verify the superiority of the proposed NOMA scheme over orthogonal multiple access (OMA) and other NOMA schemes.

Index Terms—Channel assignment, effective throughput, imperfect successive interference cancellation (SIC), quadrature amplitude modulation (QAM), non-orthogonal multiple access (NOMA), power allocation.

I. INTRODUCTION

Non-orthogonal multiple access (NOMA) has been considered as a promising technology to support massive connectivity due to its non-orthogonal characteristics [1]–[4]. Compared with the conventional orthogonal multiple access (OMA),

The work of Y. Wang, J. Wang, and X. You was supported by the National Key R&D Program of China under Grant 2021YFB2900300, the National Natural Science Foundation of China under Grants 61971130 and 61720106003, and the Natural Science Foundation on Frontier Leading Technology Basic Research Project of Jiangsu under Grant BK20212001. The work of V. W.S. Wong was supported by the Natural Sciences and Engineering Research Council of Canada (NSERC). (*Corresponding author: Jiaheng Wang.*)

Y. Wang, J. Wang, and X. You are with the National Mobile Communications Research Laboratory, School of Information Science and Engineering, Southeast University, Nanjing 210018, China, and also with Purple Mountain Laboratories, Nanjing 210018, China (e-mail: wang_yuan@seu.edu.cn; jhwang@seu.edu.cn; xhyu@seu.edu.cn).

V. Wong is with the Department of Electrical and Computer Engineering, The University of British Columbia, Vancouver, Canada (e-mail: vincentw@ece.ubc.ca).

NOMA can provide higher spectrum efficiency, better user fairness, and lower signalling cost [5]–[7]. In a power-domain NOMA system¹, the NOMA scheme is performed within each orthogonal channel by employing superposition coding at the transmitter and successive interference cancellation (SIC) at the receiver, so that users' signals can be multiplexed in the power domain. Hence, the resource allocation scheme, including power allocation and channel assignment, is a critical issue to fulfill the benefits of NOMA.

Many different approaches have been proposed to optimize the power allocation and channel assignment of NOMA systems [8]–[16]. The authors in [8] develop power allocation schemes under various performance criteria and a channel assignment algorithm based on the deferred acceptance method. Branch and bound approach is employed in [9] to jointly allocate resources for minimizing the transmit power subject to the data rate and outage probability constraints. In [10], the power weights and channel assignment are obtained through dynamic programming. In [11], energy-efficient resource allocation algorithms are developed by transforming and approximating the formulated problems into convex subproblems. A deep reinforcement learning framework is proposed in [12] for channel assignment to maximize the sum rate or minimal rate. In [13], the power allocation is pre-defined to follow the inverse proportional fairness criterion, while the channel assignment problem is formulated as a cooperative multi-agent game and is solved using deep deterministic policy gradient method. The power allocation of multiple-input multiple-output (MIMO)-NOMA is addressed through a communication deep neural network in [14] to maximize the sum rate and energy efficiency. An energy and delay cost minimization problem of Internet-of-things (IoT) systems is formulated in [15], where the subcarrier allocation and task scheduling are solved based on matching theory and reinforcement learning. In [16], the authors optimize the altitudes of unmanned aerial vehicles (UAVs) and improve the channel access to maximize the sum rate via constrained deep reinforcement learning.

Most of the existing works of NOMA, e.g., [8]–[17], aim at maximizing, e.g., the ideal sum rate, weighted sum rate, max-min fairness, or energy efficiency, based on Shannon capacity, which implicitly considered Gaussian transmit signals and perfect SIC without error propagation. However, in practical systems, signals are generally constrained to discrete modulation constellations and error propagation generally exists since

¹According to the multiplexing method, NOMA can be divided into power-domain NOMA and code-domain NOMA. In this paper, we focus on the former one. Hereafter, NOMA is used to refer to power-domain NOMA.

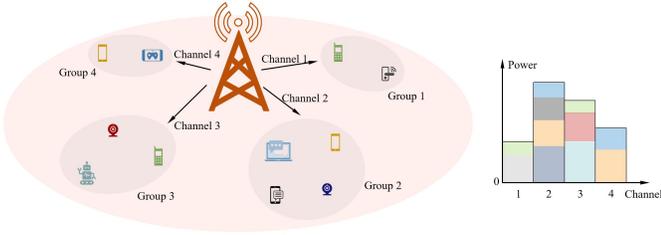


Figure 1: System model of a downlink NOMA system.

the interference cannot always be perfectly removed [2]. Some works have considered imperfect SIC for NOMA [18]–[24]. However, in those works, the residual interference caused by imperfect SIC is modeled as a continuously-valued Gaussian signal, which is an approximation without considering practical modulation schemes.

In this work, we consider, from a general point of view, multi-user multi-channel NOMA systems using the quadrature amplitude modulation (QAM) [25] with imperfect SIC. In the literature, only a few related works [26], [27] considered practical modulations along with imperfect SIC. Specifically, the work [27] focused on a two-user NOMA system and attempted to achieve the minimum error probability of the two users. In this paper, we consider a more general NOMA framework where multiple channels are available to multiple users who may adopt different QAM schemes and transmit using different power levels. More importantly, we introduce a more practical performance metric, namely the *effective throughput*, to the NOMA system design, which accounts for the correctly transmitted bit rate and thus takes into account both the data rate and error performance.

The aim of this work is to maximize the effective throughput of a multi-user multi-channel NOMA system while taking into account practical QAMs and imperfect SIC. There are several challenges to achieve this goal. First, the expression of the effective throughput, as the objective function, has to be derived, which is, however, quite complicated, due to the consideration of QAMs and imperfect SIC. Second, the effective throughput maximization problem turns out to be a joint resource optimization problem of the power allocation, channel assignment and modulation selection, which belongs to the class of mixed integer programs (MIPs). Third, the nonlinear and complex nature of the effective throughput expression makes each part, i.e., power allocation, channel assignment, or modulation selection, a difficult job. In the rest of this paper, we address these issues and develop efficient methods to maximize the effective throughput of the practical NOMA system. The main contributions of this paper are summarized as follows.

- We consider the effective throughput maximization of the multi-user multi-channel NOMA system considering both the data rate and error performance along with imperfect SIC and practical QAM schemes, which is formulated as a joint power allocation, channel assignment, and modulation selection problem.
- We derive the exact expression of the effective throughput

as a function of the transmit power, channel gain, and modulation scheme, and also provide a simpler lower bound of it for facilitating the system design.

- We decompose the power allocation problem into two subproblems. We obtain an analytical power allocation solution by providing a closed-form power allocation within channels and a waterfilling-form power budget allocation among channels. The proposed power allocation schemes can avoid allocating all power to the strongest user and thus can improve user fairness.
- We transform the channel assignment into a two-sided matching problem with peer effects and develop efficient channel assignment algorithms based on swap matching policies.
- We transform the modulation selection into a classification problem and propose a deep neural network (DNN)-based modulation selection method.
- A joint resource allocation algorithm is proposed via iteratively optimizing power allocation, channel assignment, and modulation selection for solving the effective throughput maximization problem.
- Representative numerical results are provided, showing that the proposed NOMA scheme outperforms the existing OMA and other NOMA schemes in terms of the effective throughput.

The rest of this paper is organized as follows. The system model is introduced in Section II. In Section III, we derive the expression of the effective throughput and address the power allocation subproblem. In Section IV, we investigate the channel assignment and modulation selection subproblems. Numerical results are presented in Section V. Conclusions are drawn in Section VI.

II. SYSTEM MODEL

A. Signal Model

In this paper, we consider a downlink NOMA system, where a single-antenna base station (BS) serves N single-antenna users² through K wireless channels, which can be frequency bands, time slots or spreading codes. We denote the n th user by U_n ($n = 1, \dots, N$) and the k th channel by C_k ($k = 1, \dots, K$). As depicted in Fig. 1, the users are divided into K groups with each group assigned with a channel. Signals transmitted on different channels are assumed to have no interference with each other, while signals on the same channel are multiplexed in the power domain. We assume that the BS has full knowledge of the channel state information (CSI).

At the BS, after superposition coding, the transmit signal on channel C_k can be written as

$$x_k = \sum_{n=1}^N d_{k,n} \sqrt{p_{k,n}} s_n, \quad (1)$$

²According to [28], by carefully designing the precoding and detection matrices, the inter-cluster interference of MIMO-NOMA can be canceled completely. The MIMO-NOMA system can be reduced into several single-input single-output (SISO) NOMA systems. Thus, in this paper, we focus on the single-antenna NOMA systems.

where $s_n \in \mathbb{C}$ is the intended signal of user U_n , $d_{k,n}$ is a binary variable and represents whether channel C_k is assigned to user U_n , i.e.,

$$d_{k,n} = \begin{cases} 0, & C_k \text{ is not assigned to } U_n \\ 1, & C_k \text{ is assigned to } U_n. \end{cases} \quad (2)$$

Moreover, $p_{k,n} \geq 0$ is the power allocated to user U_n on channel C_k , subject to the total power constraint $\sum_{k=1}^K \sum_{n=1}^N d_{k,n} p_{k,n} \leq P$, where P is the total power budget. At the receiver, if assigned with channel C_k , user U_n receives the signal

$$y_{k,n} = g_{k,n} x_k + z_n, \quad (3)$$

where $g_{k,n} \in \mathbb{C}$ is the channel coefficient of channel C_k between the BS and user U_n , which includes both large-scale path loss and small-scale fading, and $z_n \in \mathbb{C}$ is the additive white Gaussian noise (AWGN) of user U_n , with zero mean and variance σ_n^2 . To facilitate further analysis, we define the normalized channel gain as $h_{k,n} \triangleq |g_{k,n}|^2 / \sigma_n^2 > 0$.

SIC is employed at the receiver to decode the intended signal s_n from $y_{k,n}$. For those users which are assigned with the same channel, each user decodes and removes the signals of the other users with higher power successively, and then decodes its own signal, while treating the signals of those users with lower power as interference. Considering the complexity of SIC, the decoding delay, and the co-channel interferences are proportional to the number of users on each channel [8], [11], [20], [29], we limit the number of users within each channel to be at most two. The total number of users N is less than or equal to $2K$. Specifically, we assume that each user occupies one channel. The situation where each user can occupy more than one channel can be extended from our system model by treating the user assigned with multiple channels as multiple virtual users assigned with single channel.

B. Effective Throughput

To evaluate the NOMA system performance, most of the previous works adopt metrics based on the ideal information rate, assuming perfect SIC without considering the error performance [8]–[17]. Although some works take into account imperfect SIC, they do not consider the fact that the transmit signals in practice are generally constrained to discrete modulation constellations [18]–[24]. Thus, the resulting schemes of the models in [8]–[24] may be different when the algorithms are deployed in practical systems.

To tackle this issue, in this paper, we use the effective throughput as the performance metric, considering both the error performance and the data rate along with imperfect SIC and practical modulation schemes. Suppose the bandwidth occupied by each channel is equal to B . Consider that user U_n assigned with channel C_k employs M_n^2 -QAM. According to [30], the symbol rate of user U_n is also equal to B . Thus, the effective symbol rate (which measures the correctly transmitted symbol rate [23], [31], [32]) of user U_n can be expressed as $B(1 - \epsilon_{k,n})$, where $\epsilon_{k,n}$ is the symbol error rate

(SER) of user U_n on channel C_k . On this basis, we define the effective throughput of user U_n on channel C_k as

$$J_{k,n} \triangleq B(1 - \epsilon_{k,n}) \log_2 M_n^2, \quad (4)$$

to measure the correctly transmitted bit rate.

Note that it is reasonable to use the SER, and not the bit error rate (BER), in determining the effective throughput. If the BER is used in (4), the effective throughput of M_n^2 -QAM has the range $\frac{B}{2} \log_2 M_n^2 \leq J_{k,n} < B \log_2 M_n^2$. Then, the maximum effective throughput of M_n^2 -QAM, i.e., $B \log_2 M_n^2$, is equal to the minimum effective throughput of M_n^4 -QAM, i.e., $\frac{B}{2} \log_2 M_n^4 = B \log_2 M_n^2$. Under this assumption, the scheduler will select the highest modulation order to maximize the effective throughput no matter how poor the real error performance is.

C. Problem Formulation

From (4), the effective throughput is a function of the bandwidth, the modulation scheme, and the SER, which further depends on the channel quality and transmit power. Thus, the effective throughput maximization problem can be formulated as a joint power allocation, channel assignment, and modulation selection problem:

$$\mathcal{P}_1 : \underset{\{p_{k,n}, d_{k,n}, M_n\}}{\text{maximize}} \quad J \triangleq \sum_{k=1}^K \sum_{n=1}^N d_{k,n} J_{k,n} \quad (5a)$$

$$\text{subject to} \quad d_{k,n} \in \{0, 1\}, \quad k = 1, \dots, K, \quad n = 1, \dots, N \quad (5b)$$

$$\sum_{k=1}^K d_{k,n} = 1, \quad n = 1, \dots, N \quad (5c)$$

$$\sum_{n=1}^N d_{k,n} \leq 2, \quad k = 1, \dots, K \quad (5d)$$

$$\sum_{k=1}^K \sum_{n=1}^N d_{k,n} p_{k,n} \leq P \quad (5e)$$

$$p_{k,n} \geq 0, \quad k = 1, \dots, K, \quad n = 1, \dots, N \quad (5f)$$

$$M_n \in \mathcal{M}, \quad n = 1, \dots, N, \quad (5g)$$

where J is the effective system throughput. Constraints (5c) and (5d) indicate that each user occupies one channel and each channel can be assigned to at most two users. We refer to the channels assigned to only one user as the *single-user channels*, the channels assigned to two users as the *two-user channels*, and the channels not being assigned to any users as the *idle channels*. Constraint (5e) is the total power constraint. \mathcal{M} represents the set of all possible modulation orders. In this paper, we consider four practical modulation schemes: 4-QAM, 16-QAM, 64-QAM, and 256-QAM. Thus, $\mathcal{M} = \{2, 4, 8, 16\}$. Problem \mathcal{P}_1 is solved for each symbol to obtain the optimal performance. For slow-fading channels, the resource allocation can be updated in a longer time scale in order to reduce the complexity.

Problem \mathcal{P}_1 is an MIP, which is NP-hard [33], [34]. To find the globally optimal solution, one has to employ the

$$\begin{aligned}
\phi_{k,I} = & \frac{1}{M_{k,I}M_{k,II}} \sum_{m=1}^{M_{k,I}} \sum_{j=1}^{M_{k,II}} \sum_{l=1}^{M_{k,II}} \left[f(m,1)Q \left(\sqrt{\frac{3h_{k,I}p_{k,I}}{M_{k,I}^2-1}} + 2(j-l)\sqrt{\frac{3h_{k,I}p_{k,II}}{M_{k,II}^2-1}} \right) \right. \\
& \left. + f(m,M_{k,I})Q \left(\sqrt{\frac{3h_{k,I}p_{k,I}}{M_{k,I}^2-1}} - 2(j-l)\sqrt{\frac{3h_{k,I}p_{k,II}}{M_{k,II}^2-1}} \right) \right] \\
& \times \left[1 - f(l,1)Q \left((2m-M_{k,I}-1)\sqrt{\frac{3h_{k,I}p_{k,I}}{M_{k,I}^2-1}} + (2j+1-2l)\sqrt{\frac{3h_{k,I}p_{k,II}}{M_{k,II}^2-1}} \right) \right. \\
& \left. - f(l,M_{k,II})Q \left((M_{k,I}+1-2m)\sqrt{\frac{3h_{k,I}p_{k,I}}{M_{k,I}^2-1}} + (2l-2j+1)\sqrt{\frac{3h_{k,I}p_{k,II}}{M_{k,II}^2-1}} \right) \right]. \tag{7}
\end{aligned}$$

exhaustive search that has a high computational complexity. In practice, an efficient way to solve an MIP is to tackle the continuous variables and the discrete variables alternatively and iteratively. In our case, problem \mathcal{P}_1 can be decomposed into two subproblems: a) optimizing power allocation for fixed channel assignment and modulation schemes; b) optimizing channel assignment and modulation schemes for fixed power allocation. The joint solution can be obtained by iteratively solving them. However, neither of these two subproblems is easy, due to the nonlinear and complicated relation of the effective throughput with the power allocation, channel assignment, and modulation schemes.

III. EFFECTIVE THROUGHPUT AND POWER ALLOCATION

In this section, we first derive the exact expression of the effective throughput. Then, a low-complexity power allocation algorithm is proposed by adopting proper approximation of the effective throughput and exploring the hidden convexity of the resulting power allocation problem.

A. Effective Throughput

In order to derive the expression of the effective throughput $J_{k,n}$, we need to determine $\epsilon_{k,n}$ of user U_n on channel C_k . According to [35], an M_n^2 -QAM can be decoupled into two M_n -pulse amplitude modulations (PAMs). When the signal-to-noise ratios (SNRs) and the bit mapping schemes of M_n^2 -QAM and M_n -PAM are identical, the SER of the M_n^2 -QAM can be expressed as

$$\epsilon_{k,n} = 1 - (1 - \phi_{k,n})^2 = 2\phi_{k,n} - \phi_{k,n}^2, \tag{6}$$

where $\phi_{k,n}$ is the SER of M_n -PAM. Next, we determine the SER $\phi_{k,n}$ and the effective throughput $J_{k,n}$ of the two-user channels and the single-user channels, respectively.

1) Effective Throughput of the Two-User Channels

Suppose that channel C_k is a two-user channel, and then, the SER depends on the SIC decoding order. We denote the two users on channel C_k as users $U_{k,I}$ and $U_{k,II}$, where user $U_{k,I}$ has a better channel condition than user $U_{k,II}$, i.e., the channel gains satisfy $h_{k,I} \geq h_{k,II}$. Assume that $U_{k,i}$ ($i = I, II$) employs $M_{k,i}^2$ -QAM and is allocated with the power $p_{k,i} \geq 0$. User $U_{k,I}$ decodes its signal directly, treating the signal of user $U_{k,II}$ as interference. User $U_{k,II}$ first decodes and cancels the signal of user $U_{k,I}$, and then decodes the signal of itself.

According to [27], the SER of the strong channel user, i.e., $U_{k,I}$, when adopting $M_{k,I}$ -PAM, is given in (7) at the top of this page, where $Q(x) = \frac{1}{\sqrt{2\pi}} \int_x^{+\infty} \exp(-\frac{t^2}{2}) dt$ and $f(\cdot, \cdot)$ is defined as

$$f(x, y) \triangleq \begin{cases} 0, & x = y \\ 1, & x \neq y. \end{cases} \tag{8}$$

The SER of the weak channel user, i.e., user $U_{k,II}$, when adopting $M_{k,II}$ -PAM, is

$$\begin{aligned}
\phi_{k,II} = & \frac{2M_{k,II} - 2}{M_{k,I}M_{k,II}} \\
& \times \sum_{m=1}^{M_{k,I}} Q \left((M_{k,I} + 1 - 2m) \sqrt{\frac{3h_{k,II}p_{k,I}}{M_{k,I}^2 - 1}} + \sqrt{\frac{3h_{k,II}p_{k,II}}{M_{k,II}^2 - 1}} \right). \tag{9}
\end{aligned}$$

With the expressions of $\phi_{k,I}$ and $\phi_{k,II}$, we can obtain the effective throughput $J_{k,I}$ and $J_{k,II}$ via (4) and (6), as shown in the following result.

Proposition 1. *The effective throughput of user $U_{k,i}$ ($i = I, II$) on the two-user channel is*

$$J_{k,i} = B(1 - \phi_{k,i})^2 \log_2 M_{k,i}^2, \tag{10}$$

where $\phi_{k,I}$ and $\phi_{k,II}$ are given in (7) and (9), respectively.

2) Effective Throughput of the Single-User Channels

Suppose that channel C_k is a single-user channel. We denote the user assigned to channel C_k as $U_{k,0}$. According to [30], the SER of user $U_{k,0}$, when adopting $M_{k,0}$ -PAM, is

$$\phi_{k,0} = \frac{2(M_{k,0} - 1)}{M_{k,0}} Q \left(\sqrt{\frac{3h_{k,0}p_{k,0}}{M_{k,0}^2 - 1}} \right), \tag{11}$$

where $h_{k,0}$ and $p_{k,0}$ are the channel gain and the power of user $U_{k,0}$, respectively. Thus, the effective throughput $J_{k,0}$, when user $U_{k,0}$ adopts $M_{k,0}^2$ -QAM, can be obtained via (4) and (6), as shown in the following result.

Proposition 2. *The effective throughput of user $U_{k,0}$ on the single-user channel is*

$$J_{k,0} = B(1 - \phi_{k,0})^2 \log_2 M_{k,0}^2, \tag{12}$$

where $\phi_{k,0}$ is given in (11).

$$\begin{aligned}
J_k \geq J_k^L &= 2B \log_2 M_{k,I} + 2B \log_2 M_{k,II} - \frac{8B(M_{k,I} - 1) \log_2 M_{k,I}}{M_{k,I}} Q \left(\sqrt{\frac{3h_{k,I} p_{k,I}}{M_{k,I}^2 - 1}} \right) \\
&\quad - \left(\frac{8B(M_{k,II} - 1) \log_2 M_{k,II}}{M_{k,II}} + \frac{12B(M_{k,I} - 1)(M_{k,II} - 1) \log_2 M_{k,I}}{M_{k,I}} \right) \\
&\quad \times Q \left((1 - M_{k,I}) \sqrt{\frac{3h_{k,II} p_{k,I}}{M_{k,I}^2 - 1}} + \sqrt{\frac{3h_{k,II}(q_k - p_{k,I})}{M_{k,II}^2 - 1}} \right). \tag{19}
\end{aligned}$$

According to Propositions 1 and 2, the effective throughput is a function of the bandwidth, the channel quality, the modulation scheme, and the transmit power. Compared with the ideal information rate, the effective throughput considers imperfect SIC and practical modulations, and thus, is more practical. Furthermore, in (10) and (12), the terms $(1 - \phi_{k,i})^2$ and $(1 - \phi_{k,0})^2$ reflect the error performance, whereas the terms $B \log_2 M_{k,i}^2$ and $B \log_2 M_{k,0}^2$ reflect the bit rate. Consequently, the effective throughput takes into account both the reliability and transmission rate of the NOMA system.

B. Power Allocation Subproblems

Suppose that the channel assignment and modulation schemes are fixed. For convenience, we denote the index sets of the single-user channels and the two-user channels by \mathcal{K}^\dagger and \mathcal{K}^\ddagger , respectively. Then, the power allocation optimization subproblem can be formulated as

$$\mathcal{P}_2 : \begin{aligned} &\text{maximize}_{\{p_{k,0}, p_{k,I}, p_{k,II} \geq 0\}} J = \sum_{k \in \mathcal{K}^\dagger} J_{k,0} + \sum_{k \in \mathcal{K}^\ddagger} (J_{k,I} + J_{k,II}) \end{aligned} \tag{13a}$$

$$\text{subject to } \sum_{k \in \mathcal{K}^\dagger} p_{k,0} + \sum_{k \in \mathcal{K}^\ddagger} (p_{k,I} + p_{k,II}) \leq P. \tag{13b}$$

Note that problem \mathcal{P}_2 is a nonconvex optimization problem. Next, we address the nonconvexity of problem \mathcal{P}_2 through decomposition and approximation.

Specifically, we define the effective throughput of channel C_k as

$$J_k \triangleq \begin{cases} J_{k,0}, & k \in \mathcal{K}^\dagger \\ J_{k,I} + J_{k,II}, & k \in \mathcal{K}^\ddagger, \end{cases} \tag{14}$$

and introduce the power budget of channel C_k as

$$q_k \triangleq \begin{cases} p_{k,0}, & k \in \mathcal{K}^\dagger \\ p_{k,I} + p_{k,II}, & k \in \mathcal{K}^\ddagger, \end{cases} \tag{15}$$

which satisfies $q_k \geq 0$ and $\sum_{k=1}^K q_k \leq P$. Consequently, problem \mathcal{P}_2 can be equivalently decomposed into a power budget allocation subproblem among channels:

$$\mathcal{P}_3 : \text{maximize}_{\{q_k \geq 0\}} J = \sum_{k=1}^K J_k \tag{16a}$$

$$\text{subject to } \sum_{k=1}^K q_k \leq P, \tag{16b}$$

and a series of power allocation subproblems within the two-user channels:

$$\mathcal{P}_{4,k} (k \in \mathcal{K}^\ddagger) : \begin{aligned} &\text{maximize}_{\{p_{k,I}, p_{k,II} \geq 0\}} J_k = J_{k,I} + J_{k,II} \end{aligned} \tag{17a}$$

$$\text{subject to } p_{k,I} + p_{k,II} = q_k. \tag{17b}$$

In the following subsection, we provide a closed-form solution to problem $\mathcal{P}_{4,k}$, and show that problem \mathcal{P}_3 can be transformed into a convex problem and be solved analytically.

C. Power Allocation within Channels

Substituting $p_{k,II}$ by $q_k - p_{k,I}$, problem $\mathcal{P}_{4,k}$ can be equivalently reformulated as

$$\mathcal{P}_{5,k} (k \in \mathcal{K}^\ddagger) : \begin{aligned} &\text{maximize}_{\{p_{k,I} \in [0, q_k]\}} J_k = J_{k,I} + J_{k,II}. \end{aligned} \tag{18}$$

From Proposition 1, one can find that the expressions of $J_{k,I}$ and $J_{k,II}$ are nonlinear. The exhaustive search³ is required to find the optimal solution to $\mathcal{P}_{5,k}$. Next, we derive the lower bound of J_k as the approximation and obtain the near-optimal solution via maximizing the lower bound of J_k .

Proposition 3. *The effective throughput J_k is lower bounded by (19) at the top of this page.*

Proof. The proof is provided in Appendix A. \square

In Proposition 3, the lower bound of J_k is given in a simpler expression compared with the exact expression of J_k . In the following result, we take the lower bound as an approximation of J_k and obtain a closed-form power allocation scheme via maximizing the lower bound.

Theorem 1. *For a sufficiently large power budget q_k , the power allocation within the two-user channels that maximizes the lower bound of J_k is given by*

$$p_{k,I}^* = \frac{(M_{k,I}^2 - 1)q_k}{(M_{k,I} - 1 + t_k)^2(M_{k,II}^2 - 1) + M_{k,I}^2 - 1}, \tag{20}$$

$$p_{k,II}^* = q_k - p_{k,I}^*, \tag{21}$$

where $t_k \triangleq \sqrt{h_{k,I}/h_{k,II}} \geq 1$ is defined as the channel gain ratio of channel C_k .

Proof. The proof is provided in Appendix B. \square

In Theorem 1, we provide a closed-form power allocation scheme within the two-user channels, which can be determined

³By discretizing the power into different levels, the optimal power allocation can be found via an exhaustive search.

by the modulation orders and the channel gain ratio. One can observe that more power should be allocated to a user when its modulation order increases or when its channel gain decreases. Numerical results demonstrate that the proposed power allocation scheme can well approximate the optimal one obtained via an exhaustive search.

D. Power Budget Allocation among Channels

With the closed-form power allocation scheme within channels, the power budget allocation problem among channels, i.e., problem \mathcal{P}_3 , can be simplified into

$$\mathcal{P}_6 : \underset{\{q_k \geq 0\}}{\text{maximize}} \quad \tilde{J} \triangleq \sum_{k=1}^K \tilde{J}_k \quad (22a)$$

$$\text{subject to} \quad \sum_{k=1}^K q_k \leq P, \quad (22b)$$

where, for the two-user channels ($k \in \mathcal{K}^\ddagger$), $\tilde{J}_k \triangleq J_k^L|_{p_{k,I}=p_{k,I}^*, p_{k,II}=p_{k,II}^*}$ is the lower bound of the effective throughput J_k . According to (19) and (20), we have

$$\begin{aligned} \tilde{J}_k &= 2B \log_2 M_{k,I} + 2B \log_2 M_{k,II} \\ &\quad - Q \left(\sqrt{\frac{3h_{k,I}q_k}{(M_{k,I}-1+t_k)^2(M_{k,II}^2-1)+M_{k,I}^2-1}} \right) \\ &\quad \times \left(\frac{4B(M_{k,I}-1)(3M_{k,II}-1)\log_2 M_{k,I}}{M_{k,I}} \right. \\ &\quad \left. + \frac{8B(M_{k,II}-1)\log_2 M_{k,II}}{M_{k,II}} \right), \quad k \in \mathcal{K}^\ddagger. \end{aligned} \quad (23)$$

For the single-user channels ($k \in \mathcal{K}^\dagger$), the effective throughput $J_{k,0}$ can be obtained via (12), which, however, is still complicated. According to [30], $J_{k,0}$ has the lower bound

$$J_{k,0} = B(1 - \phi_{k,0})^2 \log_2 M_{k,0}^2 \geq B(1 - 2\phi_{k,0}) \log_2 M_{k,0}^2, \quad (24)$$

which is a tight approximation for a sufficiently large power budget. Thus, we use the lower bound of $J_{k,0}$ to approximate \tilde{J}_k as

$$\begin{aligned} \tilde{J}_k &\triangleq B(1 - 2\phi_{k,0}) \log_2 M_{k,0}^2 \\ &= 2B \log_2 M_{k,0} - \frac{8B(M_{k,0}-1)\log_2 M_{k,0}}{M_{k,0}} \\ &\quad \times Q \left(\sqrt{\frac{3h_{k,0}p_{k,0}}{M_{k,0}^2-1}} \right), \quad k \in \mathcal{K}^\dagger. \end{aligned} \quad (25)$$

Combining (23) and (25), \tilde{J}_k can be given by

$$\tilde{J}_k = \begin{cases} 2B \log_2 M_{k,0} - \alpha_k Q(\sqrt{\beta_k q_k}), & k \in \mathcal{K}^\dagger \\ 2B \log_2 M_{k,I} + 2B \log_2 M_{k,II} \\ - \alpha_k Q(\sqrt{\beta_k q_k}), & k \in \mathcal{K}^\ddagger, \end{cases} \quad (26)$$

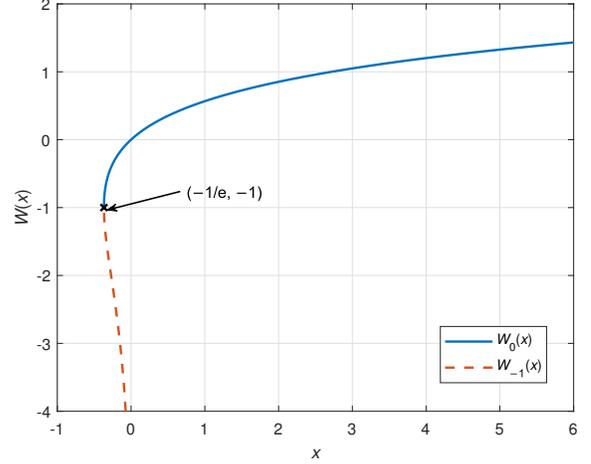


Figure 2: The Lambert W function.

where

$$\alpha_k \triangleq \begin{cases} \frac{8B(M_{k,0}-1)\log_2 M_{k,0}}{M_{k,0}}, & k \in \mathcal{K}^\dagger \\ \frac{4B(M_{k,I}-1)(3M_{k,II}-1)\log_2 M_{k,I}}{M_{k,I}} \\ + \frac{8B(M_{k,II}-1)\log_2 M_{k,II}}{M_{k,II}}, & k \in \mathcal{K}^\ddagger, \end{cases} \quad (27)$$

$$\beta_k \triangleq \begin{cases} \frac{3h_{k,0}}{M_{k,0}^2-1}, & k \in \mathcal{K}^\dagger \\ \frac{3h_{k,I}}{(M_{k,I}-1+t_k)^2(M_{k,II}^2-1)+M_{k,I}^2-1}, & k \in \mathcal{K}^\ddagger. \end{cases} \quad (28)$$

To solve problem \mathcal{P}_6 , we first provide the following useful results.

Lemma 1. \tilde{J}_k is monotonically increasing with respect to q_k for $k = 1, \dots, K$.

Lemma 2. \tilde{J}_k is a concave function of q_k for $k = 1, \dots, K$.

Proof. The proofs of Lemmas 1 and 2 are provided in Appendix C. \square

Lemma 1 complies with the intuition that the larger the power budget is, the higher the effective throughput will be. Lemma 2 indicates that problem \mathcal{P}_6 is a convex optimization problem, whose solution can be found efficiently via standard convex optimization tools, e.g., CVX. Nevertheless, we are able to provide an analytical solution to problem \mathcal{P}_6 . For this purpose, the following definition is introduced.

Definition 1. The Lambert W function [36], denoted by $W(\cdot)$, is a function satisfying

$$W(x) \exp(W(x)) = x. \quad (29)$$

As shown in Fig. 2, $W(x)$ has two possible values when $x \in [-1/e, 0)$. The function is divided into two branches, taking the point $(-1/e, -1)$ as the cut-off point. The branch satisfying $W(x) \leq -1$ is denoted by $W_{-1}(x)$. The branch satisfying $W(x) \geq -1$ is denoted by $W_0(x)$, which is called the principal branch. Note that $W_0(x)$ is a monotonically increasing function when $x \geq -1/e$.

Theorem 2. The optimal solution to problem \mathcal{P}_6 is given by

$$q_k^* = \frac{W_0\left(\frac{\alpha_k^2 \beta_k^2}{8\pi\lambda}\right)}{\beta_k}, \quad (30)$$

where λ is chosen such that $\sum_{k=1}^K q_k^* = P$ and $W_0(\cdot)$ is the principal branch of Lambert W function.

Proof. The proof of Theorem 2 is provided in Appendix D. \square

Theorem 2 provides an analytical power budget allocation scheme among channels in waterfilling form. From (30), q_k^* is monotonically decreasing with respect to λ . Thus, the value of λ can be efficiently found via a one-dimensional search. In the following corollary, we provide the upper and lower bounds of λ to accelerate the search.

Corollary 1. λ is upper bounded by

$$\lambda \leq \lambda^U = \max_{k=1, \dots, K} \left\{ \frac{\alpha_k^2 \beta_k}{8\pi P/K} \exp(-\beta_k P/K) \right\}, \quad (31)$$

and lower bounded by

$$\lambda \geq \lambda^L = \max_{k=1, \dots, K} \left\{ \frac{\alpha_k^2 \beta_k}{8\pi P} \exp(-\beta_k P) \right\}. \quad (32)$$

Proof. The proof is provided in Appendix E. \square

In Theorems 1 and 2, we provide the analytical power allocation expressions within and among channels, respectively, which are both determined by the modulation orders and the channel gains. The former one is given in closed form and the latter one is given in a waterfilling fashion. By combining Theorems 1 and 2, an efficient power allocation scheme can be analytically obtained, which, though including some proper approximation, is able to achieve near optimal performance. The proposed power allocation scheme given in waterfilling form can be efficiently obtained via a one-dimensional search. Its computational complexity linearly increases with respect to the number of channels, i.e., $\mathcal{O}(K)$ [37].

E. User Fairness of the Power Allocation

In most of the NOMA works based on the ideal information rate, e.g., [8]–[17], the system tends to allocate more power to the stronger channel user. If the fairness constraints (e.g., the minimum data rate constraint or the power order constraint) are not considered, all power will be allocated to the strongest user in an unfair manner [17]. This is a critical drawback of using the ideal information rate in NOMA designs. The reason is that the ideal information rate does not consider the error performance and assumes infinite block length. Meanwhile, SIC is assumed to be perfect, which may also contribute to allocating more power to the stronger channel user.

In this paper, we introduce the effective throughput as the performance metric by considering both the error performance and the data rate along with imperfect SIC and practical modulation schemes, and thus, can overcome the drawback of the ideal information rate and facilitate user fairness. In the following results, we analyze user fairness of the power

allocation within channels and the power budget allocation among channels, respectively.

Proposition 4. In the power allocation scheme proposed in Theorem 1, $p_{k,I}^*$ is monotonically decreasing and $p_{k,II}^*$ is monotonically increasing with respect to the channel gain ratio t_k .

Proof. The monotonicity can be obtained via taking the derivative of $p_{k,I}^*$ and $p_{k,II}^*$ with respect to t_k . \square

Proposition 5. When the channel gain of the user (or user pair) on channel C_k is very strong, i.e., $h_{k,0} \rightarrow \infty$ (or $h_{k,I} \rightarrow \infty$), the proposed power budget allocation scheme in Theorem 2 satisfies $q_k^* \rightarrow 0$.

Proof. The proof is provided in Appendix F. \square

Propositions 4 and 5 imply that, the proposed power allocation schemes tend to allocate smaller fraction of total power to the strong user (or user pair) even without the fairness constraints. This is different from the works based on the ideal information rate [8]–[17]. The proposed power allocation schemes can avoid allocating all power to the strongest user (or user pair) and thus can improve user fairness.

IV. CHANNEL ASSIGNMENT AND MODULATION SELECTION

Now, we consider the channel assignment and modulation selection for effective throughput maximization with given power allocation. In this case, problem \mathcal{P}_1 reduces to

$$\mathcal{P}_7 : \underset{\{d_{k,n}, M_n\}}{\text{maximize}} \quad J = \sum_{k=1}^K \sum_{n=1}^N d_{k,n} J_{k,n} \quad (33a)$$

$$\text{subject to} \quad \text{constraints (5b) – (5d), (5g)}. \quad (33b)$$

It is a combinatorial problem whose optimal solution has to be found via exhaustive search. When the number of users increases, the complexity of exhaustive search soars. In this section, we develop efficient discrete optimization methods based on matching theory and machine learning to determine the channel assignment and modulation selection.

A. Channel Assignment

We find that the channel assignment problem in NOMA can be regarded as a bipartite many-to-one matching game with peer effects [38]. The channels and the users can be regarded as two distinct sets of players. The users have peer effects due to the co-channel interference caused by superposition coding. Thus, it can be addressed by the swap matching method [34], [39], [40], which provides an efficient and low-complexity solution to the combinatorial problem. To describe the channel assignment scheme, we define a $K \times N$ matching matrix D with the binary elements $d_{k,n}$ to match the K channels C_k ($k = 1, \dots, K$) to the N users U_n ($n = 1, \dots, N$), satisfying the constraints (33b).

For a matching matrix D , a swap operation is to choose two users on different channels, e.g., user U_n on channel C_k and user $U_{n'}$ on channel $C_{k'}$ (where $k \neq k'$ and $n \neq n'$),

and to exchange the channels of the two users. We denote the original user-channel assignment pair by $(U_n, C_k), (U_{n'}, C_{k'})$. Then, the new pair after swap is $(U_{n'}, C_k), (U_n, C_{k'})$. The new matching matrix is denoted by D^{new} .

1) Channel Assignment Based on the Swap-Blocking Policy

Assumed that all users are selfish and aim to maximize their own effective throughput [34], [39], [40]. A swap operation will not be executed unless all the users on channels C_k and $C_{k'}$ approve. We give the definition of swap-blocking (SB) pair to decide whether a swap operation should be executed and develop a channel assignment algorithm based on the SB policy.

Definition 2. A pair $(U_n, C_k), (U_{n'}, C_{k'})$ is defined as a swap-blocking pair if and only if

- (a) $\forall \zeta \in \mathcal{U}, J_{(\zeta)}(D^{\text{new}}) \geq J_{(\zeta)}(D)$,
- (b) $\exists \zeta \in \mathcal{U}, J_{(\zeta)}(D^{\text{new}}) > J_{(\zeta)}(D)$,

where \mathcal{U} represents the set of all the users on channels C_k and $C_{k'}$. $J_{(\zeta)}(D)$ and $J_{(\zeta)}(D^{\text{new}})$ represent the effective throughput of user ζ under the matching scheme D and D^{new} , respectively.

In the above definition, condition (a) implies that the performance of all the users should not be sacrificed after the swap operation and condition (b) implies that at least the performance of one user should be improved. The swap operation is executed if and only if the pair forms an SB pair, called the SB policy. Then, the channel assignment can be optimized through searching for SB pairs and executing the swap operation until no SB pair exists.

2) Channel Assignment Based on the Centralized Swap-Blocking Policy

In the channel assignment algorithm based on the SB policy, the swap operation is executed only when the effective throughput of all the users does not decrease. This means the SB policy may miss some swap operations which increase the effective system throughput by degrading the performance of some specific users. To overcome this limitation, we define the centralized swap-blocking (CSB) pair and propose a channel assignment algorithm based on the CSB policy.

Definition 3. A pair $(U_n, C_k), (U_{n'}, C_{k'})$ is defined as a centralized swap-blocking pair if and only if $J(D^{\text{new}}) > J(D)$, where $J(D)$ and $J(D^{\text{new}})$ represent the effective system throughput under the matching scheme D and D^{new} , respectively.

Then, in the channel assignment algorithm based on the CSB policy, we execute the swap operation if and only if the pair forms a CSB pair, and the effective system throughput is increased. The details of the channel assignment algorithm based on the SB or CSB policy are shown in Algorithm 1. In the following results, we prove that the proposed CSB policy outperforms the SB policy for the channel assignment in NOMA.

Lemma 3. *If a pair forms an SB pair, then it must form a CSB pair.*

Proof. If a pair $(U_n, C_k), (U_{n'}, C_{k'})$ forms an SB pair, according to Definition 2, conditions (a) and (b) are satisfied.

Algorithm 1 Channel assignment algorithm based on the SB or CSB policy

- 1) **Initialization:**
Generate a channel assignment matrix D randomly.
- 2) **Repeat**
- 3) Choose a pair $(U_n, C_k), (U_{n'}, C_{k'})$ from D arbitrarily;
- 4) **If** the pair $(U_n, C_k), (U_{n'}, C_{k'})$ satisfies condition I*:
- 5) Execute the swap operation, update D ;
- 6) **End if**
- 7) **Until** condition II* is satisfied.
- 8) **Return** the matrix D .

*For the SB policy, condition I is that the pair $(U_n, C_k), (U_{n'}, C_{k'})$ forms an SB pair and condition II is that no SB pair exists. For the CSB policy, condition I is that the pair $(U_n, C_k), (U_{n'}, C_{k'})$ forms a CSB pair and condition II is that no CSB pair exists.

Then, we have

$$J(D^{\text{new}}) - J(D) = \sum_{\zeta \in \mathcal{U}} J_{(\zeta)}(D^{\text{new}}) - \sum_{\zeta \in \mathcal{U}} J_{(\zeta)}(D) > 0, \quad (34)$$

and this pair forms a CSB pair. Conversely, if a pair $(U_n, C_k), (U_{n'}, C_{k'})$ forms a CSB pair, this pair is not necessarily an SB pair. For example, if $J_{(U_n)}(D^{\text{new}}) - J_{(U_n)}(D) = 10$ and $J_{(\zeta)}(D^{\text{new}}) - J_{(\zeta)}(D) = -1$, for $\zeta \in \mathcal{U}$ and $\zeta \neq U_n$, then this pair is a CSB pair, but not an SB pair. \square

Theorem 3. *The effective system throughput achieved by the CSB-based channel assignment algorithm is always no less than that achieved by the SB-based algorithm.*

Proof. For an arbitrary pair, it has three possible situations: (a) it forms an SB pair, and then, according to Lemma 3, it must form a CSB pair; (b) it forms a CSB pair but is not an SB pair; (c) it neither forms an SB pair nor forms a CSB pair. For situation (a) or (c), the SB and CSB policies choose the same action, i.e., execute the swap operation for (a) and reject the swap operation for (c). For situation (b), the CSB policy executes the swap operation and J increases, but the SB policy rejects the swap operation and J remains unchanged. During the search process, the CSB policy will never perform worse than the SB policy, and once a pair is in situation (b), the CSB policy outperforms the SB policy. Theorem 3 is proven. \square

According to Lemma 3 and Theorem 3, for an arbitrary matching matrix D , the number of CSB pairs is generally more than that of SB pairs. The CSB policy outperforms the SB policy in terms of the effective system throughput at the cost of more iterations. Next, we analyze the convergence of Algorithm 1. Note that J increases each time an SB or CSB pair is being swapped. Meanwhile, J is upper bounded by a finite value. Furthermore, the number of potential SB or CSB pairs is finite since the number of users is limited. Therefore, Algorithm 1 converges and terminates after a finite number of swap operations.

The complexity of Algorithm 1 can be determined as follows. We denote the number of iterations by N_{iter} . During each iteration, at most $N(N-1)$ times of swap operations are required to check the termination condition, where N is the number of users. Therefore, the computational complexity of Algorithm 1 is $\mathcal{O}(N^2 N_{\text{iter}})$. For comparison, if

$$\begin{aligned}
& J_k|_{p_{k,I}=p_{k,I}^*, p_{k,II}=p_{k,II}^*} \\
&= 2B \log_2 M_{k,I} (1 - \phi_{k,I})^2 + 2B \log_2 M_{k,II} (1 - \phi_{k,II})^2 \\
&= 2B \log_2 M_{k,I} \left\{ 1 - \frac{1}{M_{k,I} M_{k,II}} \sum_{m=1}^{M_{k,I}} \sum_{j=1}^{M_{k,II}} \sum_{l=1}^{M_{k,II}} \left[1 - f(l, 1) Q \left((2m - 2 + t_k + (2j - 2l)(M_{k,I} - 1 + t_k)) \sqrt{\beta_k q_k} \right) \right. \right. \\
&\quad \left. \left. - f(l, M_{k,II}) Q \left((2M_{k,I} - 2m + t_k + (2l - 2j)(M_{k,I} - 1 + t_k)) \sqrt{\beta_k q_k} \right) \right] \right. \\
&\quad \left. \times \left[f(m, 1) Q \left((1 + 2(j - l)(M_{k,I} - 1 + t_k)) \sqrt{\beta_k q_k} \right) + f(m, M_{k,I}) Q \left((1 - 2(j - l)(M_{k,I} - 1 + t_k)) \sqrt{\beta_k q_k} \right) \right] \right\}^2 \\
&\quad + 2B \log_2 M_{k,II} \left[1 - \frac{2M_{k,II} - 2}{M_{k,I} M_{k,II}} \sum_{m=1}^{M_{k,I}} Q \left(\frac{2M_{k,I} - 2m + t_k}{t_k} \sqrt{\beta_k q_k} \right) \right]^2. \tag{38}
\end{aligned}$$

adopting exhaustive search, all possible channel assignment combinations has be calculated. The computational complexity becomes $\mathcal{O}(K!/(K - N)!)$ for $N \leq K$ and $\mathcal{O}(N!/2^{N-K})$ for $N > K$ due to

$$\begin{aligned}
& \mathcal{O} \left[\binom{N}{2} \binom{N-2}{2} \cdots \binom{2K-N+2}{2} \right. \\
& \quad \left. \binom{2K-N}{1} \binom{2K-N-1}{1} \cdots \binom{2}{1} \binom{1}{1} \right] \\
&= \mathcal{O}(N!/2^{N-K}), \tag{35}
\end{aligned}$$

where $\binom{j}{i} = j! / [(j-i)!i!]$. The computational complexity of exhaustive search becomes prohibitive for a large N .

B. Modulation Selection

In Subsection IV-A, we have proposed low-complexity channel assignment algorithms based on swap matching. Note that each time a swap operation is executed, the modulation schemes should be reoptimized according to the updated channel qualities. In this subsection, we consider an adaptive modulation method based on the SNR thresholds for single-channel users and propose a DNN-based modulation selection scheme for the two-user channels.

1) Modulation Selection for the Single-user Channels

If channel C_k is a single-user channel, i.e., $k \in \mathcal{K}^\dagger$, then the modulation selection problem can be formulated as

$$\mathcal{P}_{8,k} (k \in \mathcal{K}^\dagger) : \underset{M_{k,0}}{\text{maximize}} \quad J_k \tag{36a}$$

$$\text{subject to} \quad M_{k,0} \in \mathcal{M}, \tag{36b}$$

where $M_{k,0}$ is the modulation order of user $U_{k,0}$ on channel C_k and $J_k = J_{k,0}$ is the effective throughput of user $U_{k,0}$ as given in (12). According to (12), $J_{k,0}$ is a function of $M_{k,0}$ and $h_{k,0} p_{k,0}$, where $h_{k,0} p_{k,0}$ is the SNR of user $U_{k,0}$. We consider an adaptive modulation method for conventional OMA systems [41]–[43]. First, we plot the effective throughput curves of different modulation schemes versus the SNR as shown in Fig. 3. From Fig. 3, the SNR thresholds can be obtained as 35, 45, and 52.5 dBm, and a modulation selection table can be obtained as shown in Table I. Thus, to maximize

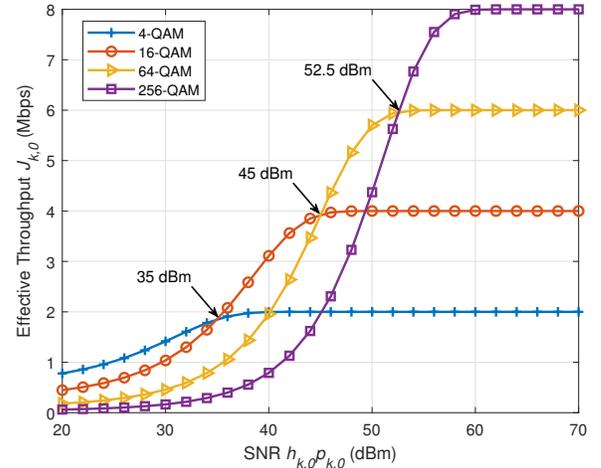


Figure 3: Effective throughput of the single-user channels versus the SNR.

Table I: The adaptive modulation selection of the single-user channels based on the SNR thresholds.

SNR range (dBm)	Modulation scheme
$\text{SNR} \leq 35$	4-QAM
$35 < \text{SNR} \leq 45$	16-QAM
$45 < \text{SNR} \leq 52.5$	64-QAM
$\text{SNR} > 52.5$	256-QAM

the effective throughput in problem $\mathcal{P}_{8,k}$, we can choose the optimal modulation scheme according to the SNR by referring to Table I.

2) Modulation Selection for the Two-User Channels

If channel C_k is a two-user channel, i.e., $k \in \mathcal{K}^\ddagger$, then the modulation selection problem can be formulated as

$$\mathcal{P}_{9,k} (k \in \mathcal{K}^\ddagger) : \underset{\{M_{k,I}, M_{k,II}\}}{\text{maximize}} \quad J_k \tag{37a}$$

$$\text{subject to} \quad M_{k,I}, M_{k,II} \in \mathcal{M}, \tag{37b}$$

where J_k is the effective throughput of channel C_k and $M_{k,i}$ ($i = I, II$) is the modulation order of user $U_{k,i}$ on channel

$$\langle M_{k,I}^{(i)}, M_{k,II}^{(i)} \rangle = \underset{\langle M_{k,I}, M_{k,II} \rangle \in \mathcal{M}^2}{\operatorname{argmax}} J_k \Big|_{p_{k,I}=p_{k,I}^*, p_{k,II}=p_{k,II}^*, h_{k,I}=h_{k,I}^{(i)}, t_k=t_k^{(i)}, q_k=q_k^{(i)}}. \quad (40)$$

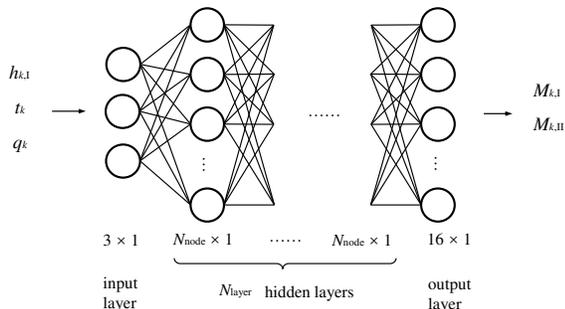


Figure 4: Architecture of a DNN-based modulation selection scheme.

C_k . Substituting $p_{k,I}$ and $p_{k,II}$ by (20) and (21), the effective throughput $J_k|_{p_{k,I}=p_{k,I}^*, p_{k,II}=p_{k,II}^*}$ is given in (38) at the top of the last page, where β_k is defined in (28).

Different from the single-user channels, the modulation schemes of the two users on the same channel are coupled and should be jointly designed. Most of the existing works on NOMA employ adaptive modulation methods based on the single user's SNR thresholds [44]–[46], which are suboptimal for the two-user channels (as shown in Section V). To obtain the optimal modulation scheme which maximizes J_k in (38), one has to exhaustively search all the feasible modulation combinations, i.e.,

$$\begin{aligned} & \langle M_{k,I}, M_{k,II} \rangle \in \mathcal{M}^2 \\ & = \{ \langle 2, 2 \rangle, \langle 2, 4 \rangle, \langle 2, 8 \rangle, \langle 2, 16 \rangle, \\ & \quad \langle 4, 2 \rangle, \langle 4, 4 \rangle, \langle 4, 8 \rangle, \langle 4, 16 \rangle, \\ & \quad \langle 8, 2 \rangle, \langle 8, 4 \rangle, \langle 8, 8 \rangle, \langle 8, 16 \rangle, \\ & \quad \langle 16, 2 \rangle, \langle 16, 4 \rangle, \langle 16, 8 \rangle, \langle 16, 16 \rangle \}. \quad (39) \end{aligned}$$

Although the range of search contains only 16 candidates, each round of search requires approximately $\sum_{\langle M_{k,I}, M_{k,II} \rangle \in \mathcal{M}^2} M_{k,I} + 4M_{k,I}M_{k,II}^2 \approx 4 \times 10^4$ times of $Q(\cdot)$ function calculation including integral operation, which still leads to a considerable complexity.

Note that the modulation selection problem can be regarded as a classification problem, mapping the channel quality and the allocated power into the modulation scheme. Inspired by the idea that DNN is a useful tool which can learn from complex data structures and derive highly nonlinear decision boundaries for solving classification problems [47], [48], we next propose a low-complexity modulation selection scheme via a DNN-based classifier. According to (38), we extract $h_{k,I}$, t_k , and q_k as the features and take a 16-dimensional one-hot vector as the label to indicate which modulation combination in (39) is selected. The overview of the fully-connected DNN architecture is shown in Fig. 4, consisting of an input layer with 3 neurons, N_{layer} hidden layers with N_{node} nodes, and an output layer with 16 neurons.

Algorithm 2 Joint power allocation, channel assignment and modulation selection algorithm

- 1) **Initialization:**
 - a) Initiate the power allocation;
 - b) Set precision δ and the parameter N_{converge} .
 - 2) **Repeat**
 - 3) Obtain the channel assignment according to Algorithm 1, and select the modulation scheme through DNN;
 - 4) Update the power allocation according to (20), (21) and (30);
 - 5) **Until** the change of the effective throughput is less than δ for consecutive N_{converge} iterations.
 - 6) **Return** the optimal power allocation, channel assignment, and modulation scheme, and the corresponding effective throughput.
-

The parameters of the neural network, i.e., the weights and biases, are trained in a supervised manner based on backpropagation and the stochastic gradient descent method. The training dataset is generated by exhaustively searching the optimal modulation schemes that maximize $J_k|_{p_{k,I}=p_{k,I}^*, p_{k,II}=p_{k,II}^*}$ in (38). For example, assume that the i th training example has the input $h_{k,I}^{(i)}$, $t_k^{(i)}$, and $q_k^{(i)}$. Then, the output of the i th training example is obtained from (40) at the top of this page.

For the DNN-based modulation selection scheme, the dataset generation and the DNN training process can be completed offline. The computational complexity during real-time scheduling mainly depends on the size of the network. Calculating the output of the DNN requires $N_{\text{node}}^2(N_{\text{layer}} - 1) + 19N_{\text{node}}$ times of multiplication and addition, and $N_{\text{node}}N_{\text{layer}} + 16$ times of activation operation. Thus, the complexity is $\mathcal{O}(N_{\text{node}}^2N_{\text{layer}})$.

C. Joint Resource Allocation Algorithm

In Section III, we proposed the power allocation scheme for fixed channel assignment and modulation schemes. In Subsections IV-A and IV-B, we developed the channel assignment algorithm and the modulation selection scheme for fixed power allocation. Naturally, a joint resource allocation scheme, as the solution to the original problem \mathcal{P}_1 , can be achieved by iteratively optimizing power allocation, channel assignment, and modulation selection.

Specifically, we first initialize the power allocation, e.g., equally allocate power to each channel. Then, the channel assignment, modulation schemes, and power allocation are iteratively optimized according to Algorithm 1, the DNN scheme, and Theorems 1 and 2, respectively, until convergence. The details are described in Algorithm 2.

The computational complexity depends on the three parts, which are $\mathcal{O}(K)$ for power allocation, $\mathcal{O}(N^2N_{\text{iter}})$ for channel assignment, and $\mathcal{O}(N_{\text{node}}^2N_{\text{layer}})$ for modulation selection, respectively. Note that during each iteration of Algorithm 2,

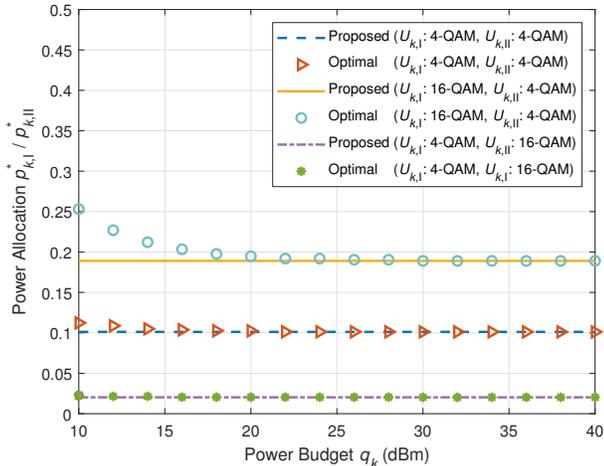


Figure 5: Comparison of the proposed and optimal power allocation schemes within channels.

Table II: Simulation setting of the DNN-based modulation selection scheme.

Parameter	Value
Number of epochs	100
Batch size	1000
Training set size	8×10^5
Validation set size	10^5
Testing set size	10^5

both the power allocation and the channel assignment are executed once. The modulation selection module is invoked each time a swap operation is executed in the channel assignment algorithm. Thus, the computational complexity of Algorithm 2 is $\mathcal{O}(N_{\text{outer}}(K + N^2 N_{\text{iter}} N_{\text{node}}^2 N_{\text{layer}}))$, where N_{outer} is the number of outer iterations.

Algorithm 2 converges quickly after a few iterations with a relatively low computational complexity and it achieves almost the same effective throughput as exhaustive search, which is demonstrated by numerical results. Therefore, Algorithm 2 can efficiently address the effective throughput maximization of NOMA.

V. NUMERICAL RESULTS

In this section, numerical results are presented to evaluate the aforementioned algorithms along with some insightful discussions. We employ the channel model with the fading coefficient being a zero-mean unit-variance complex Gaussian variable and the path loss being $-(128.1 + 37.6 \log_{10} \tau_n)$ dB [49], where τ_n is the distance between the BS and user U_n in kilometers. The noise power spectral density is assumed to be -174 dBm/Hz for all users.

A. Power Allocation and Modulation Selection

In Fig. 5, we compare the proposed power allocation scheme in Theorem 1 with the optimal one obtained via exhaustive search by illustrating $p_{k,I}^*/p_{k,II}^*$ versus the power budget q_k . The distances between the BS and the users are 0.2 km

Table III: Validation accuracy of the DNN-based modulation selection scheme.

N_{layer}	N_{node}	Validation accuracy
1	5	0.9116
1	10	0.9785
1	20	0.9868
2	5	0.9783
2	10	0.9871
2	20	0.9941

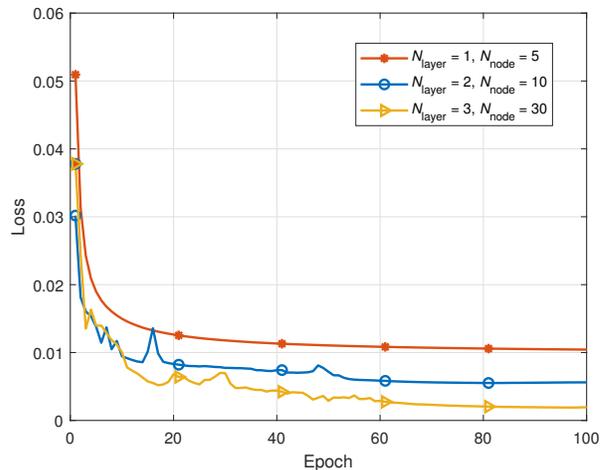


Figure 6: Convergence of the proposed DNN model.

for $U_{k,I}$ and 0.3 km for $U_{k,II}$. Three modulation modes are considered: 1) both $U_{k,I}$ and $U_{k,II}$ employ 4-QAM; 2) $U_{k,I}$ employs 16-QAM and $U_{k,II}$ employs 4-QAM; 3) $U_{k,I}$ employs 4-QAM and $U_{k,II}$ employs 16-QAM. The results show that the proposed scheme well approximates the optimal one even if the power budget is not large.

Next, the performance of the DNN-based modulation selection scheme is evaluated. We preprocess the inputs $h_{k,I}$, t_k , and q_k by $\log_{10}(\cdot)$ and normalize them in the range $[0, 1]$ before feeding them into the neural network such that the differences between the samples can be captured accurately. The activation functions for the hidden layers and the output layer are rectified linear unit (ReLU) and softmax, respectively, so that the output values are scaled to the range $[0, 1]$. The learning rate is initially set as 0.01 and multiplied by a discount factor 0.98 in each epoch. Other detailed parameters are presented in Table II.

In Fig. 6, we illustrate the loss values versus training epochs to show the convergence of the DNN training when adopting different numbers of hidden layers N_{layer} and neurons N_{node} . Results in Table III show that the DNN models with different network sizes achieve an accuracy over 0.9 and have a lower complexity than the exhaustive search, which is approximately 4×10^4 times of $Q(\cdot)$ function calculation according to the analysis in Subsection IV-B. To achieve a good accuracy with a relatively low complexity, we choose $N_{\text{layer}} = 2$ and $N_{\text{node}} = 10$ for the following joint algorithm simulations.

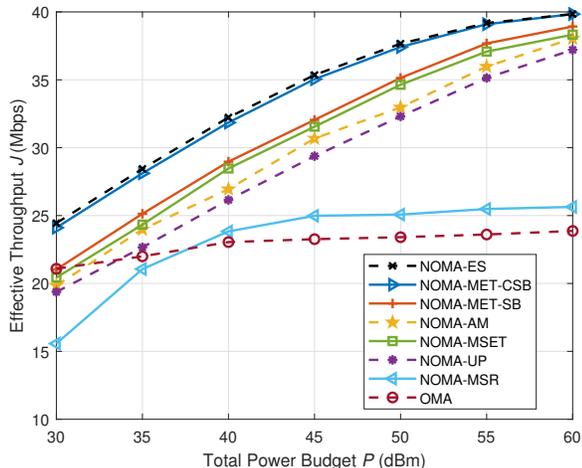


Figure 7: Comparison of the effective throughput versus the total power budget. We set $K = 3$ and $N = 5$.

B. Joint Resource Allocation Algorithm

In this subsection, we evaluate the performance of the proposed joint resource allocation algorithm of NOMA, i.e., Algorithm 2. The total bandwidth is set as K MHz. N users are located uniformly in a cell with the radius of 1 km and the BS is located in the center. Two users within a channel perform the NOMA scheme, while the users on different channels are assumed to be interference-free. The following effective throughput is obtained by averaging 100 randomly-generated channel profiles.

For convenience, the proposed maximum effective throughput NOMA schemes based on the CSB policy and the SB policy are denoted by NOMA-MET-CSB and NOMA-MET-SB, respectively. For comparison, six benchmarks are considered as follows:

- 1) NOMA-ES: the exhaustive search scheme of NOMA.
- 2) NOMA-AM: the modulation orders of all the users, including the users on the two-user channels, are determined separately via the adaptive modulation selection method in Table I.
- 3) NOMA-MSET: the maximum short-packet effective throughput NOMA scheme [23].
- 4) NOMA-UP: the NOMA scheme based on the user pairing method where a near user and a far user are selected to share a channel [50].
- 5) NOMA-MSR: the maximum sum rate NOMA scheme where the residual interference caused by imperfect SIC is modeled as a continuously-valued Gaussian signal [18].
- 6) OMA: the conventional OMA scheme where K MHz bandwidth is orthogonally allocated to N users without superposition [51].

In Fig. 7, we compare the effective throughput obtained via different algorithms versus the total power budget with 3 channels and 5 users. One can find that the proposed NOMA-MET-CSB and NOMA-MET-SB schemes outperform the NOMA-AM, NOMA-MSET, NOMA-UP, NOMA-MSR, and OMA schemes, achieving almost the same performance

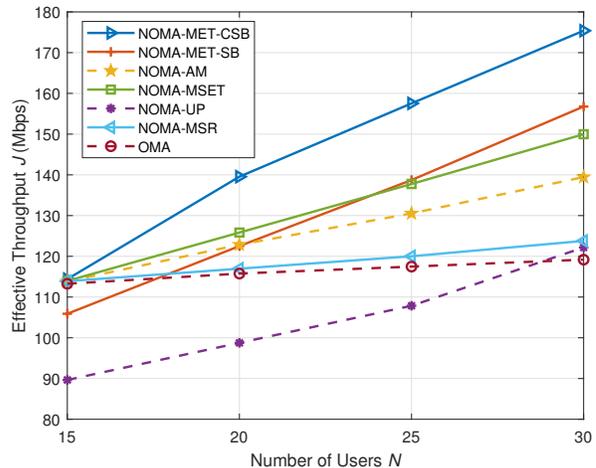


Figure 8: Comparison of the effective throughput obtained via different algorithms versus the number of users. We set $K = 15$ and $P = 40$ dBm.

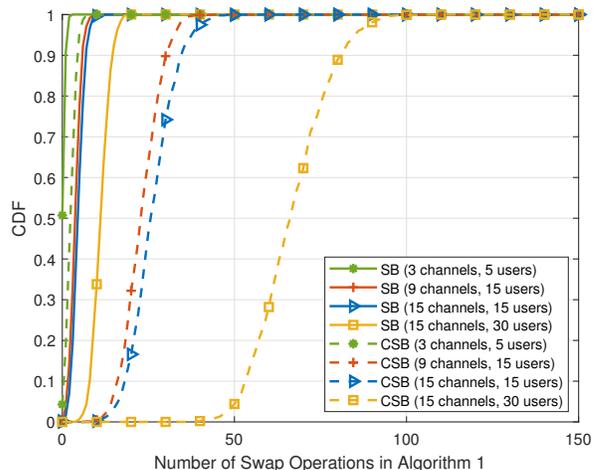


Figure 9: CDF of the number of swap operations in Algorithm 1.

as NOMA-ES. The NOMA-MSR scheme obtains a relatively lower effective throughput than other NOMA schemes when the total power budget is small.

Furthermore, Fig. 7 shows that as the total power budget increases, the effective throughput has an increasing trend due to lower error probabilities and higher modulation orders. Eventually, it reaches an upper bound. The upper bound of the effective throughput corresponds to the situation where all the users are error-free and the highest 256-QAM modulation is employed, which for NOMA is $\sum_{n=1}^N \log_2 256 = 8N = 40$ Mbps and for OMA is $\frac{K}{N} \sum_{n=1}^N \log_2 256 = 8K = 24$ Mbps, as shown in Fig. 7.

Fig. 8 illustrates the effective throughput versus the number of users when the number of channels K is equal to 15 and the total power budget is 40 dBm [52]. One can observe that the effective throughput of both NOMA and OMA schemes increases monotonically with respect to the number of users.

$$\phi_{k,I} < \frac{2(M_{k,I} - 1)}{M_{k,I}} Q \left(\sqrt{\frac{3h_{k,I}p_{k,I}}{M_{k,I}^2 - 1}} \right) + \frac{3(M_{k,I} - 1)(M_{k,II} - 1)}{M_{k,I}} Q \left((1 - M_{k,I}) \sqrt{\frac{3h_{k,II}p_{k,I}}{M_{k,I}^2 - 1}} + \sqrt{\frac{3h_{k,II}(q_k - p_{k,I})}{M_{k,II}^2 - 1}} \right). \quad (42)$$

$$\phi_{k,II} < \frac{2(M_{k,II} - 1)}{M_{k,II}} Q \left((1 - M_{k,I}) \sqrt{\frac{3h_{k,II}p_{k,I}}{M_{k,I}^2 - 1}} + \sqrt{\frac{3h_{k,II}(q_k - p_{k,I})}{M_{k,II}^2 - 1}} \right). \quad (43)$$

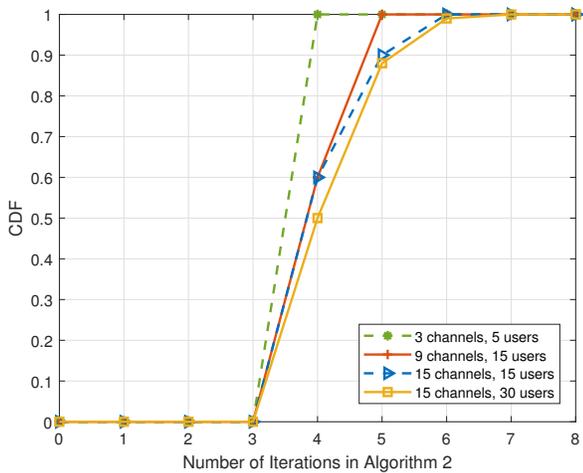


Figure 10: CDF of the number of iterations in Algorithm 2.

When the number of users is equal to the number of channels, i.e., $N = K = 15$, the NOMA system degrades to an OMA system and NOMA-MET-CSB achieves similar performance to the NOMA-AM, NOMA-MSET, NOMA-MSR, and OMA schemes. As N increases, the superiority of the NOMA-MET-CSB and NOMA-MET-SB schemes over OMA increases due to its overloading. The superiority of the NOMA-MET-CSB and NOMA-MET-SB schemes over the NOMA-AM, NOMA-MSET, and NOMA-MSR schemes increases due to its near-optimal power allocation scheme and DNN-based modulation selection scheme. Furthermore, NOMA-MET-CSB always outperforms NOMA-MET-SB and NOMA-UP, which shows the effectiveness of the channel assignment algorithm based on the CSB policy.

The cumulative distribution functions (CDFs) of the number of swap operations in Algorithm 1 and the number of iterations in Algorithm 2 are illustrated in Figs. 9 and 10, respectively. From Fig. 9, we can observe that the channel assignment algorithm based on the CSB policy requires more swap operations than that based on the SB policy. Furthermore, as the numbers of channels and users increase, more swap operations and iterations are required to achieve convergence. Nevertheless, even in a system with 15 channels and 30 users, the proposed NOMA algorithm converges quickly within 100 swap operations and 10 iterations, while the exhaustive NOMA algorithm has to search $30!/2^{15} \approx 8 \times 10^{27}$ channel assignment combinations, leading to a prohibitive complexity.

VI. CONCLUSION

In this paper, we investigated the effective throughput maximization of downlink multi-user multi-channel NOMA systems, considering both the error performance and the data rate along with imperfect SIC and practical QAMs. The exact and approximate expressions of the effective throughput were derived as functions of the transmit power, channel gain, and modulation order. We formulated a joint power allocation, channel assignment, and modulation selection problem to maximize the effective throughput. To address this problem, we developed an analytical power allocation scheme, including the closed-form power allocation within channels and the waterfilling-form power budget allocation among channels, and proposed efficient channel assignment and modulation selection methods based on matching theory and machine learning, respectively. Consequently, a joint resource allocation algorithm was proposed to maximize the effective throughput. We conducted comprehensive simulations to evaluate the performance of the proposed NOMA scheme, which is shown to outperform the existing OMA and other NOMA schemes.

APPENDIX A PROOF OF PROPOSITION 3

Due to $\phi_{k,I}, \phi_{k,II}, \epsilon_{k,I}, \epsilon_{k,II} \in [0, 1]$, we have $\epsilon_{k,I} \leq 2\phi_{k,I}$ and $\epsilon_{k,II} \leq 2\phi_{k,II}$ according to (6). Thus, J_k is lower bounded by

$$J_k \geq 2B(1 - 2\phi_{k,I}) \log_2 M_{k,I} + 2B(1 - 2\phi_{k,II}) \log_2 M_{k,II}. \quad (41)$$

According to [27], the upper bounds of $\phi_{k,I}$ and $\phi_{k,II}$ are given by (42) and (43) shown at the top of this page, respectively. Substituting (42) and (43) into (41), Proposition 3 can be obtained.

APPENDIX B PROOF OF THEOREM 1

The derivative of J_k^L with respect to $p_{k,I}$ exists when $q_k \neq 0$ and $0 < p_{k,I} < q_k$. The derivative is given by (44) on the next page. Setting $\partial J_k^L / \partial p_{k,I}$ in (44) to be equal to zero, we obtain (45), where $t_k = \sqrt{h_{k,I}/h_{k,II}}$ is the channel gain ratio of channel C_k . In (45), the term $p_{k,I}/q_k$ is bounded by $0 < p_{k,I}/q_k < 1$. We further define $\theta_k \triangleq p_{k,I}/q_k \in (0, 1)$, and (45) can be rewritten as (46) on the next page. In (46), the left-hand side and the $\ln(\cdot)$ term on the right-hand side are

$$\begin{aligned} \frac{\partial J_k^L}{\partial p_{k,I}} = & - \left(\frac{4B(M_{k,II} - 1) \log_2 M_{k,II}}{M_{k,II} q_k \sqrt{2\pi}} + \frac{6B(M_{k,I} - 1)(M_{k,II} - 1) \log_2 M_{k,I}}{M_{k,I} q_k \sqrt{2\pi}} \right) \times \\ & \exp \left[\frac{- \left((1 - M_{k,I}) \sqrt{\frac{3h_{k,II} p_{k,I}}{M_{k,I}^2 - 1}} + \sqrt{\frac{3h_{k,II}(q_k - p_{k,I})}{M_{k,II}^2 - 1}} \right)^2}{2} \right] \times \left(\sqrt{\frac{3h_{k,II} q_k^2}{(M_{k,II}^2 - 1)(q_k - p_{k,I})}} + \right. \\ & \left. \sqrt{\frac{3(M_{k,I} - 1)^2 h_{k,II} q_k^2}{(M_{k,I}^2 - 1)p_{k,I}}} + \frac{4B(M_{k,I} - 1) \log_2 M_{k,I}}{M_{k,I} q_k \sqrt{2\pi}} \sqrt{\frac{3h_{k,I} q_k^2}{(M_{k,I}^2 - 1)p_{k,I}}} \exp \left[\frac{-3h_{k,I} p_{k,I}}{2(M_{k,I}^2 - 1)} \right] \right). \end{aligned} \quad (44)$$

$$\begin{aligned} & \left((1 - M_{k,I}) \sqrt{\frac{3h_{k,II} p_{k,I}/q_k}{M_{k,I}^2 - 1}} + \sqrt{\frac{3h_{k,II}(1 - p_{k,I}/q_k)}{M_{k,II}^2 - 1}} \right)^2 - \frac{3h_{k,I} p_{k,I}/q_k}{M_{k,I}^2 - 1} \\ = & \frac{2}{q_k} \ln \left[\left(\frac{3(M_{k,II} - 1)}{2t_k} + \frac{M_{k,I}(M_{k,II} - 1) \log_2 M_{k,II}}{t_k M_{k,II}(M_{k,I} - 1) \log_2 M_{k,I}} \right) \left(M_{k,I} - 1 + \sqrt{\frac{M_{k,I}^2 - 1}{M_{k,II}^2 - 1}} \sqrt{\frac{p_{k,I}/q_k}{1 - p_{k,I}/q_k}} \right) \right]. \end{aligned} \quad (45)$$

$$\begin{aligned} & \left((1 - M_{k,I}) \sqrt{\frac{3h_{k,II} \theta_k}{M_{k,I}^2 - 1}} + \sqrt{\frac{3h_{k,II}(1 - \theta_k)}{M_{k,II}^2 - 1}} \right)^2 - \frac{3h_{k,I} \theta_k}{M_{k,I}^2 - 1} \\ = & \frac{2}{q_k} \ln \left[\left(\frac{3(M_{k,II} - 1)}{2t_k} + \frac{M_{k,I}(M_{k,II} - 1) \log_2 M_{k,II}}{t_k M_{k,II}(M_{k,I} - 1) \log_2 M_{k,I}} \right) \left(M_{k,I} - 1 + \sqrt{\frac{M_{k,I}^2 - 1}{M_{k,II}^2 - 1}} \sqrt{\frac{\theta_k}{1 - \theta_k}} \right) \right]. \end{aligned} \quad (46)$$

$$\begin{aligned} J_k^L|_{p_{k,I}=q_k} = & 2B \log_2 M_{k,I} + 2B \log_2 M_{k,II} - \frac{8B(M_{k,I} - 1) \log_2 M_{k,I}}{M_{k,I}} Q \left(\sqrt{\frac{3h_{k,I} q_k}{M_{k,I}^2 - 1}} \right) - \\ & \left(\frac{8B(M_{k,II} - 1) \log_2 M_{k,II}}{M_{k,II}} + \frac{12B(M_{k,I} - 1)(M_{k,II} - 1) \log_2 M_{k,I}}{M_{k,I}} \right) Q \left((1 - M_{k,I}) \sqrt{\frac{3h_{k,II} q_k}{M_{k,I}^2 - 1}} \right) \\ < & 2B \log_2 M_{k,I} + 2B \log_2 M_{k,II} - \frac{4B(M_{k,II} - 1) \log_2 M_{k,II}}{M_{k,II}} + \frac{6B(M_{k,I} - 1)(M_{k,II} - 1) \log_2 M_{k,I}}{M_{k,I}}. \end{aligned} \quad (51)$$

both finite. When q_k is sufficiently large, the right-hand side of (46) tends to zero. Thus, (46) can be approximated by

$$\begin{aligned} & \left((1 - M_{k,I}) \sqrt{\frac{3h_{k,II} \theta_k}{M_{k,I}^2 - 1}} + \sqrt{\frac{3h_{k,II}(1 - \theta_k)}{M_{k,II}^2 - 1}} \right)^2 \\ & - \frac{3h_{k,I} \theta_k}{M_{k,I}^2 - 1} = 0. \end{aligned} \quad (47)$$

To determine the (negative or positive) sign of $(1 - M_{k,I}) \sqrt{\frac{3h_{k,II} \theta_k}{M_{k,I}^2 - 1}} + \sqrt{\frac{3h_{k,II}(1 - \theta_k)}{M_{k,II}^2 - 1}}$, we define

$$\tilde{p}_{k,I} \triangleq \frac{(M_{k,I} + 1)q_k}{(M_{k,I} - 1)(M_{k,II}^2 - 1) + M_{k,I} + 1} \in (0, q_k). \quad (48)$$

When $p_{k,I} \leq \tilde{p}_{k,I}$, we have $(1 - M_{k,I}) \sqrt{\frac{3h_{k,II} \theta_k}{M_{k,I}^2 - 1}} + \sqrt{\frac{3h_{k,II}(1 - \theta_k)}{M_{k,II}^2 - 1}} \geq 0$ and equation (47) has a root

$$p_{k,I}^{(1)} = \frac{(M_{k,I}^2 - 1)q_k}{(M_{k,I} - 1 + t_k)^2(M_{k,II}^2 - 1) + M_{k,I}^2 - 1} \in (0, \tilde{p}_{k,I}). \quad (49)$$

When $p_{k,I} > \tilde{p}_{k,I}$, we have $(1 - M_{k,I}) \sqrt{\frac{3h_{k,II} \theta_k}{M_{k,I}^2 - 1}} + \sqrt{\frac{3h_{k,II}(1 - \theta_k)}{M_{k,II}^2 - 1}} < 0$ and the number of solutions to (47) depends on the relationship between $M_{k,I} - 1$ and t_k . If $M_{k,I} - 1 \leq t_k$, then (47) has no root; otherwise, (47) has a root

$$p_{k,I}^{(2)} = \frac{(M_{k,I}^2 - 1)q_k}{(M_{k,I} - 1 - t_k)^2(M_{k,II}^2 - 1) + M_{k,I}^2 - 1} \in (\tilde{p}_{k,I}, q_k). \quad (50)$$

Therefore, if $M_{k,I} - 1 \leq t_k$, then J_k^L is monotonically increasing when $0 < p_{k,I} \leq p_{k,I}^{(1)}$ and is decreasing when $p_{k,I}^{(1)} < p_{k,I} < q_k$, with the peak point $p_{k,I}^{(1)}$. If $M_{k,I} - 1 > t_k$, then J_k^L is monotonically increasing with respect to $p_{k,I}$ when $0 < p_{k,I} \leq p_{k,I}^{(1)}$, is decreasing when $p_{k,I}^{(1)} < p_{k,I} \leq p_{k,I}^{(2)}$, and is increasing again when $p_{k,I}^{(2)} < p_{k,I} < q_k$.

To pinpoint the optimum point, we next compare $J_k^L|_{p_{k,I}=p_{k,I}^{(1)}}$ with $J_k^L|_{p_{k,I}=q_k}$. According to (19), we have (51) at the top of this page. On the other hand, according to (52) shown on the next page, $J_k^L|_{p_{k,I}=p_{k,I}^{(1)}}$ approaches $2B \log_2 M_{k,I} + 2B \log_2 M_{k,II}$ as q_k increases. Hence, for

$$J_k^L|_{p_{k,I}=p_{k,I}^{(1)}} = 2B \log_2 M_{k,I} + 2B \log_2 M_{k,II} - Q \left(\sqrt{\frac{3h_{k,I}q_k}{(M_{k,I} - 1 + t_k)^2(M_{k,II}^2 - 1) + M_{k,I}^2 - 1}} \right) \\ \times \left(\frac{4B(M_{k,I} - 1)(3M_{k,II} - 1) \log_2 M_{k,I}}{M_{k,I}} + \frac{8B(M_{k,II} - 1) \log_2 M_{k,II}}{M_{k,II}} \right). \quad (52)$$

a sufficiently large power budget, we have $J_k^L|_{p_{k,I}=p_{k,I}^{(1)}} > J_k^L|_{p_{k,I}=q_k}$, and thus, $p_{k,I}^{(1)}$ is the maximum point of the lower bound of J_k . Denoting $p_{k,I}^{(1)}$ by $p_{k,I}^*$, Theorem 1 is proven.

APPENDIX C

PROOF OF LEMMAS 1 AND 2

When $q_k > 0$, the partial derivative of \tilde{J}_k with respect to q_k is

$$\frac{\partial \tilde{J}_k}{\partial q_k} = \frac{\alpha_k \sqrt{\beta_k}}{2\sqrt{2\pi}q_k} \exp\left(-\frac{\beta_k q_k}{2}\right) > 0, \quad k = 1, \dots, K, \quad (53)$$

which shows that \tilde{J}_k is monotonically increasing with respect to q_k and Lemma 1 is proven. When $q_k > 0$, the second-order partial derivative of \tilde{J}_k is

$$\frac{\partial^2 \tilde{J}_k}{\partial q_k^2} = -\frac{\alpha_k \sqrt{\beta_k}}{4\sqrt{2\pi}q_k} \left(\beta_k + \frac{1}{q_k} \right) \exp\left(-\frac{\beta_k q_k}{2}\right) < 0, \quad k = 1, \dots, K, \quad (54)$$

which shows that \tilde{J}_k is a concave function of q_k and Lemma 2 is proven.

APPENDIX D

PROOF OF THEOREM 2

The optimization problem \mathcal{P}_6 is equivalent to the following problem:

$$\mathcal{P}_9 : \text{minimize}_{\{q_k \geq 0\}} \sum_{k=1}^K \alpha_k Q(\sqrt{\beta_k q_k}) \quad (55a)$$

$$\text{subject to} \quad \sum_{k=1}^K q_k \leq P. \quad (55b)$$

The Lagrangian is given by

$$L(q_1, \dots, q_K, \mu) = \sum_{k=1}^K \alpha_k Q(\sqrt{\beta_k q_k}) + \mu \left(\sum_{k=1}^K q_k - P \right), \quad (56)$$

where μ is the Lagrange multiplier. Since \mathcal{P}_9 is a convex optimization problem, its optimal solution is characterized by satisfying Karush-Kuhn-Tucker (KKT) conditions:

$$\frac{\partial L}{\partial q_k} = \mu - \frac{\alpha_k \sqrt{\beta_k}}{2\sqrt{2\pi}q_k} \exp\left(-\frac{\beta_k q_k}{2}\right) = 0, \quad (57)$$

$$\mu \left(\sum_{k=1}^K q_k - P \right) = 0, \quad (58)$$

$$\sum_{k=1}^K q_k - P \leq 0, \quad (59)$$

$$\mu \geq 0. \quad (60)$$

From (57), it is obvious that $\mu \neq 0$. Hence, due to (58) and (60), $\mu > 0$ and $\sum_{k=1}^K q_k - P = 0$ hold. To obtain the root of (57), we have

$$(57) \iff \frac{2\mu\sqrt{2\pi}q_k}{\alpha_k\sqrt{\beta_k}} \exp\left(\frac{\beta_k q_k}{2}\right) = 1 \\ \iff \frac{8\pi\mu^2 q_k}{\alpha_k^2 \beta_k} \exp(\beta_k q_k) = 1 \iff \frac{\alpha_k^2 \beta_k^2}{8\pi\mu^2} = \beta_k q_k \exp(\beta_k q_k) \\ \iff W_0\left(\frac{\alpha_k^2 \beta_k^2}{8\pi\mu^2}\right) = \beta_k q_k \iff q_k^* = \frac{W_0\left(\frac{\alpha_k^2 \beta_k^2}{8\pi\mu^2}\right)}{\beta_k} > 0, \quad (61)$$

where condition 1 \iff condition 2 means that conditions 1 and 2 are sufficient and necessary conditions for each other. For convenience, define $\lambda \triangleq \frac{\mu^2}{8\pi}$, and thus, Theorem 2 is proven. λ is chosen such that $\sum_{k=1}^K q_k^* = P$, i.e.,

$$\sum_{k=1}^K \frac{W_0\left(\frac{\alpha_k^2 \beta_k^2}{8\pi\lambda}\right)}{\beta_k} = P. \quad (62)$$

The left-hand side is a monotonically decreasing function of λ and has the range $(0, +\infty)$. Thus, equation (62) has a unique solution.

APPENDIX E

PROOF OF COROLLARY 1

To prove the upper bound of λ , we first assume oppositely that

$$\lambda > \max_{k=1, \dots, K} \left\{ \frac{\alpha_k^2 \beta_k}{8\pi P/K} \exp(-\beta_k P/K) \right\}, \quad (63)$$

and (63) leads to

$$\lambda = \frac{\alpha_k^2 \beta_k}{8\pi q_k} \exp(-\beta_k q_k) \\ > \max_{k=1, \dots, K} \left\{ \frac{\alpha_k^2 \beta_k}{8\pi P/K} \exp(-\beta_k P/K) \right\} \\ \geq \frac{\alpha_k^2 \beta_k}{8\pi P/K} \exp(-\beta_k P/K), \quad k = 1, \dots, K. \quad (64)$$

It is clear that $\frac{\alpha_k^2 \beta_k}{8\pi q_k} \exp(-\beta_k q_k)$ is a monotonically decreasing function of q_k , so it follows from (64) that $q_k < \frac{P}{K}$, $k = 1, \dots, K$, and $\sum_{k=1}^K q_k < P$, which is contradictory to $\sum_{k=1}^K q_k = P$. Therefore, the assumption (63) is overturned and the upper bound in (31) is proven.

We next prove the lower bound of λ . Since $\sum_{k=1}^K q_k = P$ and $q_k \geq 0$, we have $q_k \leq P$. In addition,

$$\lambda = \frac{\alpha_k^2 \beta_k}{8\pi q_k} \exp(-\beta_k q_k) \geq \frac{\alpha_k^2 \beta_k}{8\pi P} \exp(-\beta_k P), \quad k = 1, \dots, K, \quad (65)$$

and thus, the lower bound in (32) is obtained and Corollary 1 is proven.

APPENDIX F
PROOF OF PROPOSITION 5

When the channel gain of the user (or user pair) on channel C_k is extremely strong, i.e., $h_{k,0} \rightarrow \infty$ (or $h_{k,1} \rightarrow \infty$), we obtain $\beta_k \rightarrow \infty$. We next prove $\lim_{\beta_k \rightarrow \infty} q_k^* \rightarrow 0$. According to Theorem 2, we have

$$\lim_{\beta_k \rightarrow \infty} q_k^* = \lim_{\beta_k \rightarrow \infty} \frac{W_0\left(\frac{\alpha_k^2 \beta_k^2}{8\pi\lambda}\right)}{\beta_k} = \lim_{\beta_k \rightarrow \infty} \frac{\partial W_0\left(\frac{\alpha_k^2 \beta_k^2}{8\pi\lambda}\right)}{\partial \beta_k}, \quad (66)$$

where the latter equation is due to L'Hôpital's rule. According to the definition of the function $W_0(\cdot)$, we have

$$\begin{aligned} & \frac{\partial W_0\left(\frac{\alpha_k^2 \beta_k^2}{8\pi\lambda}\right)}{\partial \beta_k} \\ &= \frac{\partial \left[\frac{\alpha_k^2 \beta_k^2}{8\pi\lambda} \exp\left(-W_0\left(\frac{\alpha_k^2 \beta_k^2}{8\pi\lambda}\right)\right) \right]}{\partial \beta_k} \\ &= \frac{2\alpha_k^2 \beta_k}{8\pi\lambda} \exp\left(-W_0\left(\frac{\alpha_k^2 \beta_k^2}{8\pi\lambda}\right)\right) \\ & \quad - \frac{\alpha_k^2 \beta_k}{8\pi\lambda} \exp\left(-W_0\left(\frac{\alpha_k^2 \beta_k^2}{8\pi\lambda}\right)\right) \frac{\partial W_0\left(\frac{\alpha_k^2 \beta_k^2}{8\pi\lambda}\right)}{\partial \beta_k}. \quad (67) \end{aligned}$$

By re-arranging the terms, we obtain

$$\frac{\partial W_0\left(\frac{\alpha_k^2 \beta_k^2}{8\pi\lambda}\right)}{\partial \beta_k} = \frac{\frac{2\alpha_k^2 \beta_k}{8\pi\lambda} \exp\left(-W_0\left(\frac{\alpha_k^2 \beta_k^2}{8\pi\lambda}\right)\right)}{1 + \frac{\alpha_k^2 \beta_k}{8\pi\lambda} \exp\left(-W_0\left(\frac{\alpha_k^2 \beta_k^2}{8\pi\lambda}\right)\right)} < \frac{2}{\beta_k}, \quad (68)$$

and

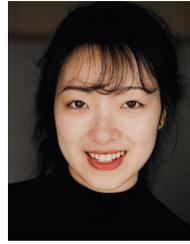
$$0 \leq \lim_{\beta_k \rightarrow \infty} q_k^* = \lim_{\beta_k \rightarrow \infty} \frac{\partial W_0\left(\frac{\alpha_k^2 \beta_k^2}{8\pi\lambda}\right)}{\partial \beta_k} < \lim_{\beta_k \rightarrow \infty} \frac{2}{\beta_k} = 0. \quad (69)$$

Thus, $\lim_{\beta_k \rightarrow \infty} q_k^* \rightarrow 0$ is proved.

REFERENCES

- [1] O. Maraqa, A. S. Rajasekaran, S. Al-Ahmadi, H. Yanikomeroğlu, and S. M. Sait, "A survey of rate-optimal power domain NOMA with enabling technologies of future wireless networks," *IEEE Commun. Surveys & Tuts.*, vol. 22, no. 4, pp. 2192–2235, 4th Quart. 2020.
- [2] Y. Liu, Z. Qin, M. El-kashlan, Z. Ding, A. Nallanathan, and L. Hanzo, "Nonorthogonal multiple access for 5G and beyond," *Proceedings of the IEEE*, vol. 105, no. 12, pp. 2347–2381, Dec. 2017.
- [3] V. W. S. Wong, R. Schober, D. W. K. Ng, and L.-C. Wang, *Key Technologies for 5G Wireless Systems*. Cambridge University Press, 2017.
- [4] M. Shirvanimoghaddam, M. Condoluci, M. Dohler, and S. J. Johnson, "On the fundamental limits of random non-orthogonal multiple access in cellular massive IoT," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 10, pp. 2238–2252, Oct. 2017.
- [5] 3GPP TR 38.812, "Study on non-orthogonal multiple access (NOMA) for NR (Release 16)," Dec. 2018.
- [6] L. Dai, B. Wang, Y. Yuan, S. Han, C. L. I, and Z. Wang, "Non-orthogonal multiple access for 5G: Solutions, challenges, opportunities, and future research trends," *IEEE Commun. Mag.*, vol. 53, no. 9, pp. 74–81, Sep. 2015.
- [7] Z. Ding, Y. Liu, J. Choi, Q. Sun, M. El-kashlan, C.-L. I, and H. V. Poor, "Application of non-orthogonal multiple access in LTE and 5G networks," *IEEE Commun. Mag.*, vol. 55, no. 2, pp. 185–191, Feb. 2017.
- [8] J. Zhu, J. Wang, Y. Huang, S. He, X. You, and L. Yang, "On optimal power allocation for downlink non-orthogonal multiple access systems," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 12, pp. 2744–2757, Dec. 2017.
- [9] Z. Wei, D. W. K. Ng, J. Yuan, and H. M. Wang, "Optimal resource allocation for power-efficient MC-NOMA with imperfect channel state information," *IEEE Trans. Commun.*, vol. 65, no. 9, pp. 3944–3961, Sep. 2017.
- [10] L. Lei, D. Yuan, C. K. Ho, and S. Sun, "Joint optimization of power and channel allocation with non-orthogonal multiple access for 5G cellular systems," in *Proc. IEEE Global Commun. Conf.*, Austin, TX, Dec. 2014, pp. 1–6.
- [11] F. Fang, H. Zhang, J. Cheng, and V. C. M. Leung, "Energy-efficient resource allocation for downlink non-orthogonal multiple access network," *IEEE Trans. Commun.*, vol. 64, no. 9, pp. 3722–3732, Sep. 2016.
- [12] C. He, Y. Hu, Y. Chen, and B. Zeng, "Joint power allocation and channel assignment for NOMA with deep reinforcement learning," *IEEE J. Sel. Areas Commun.*, vol. 37, no. 10, pp. 2200–2210, Oct. 2019.
- [13] S. Wang, T. Lv, and X. Zhang, "Multi-agent reinforcement learning-based user pairing in multi-carrier NOMA systems," in *Proc. IEEE ICC Workshops*, Shanghai, China, May 2019, pp. 1–6.
- [14] H. Huang, Y. Yang, Z. Ding, H. Wang, H. Sari, and F. Adachi, "Deep learning-based sum data rate and energy efficiency optimization for MIMO-NOMA systems," *IEEE Trans. Wireless Commun.*, vol. 19, no. 8, pp. 5373–5388, Aug. 2020.
- [15] K. Wang, Y. Zhou, Z. Liu, Z. Shao, X. Luo, and Y. Yang, "Online task scheduling and resource allocation for intelligent NOMA-based industrial Internet of things," *IEEE J. Sel. Areas Commun.*, vol. 38, no. 5, pp. 803–815, May 2020.
- [16] S. Khairy, P. Balaprakashy, L. X. Cai, and Y. Cheng, "Constrained deep reinforcement learning for energy sustainable multi-UAV based random access IoT networks with NOMA," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 4, pp. 1101–1115, Apr. 2021.
- [17] J. Wang, Q. Peng, Y. Huang, H. Wang, and X. You, "Convexity of weighted sum rate maximization in NOMA systems," *IEEE Signal Process. Lett.*, vol. 24, no. 9, pp. 1323–1327, Sep. 2017.
- [18] X. Chen, R. Jia, and D. W. K. Ng, "On the design of massive non-orthogonal multiple access with imperfect successive interference cancellation," *IEEE Trans. Commun.*, vol. 67, no. 3, pp. 2539–2551, Mar. 2019.
- [19] Q. Qi and X. Chen, "Wireless powered massive access for cellular Internet of Things with imperfect SIC and nonlinear EH," *IEEE Internet of Things J.*, vol. 6, no. 2, pp. 3110–3120, Apr. 2019.
- [20] Y. Sun, D. W. K. Ng, Z. Ding, and R. Schober, "Optimal joint power and subcarrier allocation for full-duplex multicarrier non-orthogonal multiple access systems," *IEEE Trans. Commun.*, vol. 65, no. 3, pp. 1077–1091, Mar. 2017.
- [21] J. Chen, L. Zhang, Y.-C. Liang, and S. Ma, "Optimal resource allocation for multicarrier NOMA in short packet communications," *IEEE Trans. Veh. Technol.*, vol. 69, no. 2, pp. 2141–2156, Feb. 2020.
- [22] S. Han, X. Xu, Z. Liu, P. Xiao, K. Moessner, X. Tao, and P. Zhang, "Energy-efficient short packet communications for uplink NOMA-based massive MTC networks," *IEEE Trans. Veh. Technol.*, vol. 68, no. 12, pp. 12 066–12 078, Dec. 2019.
- [23] X. Sun, S. Yan, N. Yang, Z. Ding, C. Shen, and Z. Zhong, "Short-packet downlink transmission with non-orthogonal multiple access," *IEEE Trans. Wireless Commun.*, vol. 17, no. 7, pp. 4550–4564, Jul. 2018.
- [24] S. Schiessl, M. Skoglund, and J. Gross, "NOMA in the uplink: Delay analysis with imperfect CSI and finite-length coding," *IEEE Trans. Wireless Commun.*, vol. 19, no. 6, pp. 3879–3893, Jun. 2020.
- [25] 3GPP TS 38.211, V16.7.0, "Physical channels and modulation (Release 16)," Sep. 2021.
- [26] Y. Wang, J. Wang, L. Ma, Y. Huang, and C. Zhao, "Minimum error performance of downlink non-orthogonal multiple access systems," in *Proc. IEEE Vehic. Technol. Conf. (VTC-Fall)*, Honolulu, HI, Sep. 2019, pp. 1–5.
- [27] Y. Wang, J. Wang, D. W. K. Ng, R. Schober, and X. Gao, "A minimum error probability NOMA design," *IEEE Trans. Wireless Commun.*, vol. 20, no. 7, pp. 4221–4237, Jul. 2021.
- [28] Z. Ding, F. Adachi, and H. V. Poor, "The application of MIMO to non-orthogonal multiple access," *IEEE Trans. Wireless Commun.*, vol. 15, no. 1, pp. 537–552, Jan. 2016.
- [29] Z. Ding, P. Fan, and H. V. Poor, "Impact of user pairing on 5G nonorthogonal multiple-access downlink transmissions," *IEEE Trans. Veh. Technol.*, vol. 65, no. 8, pp. 6010–6023, Aug. 2016.
- [30] A. Goldsmith, *Wireless Communications*. Cambridge University Press, 2005.
- [31] Y. Polyanskiy, H. V. Poor, and S. Verdú, "Channel coding rate in the finite blocklength regime," *IEEE Trans. Inform. Theory*, vol. 56, no. 5, pp. 2307–2359, May 2010.

- [32] K. Wang, C. Pan, H. Ren, W. Xu, and A. Nallanathan, "Packet error probability and effective throughput for ultra-reliable and low-latency UAV communications," *IEEE Trans. Commun.*, vol. 69, no. 1, pp. 73–84, Jan. 2021.
- [33] L. Lei, D. Yuan, C. K. Ho, and S. Sun, "Power and channel allocation for non-orthogonal multiple access in 5G systems: Tractability and computation," *IEEE Trans. Wireless Commun.*, vol. 15, no. 12, pp. 9580–9594, Dec. 2016.
- [34] B. Di, L. Song, and Y. Li, "Sub-channel assignment, power allocation and user scheduling for non-orthogonal multiple access networks," *IEEE Trans. Wireless Commun.*, vol. 15, no. 11, pp. 7686–7698, Nov. 2016.
- [35] J. G. Proakis and M. Salehi, *Digital Communications*. 5th Edition, McGraw-Hill Education, 2007.
- [36] R. M. Corless, G. H. Gonnet, D. E. G. Hare, D. J. Jeffrey, and D. E. Knuth, "On the Lambert W function," *Adv. Comput. Math.*, vol. 5, pp. 329–359, 1996.
- [37] S. Khakurel, C. Leung, and T. Le-Ngoc, "A generalized water-filling algorithm with linear complexity and finite convergence time," *IEEE Wireless Commun. Lett.*, vol. 3, no. 2, pp. 225–228, Apr. 2014.
- [38] E. Bodine-Baron, C. Lee, A. Chong, B. Hassibi, and A. Wierman, "Peer effects and stability in matching markets," *Lecture Notes in Computer Science*, vol. 6982, pp. 117–129, 2011.
- [39] J. Zhao, Y. Liu, K. K. Chai, A. Nallanathan, Y. Chen, and Z. Han, "Spectrum allocation and power control for non-orthogonal multiple access in HetNets," *IEEE Trans. Wireless Commun.*, vol. 16, no. 9, pp. 5825–5837, Sep. 2017.
- [40] J. Cui, Y. Liu, Z. Ding, P. Fan, and A. Nallanathan, "Optimal user scheduling and power allocation for millimeter wave NOMA systems," *IEEE Trans. Wireless Commun.*, vol. 17, no. 3, pp. 1502–1517, Mar. 2018.
- [41] A. J. Goldsmith and S.-G. Chua, "Variable-rate variable-power MQAM for fading channels," *IEEE Trans. Commun.*, vol. 45, no. 10, pp. 1218–1230, Oct. 1997.
- [42] X. Qiu and K. Chawla, "On the performance of adaptive modulation in cellular systems," *IEEE Trans. Commun.*, vol. 47, no. 6, pp. 884–895, Jun. 1999.
- [43] S. Nagaraj, "Symbol-level adaptive modulation for coded OFDM on block fading channels," *IEEE Trans. Commun.*, vol. 57, no. 10, pp. 2872–2875, Oct. 2009.
- [44] K. Wang, T. Zhou, T. Xu, H. Hu, and X. Tao, "Asymmetric adaptive modulation for uplink NOMA systems," *IEEE Trans. Commun.*, early access, 2021.
- [45] W. Yu, H. Jia, and L. Musavian, "Joint adaptive M-QAM modulation and power adaptation for a downlink NOMA network," *IEEE Trans. Commun.*, early access, 2021.
- [46] W. Liang, Z. Ding, Y. Li, and L. Song, "User pairing for downlink non-orthogonal multiple access networks using matching algorithm," *IEEE Trans. Commun.*, vol. 65, no. 12, pp. 5319–5332, Dec. 2017.
- [47] F. Hussain, S. A. Hassan, R. Hussain, and E. Hossain, "Machine learning for resource management in cellular and IoT networks: Potentials, current solutions, and open challenges," *IEEE Commun. Surveys & Tuts.*, vol. 22, no. 2, pp. 1251–1275, 2nd Quart. 2020.
- [48] M. Chen, U. Challita, W. Saad, C. Yin, and M. Debbah, "Artificial neural networks-based machine learning for wireless networks: A tutorial," *IEEE Commun. Surveys & Tuts.*, vol. 21, no. 4, pp. 3039–3071, 4th Quart. 2019.
- [49] 3GPP TR 36.814, "Further advancements for E-UTRA physical layer aspects (Release 9)," Mar. 2010.
- [50] L. Qi, M. Peng, Y. Liu, and S. Yan, "Advanced user association in non-orthogonal multiple access based fog radio access networks," *IEEE Trans. Commun.*, vol. 67, no. 12, pp. 8408–8421, Dec. 2019.
- [51] C. Sun, C. She, C. Yang, T. Q. Quek, Y. Li, and B. Vucetic, "Optimizing resource allocation in the short blocklength regime for ultra-reliable and low-latency communications," *IEEE Trans. Wireless Commun.*, vol. 18, no. 1, pp. 402–415, Jan. 2019.
- [52] 3GPP Report ITU-R M.2412-0, "Guidelines for evaluation of radio interface technologies for IMT-2020," Oct. 2017.



Yuan Wang (S'19) received the B.E. degree in information engineering from the School of Information Science and Engineering, Southeast University, Nanjing, China, in 2017, where she is currently pursuing the Ph.D. degree in information and communication engineering. Her research interests include non-orthogonal multiple access, machine learning, and optimization theory for wireless communications.



Jiaheng Wang (M'10, SM'14) received the Ph.D. degree in electronic and computer engineering from the Hong Kong University of Science and Technology, Kowloon, Hong Kong, in 2010, and the B.E. and M.S. degrees from the Southeast University, Nanjing, China, in 2001 and 2006, respectively.

He is currently a Full Professor at the National Mobile Communications Research Laboratory (NCRL), Southeast University, Nanjing, China. From 2010 to 2011, he was with the Signal Processing Laboratory, KTH Royal Institute of Technology, Stockholm, Sweden. He also held visiting positions at the Friedrich Alexander University Erlangen-Nürnberg, Nürnberg, Germany, and the University of Macau, Macau. His research interests include optimization in signal processing and wireless communications.

Dr. Wang has published more than 130 articles on international journals and conferences. From 2014 to 2018, he served as an Associate Editor for the *IEEE Signal Processing Letters*. From 2018, he serves as a Senior Area Editor for the *IEEE Signal Processing Letters*. He was a recipient of the Humboldt Fellowship for Experienced Researchers and the best paper awards of *IEEE GLOBECOM* 2019, *ADHOCNETS* 2019, and *WCSP* 2014.



Vincent W.S. Wong (S'94, M'00, SM'07, F'16) received the B.Sc. degree from the University of Manitoba, Winnipeg, MB, Canada, in 1994, the M.A.Sc. degree from the University of Waterloo, Waterloo, ON, Canada, in 1996, and the Ph.D. degree from the University of British Columbia (UBC), Vancouver, BC, Canada, in 2000. From 2000 to 2001, he worked as a systems engineer at PMC-Sierra Inc. (now Microchip Technology Inc.). He joined the Department of Electrical and Computer Engineering at UBC in 2002 and is currently a Professor. His research

areas include protocol design, optimization, and resource management of communication networks, with applications to wireless networks, smart grid, mobile edge computing, and Internet of Things. Currently, Dr. Wong is the Chair of the Executive Editorial Committee of *IEEE Transactions on Wireless Communications*, an Area Editor of *IEEE Transactions on Communications* and *IEEE Open Journal of the Communications Society*, and an Associate Editor of *IEEE Transactions on Mobile Computing*. He has served as a Guest Editor of *IEEE Journal on Selected Areas in Communications*, *IEEE Internet of Things Journal*, and *IEEE Wireless Communications*. He has also served on the editorial boards of *IEEE Transactions on Vehicular Technology* and *Journal of Communications and Networks*. He was a Tutorial Co-Chair of *IEEE GLOBECOM*'18, a Technical Program Co-chair of *IEEE VTC2020-Fall* and *IEEE SmartGridComm*'14, as well as a Symposium Co-chair of *IEEE ICC*'18, *IEEE SmartGridComm* ('13, '17) and *IEEE GLOBECOM*'13. He is the Chair of the IEEE Vancouver Joint Communications Chapter and has served as the Chair of the IEEE Communications Society Emerging Technical Sub-Committee on Smart Grid Communications. He was an IEEE Communications Society Distinguished Lecturer (2019 - 2020).



Xiaohu You (Fellow, IEEE) received the B.S., M.S., and Ph.D. degrees in electrical engineering from the Nanjing Institute of Technology, Nanjing, China, in 1982, 1985, and 1989, respectively. From 1987 to 1989, he was a Lecturer with the Nanjing Institute of Technology. Since 1990, he has been with Southeast University, first as an Associate Professor and then as a Professor. He was the Premier Foundation Investigator of China National Science Foundation. From 1999 to 2002, he was the Principal Expert of the C3G Project, responsible for organizing China's

3G mobile communications research and development activities. From 2001 to 2006, he was the Principal Expert of the National 863 Future Project. His research interests include mobile communications, adaptive signal processing, artificial neural networks with applications to communications, and biomedical engineering. He is currently the Chairman of the IEEE Nanjing Section. He was selected as IEEE Fellow in 2012 for his contributions to the development of mobile communications in China. He was the recipient of the Excellent Paper Award from the China Institute of Communications in 1987 and the Elite Outstanding Young Teacher Award from Southeast University in 1990, 1991, and 1993.