# Non-Orthogonal Random Access and Data Transmission Scheme for Machine-to-Machine Communications in Cellular Networks

**YALI WU**[1], **NINGBO ZHANG**[2], **AND KAIXUAN RONG**[3]

[1]Institute of Information Technology, Hebei University of Economics and Business, Shijiazhuang 050061, China
[2]Key Laboratory of Universal Wireless Communications, Beijing University of Posts and Telecommunications, Beijing 100876, China
[3]The 54th Research Institute of China Electronics Technology Group Corporation, Shijiazhuang 050081, China

Corresponding author: Yali Wu (arlyarly@heuet.edu.cn)

**ABSTRACT** In order to address the signalling overhead and resource allocation problems of Machine-to-Machine (M2M) communications with non-orthogonal multiple access (NOMA), we propose a hybrid non-orthogonal random access and data transmission (NORA-DT) scheme. A novel design of NORA-DT protocol for M2M communications in cellular networks is firstly proposed. A power back-off scheme is introduced to adjust machine-type communications device (MTCD)'s target arrived power, and a closed-form analytic formula for the relation of MTCD's transmission power is derived. Based on the transmission power relation, the devices are clustered into a set of NOMA clusters. In the hybrid NORA-DT protocol, the cluster center MTCD transmits a extended preamble on behalf of the MTCDs in the same NOMA cluster on the physical random access channel (PRACH) for connection request. Base station (BS) can perfectly detect the preamble collisions in advance and schedules physical uplink shared channel (PUSCH) only to the NOMA clusters without collision. Then the MTCDs in the same NOMA clusters transmit data packets right after preamble transmission on the PUSCH to reduce the signalling overhead. By finding the optimal power allocation, we propose a low-complexity energy efficiency maximization problem for NORA-DT scheme. Due to the relation of MTCD's transmission power, we transform the problem into the function of cluster center MTCD's transmission power and solve it by difference of convex (DC) programming under the maximum transmission power constraints and minimum rate requirements at the MTCDs. A computationally efficient adaptive resource allocation scheme is finally proposed to improve the system throughput and resource usage. The optimal resource allocation between PRACH and PUSCH for any number of MTCDs can be learned by BS in advance, which avoids frequent computation. The analytic model is validated by simulation results. We demonstrate that the proposed NORA-DT scheme can significantly improve the system throughput, resource efficiency and energy efficiency performance.

**INDEX TERMS** Machine-to-machine communications, non-orthogonal multiple access, system throughput, resource efficiency, energy efficiency.

## I. INTRODUCTION

### A. MOTIVATION

As a key component of Internet of Things (IoT), machine-to-machine (M2M) communications, also known as machine-type communications (MTC) in the third-generation

The associate editor coordinating the review of this manuscript and approving it for publication was Maurizio Murroni.

partnership project (3GPP), is being regarded as one of the promising technologies in 5th-generation (5G) wireless communications [1]. With the wide application in smart grids, environment monitoring, cargo tracking, intelligent pay, health care, and so on, the number of MTC devices (MTCDs) is massively expanding [2]–[6]. Cisco Visual Networking Index report predicts that 14.6 billion MTCDs will connect the cellular networks by 2022 [7]. Random access (RA)

procedure is identified as a key step for initial access [8]. It is known that in conventional RA procedure, a MTCD transmits a preamble on the physical random access channel (PRACH) in the first step of RA procedure to inform the base station (BS) of its connection request. However, massive connections bring preamble collisions on PRACH, which seriously increase the network congestion and decreases the capacity of cellular network access.

Several studies mainly focus on conventional orthogonal access methods to alleviate the preamble collision and improve the capacity of access [9]–[14]. There have been some orthogonal multiple access (OMA) techniques, e.g. orthogonal frequency division multiple access (OFDMA), is adopted in the long term evolution (LTE) systems, where each resource block (RB) serves one MTCD to keep MTCD's orthogonality [15]. However, the continuous expansion of MTCDs still results in a shortage problem of radio resources. Due to the the scarcity of spectrum, how to improve the network access capacity is essentially a problem of how to make more efficient use of the radio resource.

As a key technology in 5G wireless communications, non-orthogonal multiple access (NOMA) can simultaneously enable multiple users to transmit in the same channels by splitting different users in the power domain, which yields a significant gain in spectrum efficiency [16]–[20]. This favorable character makes NOMA to be a promising access solution for supporting the massive MTCDs in cellular enabled M2M networks. To the best of our knowledge, recent studies on NOMA mainly focus on performance analysis of data transmission process. In practical networks such as LTE-Advanced (LTE-A) and future 5G networks, the introduction and realization of NOMA in the RA process could be very challenging.

### B. BACKGROUND AND RELATED WORKS

As in MTC networks, the traffic generated at uplink is heavier, the focus of this section is on the existing works that deploy NOMA for uplink scenarios. Resource allocation in NOMA systems is more challenging than that in OMA systems since power allocation among paired users needs to be carefully optimized to mitigate the co-channel interference among these users. The resource allocation in the NOMA-enabled communication network has been well studied in the literature.

In the light of different applications of NOMA, several allocation schemes have been summarized and suggested. Choi [21] propose a power allocation algorithm for coded uplink NOMA in a multicarrier system to enlarge the system throughput under the constraints of code word error probability. Mostafa *et al.* [22] propose a power-domain uplink NOMA scheme for narrowband IoT (NB-IoT) systems to overcome the challenge of providing connectivity to a large number of IoT devices. Tan *et al.* [23] investigate a dedicated millimeter-wave-based hybrid energy harvesting mechanism with NOMA transmission to increase network throughput. Liu *et al.* [24] investigate a joint

power allocation and user scheduling scheme for Device-to-Device(D2D) communications-enabled heterogeneous networks with NOMA to maximize the ergodic sum rate of small-cell near users. To analyze the performance of the NOMA scheme, the outage performance and the achievable sum data rate are theoretically analyzed [25]–[27]. Tweed *et al.* [25] derive an expression for the user outage probability as a function of successive interference cancelation (SIC) error variance. This result is used to a robust joint resource allocation problem is formulated to minimize user transmit power subject to rate and outage constraints for power-restricted but high priority devices. Considering a stepped level back-off based power allocation, Zhang *et al.* [26] derive an expression for sum rate and outage probability for two-user uplink NOMA systems. The integration of NOMA concepts in multiple-input multiple-output (MIMO) systems can support extra users and enhance the performance gains compared with the existing MIMO-OMA schemes. Therefore, several allocation schemes have been suggested MIMO-NOMA for uplink channels. Aghdam *et al.* [27] propose a new method to achieve lower outage probability for cellular M2M communication system with the mmWave massive-MIMO-NOMA transmission scheme. Ding *et al.* [28] provide an overview on the latest progress of MIMO-NOMA and Cognitiveradio NOMA (CR-NOMA). To achieve the promised gains, the allowed users are usually grouped in clusters based on their propagation channel conditions, and different strategies have been proposed for cluster/group formation [27], [29]–[31]. In [27], a random paring scheme is introduced to reduce the overhead of the system and achieve the quality of service (QoS). Ali *et al.* [29] consider the channel gain difference among users to form clusters and optimize their respective power allocations to increase throughput. Cabrera and Vesilo [30] propose an enhanced K-means clustering algorithm accompanied by NOMA, where each strong channel gain device is allocated to the appropriate cluster as a cluster head to enhance the network sum throughput. Aghdam *et al.* [31] propose a random user grouping with optimal beamforming coefficients in mmWave MIMO-NOMA transmission systems to reduce the system overhead for massive MTCDs connection, where MTCDs are grouped according to their distances from BS. By adjusting the power allocation coefficient, the fairness between the users in a pair can be enhanced [32]–[33]. Furthermore, by using user pairing algorithm, the fairness among use pairs can be further prompted. Pischella and Ruyet [32] focus on clustering and resource block (RB) allocation in multi-carrier uplink networks using power-domain NOMA. The optimization objective is time-based proportional fairness. Hojeij *et al.* [33] propose a proportional fairness based joint power allocation and user scheduler algorithm.

Apart from system throughput, energy efficiency is also an important factor in NOMA. With the rise in desire for green communications in recent years, reducing energy consumption has become of prime importance for researchers,

and 5G has also targeted energy efficiency as one of the major parameters to be achieved. Yang *et al.* [34] investigate an optimal power control and time scheduling scheme to achieve the optimal energy consumption for M2M communications with NOMA. Zeng *et al.* [35] propose an energy efficiency maximization problem for an uplink millimeter wave massive MIMO system with NOMA, and introduce an iterative algorithm to allocate the power for energy efficiency maximization. Gu *et al.* [36] study the power control for cognitive M2M communications underlaying cellular network, where MTCDs reuse the licensed spectrum of users in an opportunistic and fair manner. Yang *et al.* [37] study and compare two energy efficient resource allocation schemes with nonlinear energy harvesting and two different multiple access strategies, i.e., NOMA and time division multiple access (TDMA) for the cellular-enabled M2M network, where the energy consumption is reduced. Li and Gui [38] provide an energy-efficient resource allocation with hybrid TDMA-NOMA for cellular-enabled M2M networks, where MTCDs reuse the time slots of user equipments to further exploit the character of low power of MTCDs. Rozario and Hossain [39] apply *k*-mean clustering for machines as well as cluster head reselection method to balance the power consumption within the machines to increase their battery life. Alemaishat *et al.* [40] propose a joint sub-channel and power allocation algorithm for D2D communication based on NOMA to maximize the uplink energy efficiency and throughput of the mobile communication system. Na  *et al.* [41] investigate a joint uplink and downlink power allocation scheme for IoT based on NOMA to improve the energy and spectrum utilization.

Different form the aforementioned studies which mainly focus on performance analysis of data transmission process, Liang *et al.* [42] propose a non-orthogonal random access (NORA) scheme based on SIC, which utilizes the spatial distribution of MTCDs to allow multiple devices to use the same RBs for preamble transmission. However, the analysis in these schemes are mostly based on the assumption that physical uplink shared channel (PUSCH) resources are sufficient. Since the uplink available resources are limited. The more the PRACH resources allocated to alleviate preamble collision problem, the less the radio resources available for uplink data transmission. While if more RBs are allocated to PUSCH, preamble collision may be increased. Thereby, we propose a novel design of non-orthogonal random access and data transmission (NORA-DT) scheme for M2M communications in cellular networks, both the congestion in PRACH and data channels are emphasized.

## C. CONTRIBUTIONS

In this paper, we propose a hybrid non-orthogonal random access and data transmission (NORA-DT) scheme optimized for M2M communications, both the congestion in PRACH and data channels are emphasized. We investigate the performance of NORA-DT in terms of throughput, resource efficiency and energy efficiency. To provide comparisons,

we consider two benchmarks, which are referred to as traditional orthogonal random access (ORA) and orthogonal random access and data transmission (ORA-DT) [43]. Simulation results show that compared with ORA and ORA-DT scheme, our NORA-DT scheme can significantly improve the system throughput, resource efficiency and energy efficiency performance. The contribution of this paper is summarized as follows.

1) To reduce the signaling overhead, we propose a novel design of a hybrid NORA-DT protocol. The access procedure is simplified by allowing multiple MTCDs to send data on the same PUSCH right after preamble transmission on the PRACH without explicitly establishing a connection. Furthermore, the cluster center MTCD transmits a extended preamble on behalf of the MTCDs in the same cluster for data transmission. By the extended preamble, BS can perfectly detect the preamble collisions in advance and schedules PUSCH only to the NOMA clusters without collision. Therefore, resources wasting can be avoided.

2) To guarantee that BS receives different received power of multiplex MTCDs, we introduce a power back-off scheme to adjust MTCDs' target arrived power. Firstly, for a given set of NOMA clusters, power back-off step size is introduced to adjust MTCD's target arrived power and control the order of interference cancelation (IC). Secondly, the closed-form solution for the relation of MTCD's transmission power is formulated. Then based on the relation of MTCD's transmission power, the MTCDs are clustered into a set of NOMA clusters.

3) Considering the high computational complexity in solving the energy efficient power allocation problem, we derive a closed-form formula for energy efficiency as functions of the center cluster MTCD's transmission power. Since the optimization problem for energy efficiency maximization is non-convex, difference of convex (DC) programming is used to resolve the energy efficient power allocation under the maximum transmission power constraints and minimum rate requirements at the MTCDs. Then the transmission power of other MTCDs in the same cluster can be obtained by the relation with the cluster center MTCD.

4) To improve the system throughput and resource efficiency, we propose a low-complexity adaptive resource allocation scheme based on device number intervals in NORA-DT. A reasonable resource tradeoff between PRACH and PUSCH is achieved, and the resource allocation between PRACH and PUSCH for any number of MTCDs can be derived from the resource allocation for these device number intervals.

5) The analytic model is validated by simulation results. We demonstrate that the proposed scheme can achieve performance improvement in system throughput, resource efficiency and energy efficiency performance.

This paper is organized as follows: Section II presents system model and we provide a detailed description of the

proposed hybrid NORA-DT protocol for M2M communications in section III. In section IV, we respectively analyze the system throughput and present the adaptive resource allocation scheme based on device number intervals for NORA-DT. The performance is evaluated in section V, and section VI concludes this paper.

## II. SYSTEM MODEL

We consider a single-cell uplink transmission scenario. All the MTCDs and BS are equipped with a single antenna. The frequency resource is consist of $K$ subchannels. The bandwidth of each subchannel is $B$. Denote $k$ as index for the $k$-th subchannel, where $k \in \mathcal{K} = \{1, \ldots, K\}$. Assuming $M$ MTCDs' signals transmit with different transmission power on the $k$-th subchannel simultaneously. Denote $m$ as index for the $m$-th device where $m \in \mathcal{M} = \{1, \ldots, M\}$. Let $s_{k,m}$ be the $m$-th device's signal transmitted on the $k$-th subchannel with $E\left[s_{k,m}^2\right] = 1$. $p_{k,m}$ is the transmission power of the $m$-th device on the $k$-th subchannel. The received signal at BS on the $k$-th subchannel is given by

$$y_k = \sum_{m=1}^{M} \sqrt{p_{k,m}} s_{k,m} h_{k,m} + n_k \tag{1}$$

where $h_{k,m} = g_{k,m} l_{k,m}$ is the channel coefficient of the $k$-th subchannel from the $m$-th device to BS, and where $l_{k,m}$ is the large-scale path loss, and $g_{k,m}$ is small-scale fading coefficient. $n_k$ is the additive white Gaussian noise observed at BS with the noise power spectral density $N_0$.

To simplify the analysis, $l_{k,m}$ is modeled by Free-Space path loss model [44], i.e., $l_{k,m} = \frac{\sqrt{G_l}\lambda}{4\pi d_{k,m}}$. Where $G_l$ is product of the transmit and receive antenna field radiation patterns in the line-of-sight (LOS) direction. $\lambda$ is the signal weavelength and $d_{k,m}$ denotes the distance between the $i$-th device on the $k$-th subband and BS. Define $g_{k,m}$ as the Rayleigh fading channel gain of the $k$-th subchannel from the $m$-th device to BS. The probability density function (PDF) of $\left|g_{k,m}\right|^2$ can be written as

$$f_{\left|g_{k,m}\right|^2}(x) = \frac{1}{2\mu^2} \exp\left(-\frac{x}{2\mu^2}\right) \tag{2}$$

where $\mu$ is is the variance of the normal distribution $\mathcal{N}(0, \mu)$.

In uplink NOMA, SIC receiver is carried out at BS to split the overlapped signals [26]. The order of detection are usually based on the arrived power. Assuming before BS detects the $m$-th device's signal, it decodes the prior $i$-th ($i < m$) devices' signal first, then remove the signal from its observation, in a successive manner. However, the interference symbol from device $i > m$ cannot be removed and will be treated as noise. Therefore, the received SINR of the $m$-th device on the $k$-th subchannel can be written as

$$SINR_{k,m} = \frac{p_{k,m}\left|h_{k,m}\right|^2}{\sum_{i=m+1}^{M} p_{k,i}\left|h_{k,i}\right|^2 + \sigma_k^2} \tag{3}$$

where $\sigma_k^2$ is the noise power on the $k$-th subchannel at BS, and $\sigma_k^2 = N_0 B$. Denote $R_{k,m}$ as the achievable data rate of the

$m$-th device on the $k$-th subchannel. $R_{k,m}$ is given by

$$R_{k,m} = B\log_2\left(1 + SINR_{k,m}\right) \tag{4}$$

### 1) POWER BACK-OFF SCHEME

In this subsection, a power back-off scheme is adopted to guarantee that BS receives diverse received power of multiplex devices [26]. BS controls the order of detection by adjusting MTCD's target arrived power. In [26], the performance of outage probability as well as the achievable sum data rate for uplink NOMA has been investigated. This paper mainly focuses on energy-efficient power allocation under minimum data rate and maximum transmission power constraints as an optimization problem for the data transmission process. In LTE systems, BS requires that the arrived power from different users equals the same target power to avoid inter-user interferences. The basic transmit power for the respective user is determined [45]. Define $\rho_k$ as the difference of received power at BS among multiplex devices. The target arrived power of the $m$-th device on the $k$-th NOMA set is $p_{k,u} - (m-1)\rho_k$, where $p_{k,u}$ is the target arrived power of the first MTCD (i.e., the cluster center MTCD) in the $k$-th NOMA cluster. This design is beneficial to cancel the co-channel interferences successively. $p_{k,m}$ is expressed as

$$p_{k,m} = p_{k,u} - (m-1)\rho_k + 10\log_{10}(\varsigma_k) + wPL_{k,m} \tag{5}$$

where $\varsigma_k$ is the number of PUSCH RBs for the $k$-th NOMA cluster, $PL_{k,m}$ is the downlink path loss estimated by the $m$-th device on the $k$-th NOMA cluster. The factor $w$ denotes to compensate the path loss difference between downlink and uplink. Based on (5), the difference between $p_{k,m}$ and $p_{k,1}$ can be written in dB by

$$p_{k,m} - p_{k,1} = wPL_{k,m} - wPL_{k,1} - (m-1)\rho_k \tag{6}$$

Given that $x[dB] = 10\log_{10} x[watt]$, (6) can be expressed in watt by

$$10\log_{10}\frac{p_{k,m}}{p_{k,1}} = 10\log_{10}\frac{wPL_{k,m}}{wPL_{k,1}} - (m-1)\rho_k \tag{7}$$

$pL_{k,m}$ is modeled by Free-Space path loss model [44]. Therefore, the transmit power of the $m$-th MTCD on the $k$-th NOMA set is expressed by

$$p_{k,m} = p_{k,1}\frac{l_{k,1}^2}{l_{k,m}^2}10^{\frac{-(m-1)\rho_k}{10}} \tag{8}$$

For uplink NOMA, the key challenge is how BS derives diverse power among NOMA MTCDs as SIC receiver needs to separate them in power domain. MTCDs are needed to be ordered as SIC receiver needs to detect overlapped MTCDs in a successive manner. For SIC receiver, the order of detection are usually based on the arrived power, so BS controls the order of detection by adjusting MTCD's target arrived power. Fig. 1 illustrates a two-MTCD uplink NOMA cluster and gives the specific signal detecting process of SIC receiver. MTCD1 and MTCD2 experience channel gains of $h_{k,1}$ and $h_{k,2}$, respectively, where $p_{k,1}\left|h_{k,1}\right|^2 > p_{k,2}\left|h_{k,2}\right|^2$.
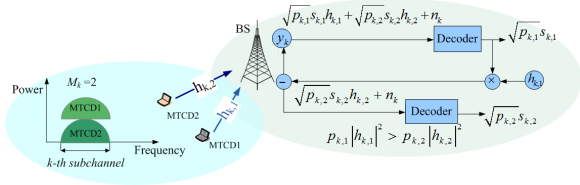
**FIGURE 1.** Illustration of a 2-MTCD uplink NOMA cluster with SIC at the BS.

The MTCD's signal with the higher arrived power is decoded first at BS. Before BS detects MTCD2's signal, it decodes MTCD1's signal first, then remove the signal from its observation, in a successive manner. However, the interference symbol from MTCD2 cannot be removed and will be treated as noise. Thus, the achievable data rate of MTCD1 depends on the interferences from MTCD2, whereas MTCD2 achieves an intra-cell interference-free data rate.

Based on (8) and $|h_{k,m}|^2 = |g_{k,m}|^2 l_{k,m}^2$, we can get

$$\frac{p_{k,m}|h_{k,m}|^2}{p_{k,m+1}|h_{k,m+1}|^2} = \frac{|g_{k,m}|^2}{|g_{k,m+1}|^2} 10^{\frac{\rho_k}{10}} \quad (9)$$

When $|g_{k,m}|^2 \geq |g_{k,m+1}|^2$, we have $p_{k,m}|h_{k,m}|^2 > p_{k,m+1}|h_{k,m+1}|^2$. In order to save MTCD's power consumption, the order of MTCDs could be based on the Rayleigh fading coefficient between MTCD and BS, i.e., the MTCD with a smaller Rayleigh fading coefficient would be assigned to a larger order, which means the corresponding target arrived power is smaller than other MTCDs. Therefore, $M$ MTCDs are allocated on the $k$-th subchannel with order

$$|g_{k,1}|^2 \geq \cdots \geq |g_{k,m}|^2 \geq |g_{k,m+1}|^2 \geq \cdots \geq |g_{k,M}|^2 \quad (10)$$

In order to save MTCD's power consumption, the order of MTCDs could be based on the pathloss between MTCD and BS, i.e., the MTCD with a smaller Rayleigh fading channel gain would be assigned to a larger order, which means the corresponding target arrived power is smaller than other MTCDs.

### 2) MINIMUM DATA RATE CONSTRAINT

Our design is based on providing QoS guarantees that each device meets the corresponding minimum rate requirement. Each device has a minimum data rate, denoted as $\hat{R}_{k,m}$. When $R_{k,m} \geq \hat{R}_{k,m}$, BS can successfully detect the signal of the $m$-th device on the $k$-th subchannel, otherwise, outage happens. Substituting (3) into (4), and let $R_{k,m} \geq \hat{R}_{k,m}$, the constraint $R_{k,m} \geq \hat{R}_{k,m}$ can be converted to

$$p_{k,m}|h_{k,m}|^2 - \varphi_{k,m}\sum_{i=m+1}^{M} p_{k,i}|h_{k,i}|^2 \geq \varphi_{k,m}\sigma_k^2 \quad (11)$$

where $\varphi_{k,m} = 2^{\hat{R}_{k,m}/B} - 1$, $m = 1, 2, \ldots, M$.
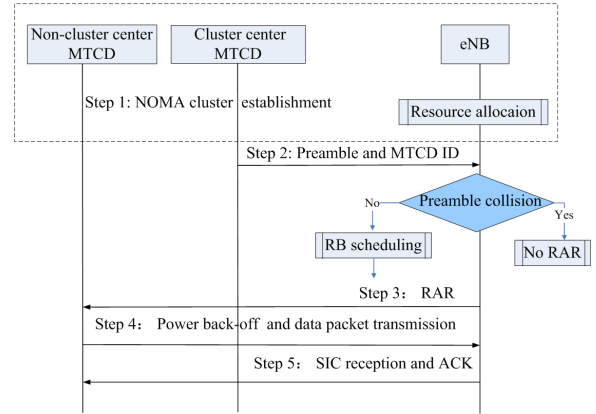


**FIGURE 2.** A new hybrid NORA-DT protocol.

Substituting (8) and $|h_{k,m}|^2 = |g_{k,m}|^2 l_{k,m}^2$ into (11), we can get

$$|g_{k,m}|^2 - \varphi_{k,m}\sum_{i=m+1}^{M} \frac{|g_{k,i}|^2}{10^{\frac{(i-m)\rho_k}{10}}} \geq \frac{\varphi_{k,m}\sigma_k^2}{p_{k,1}l_{k,1}^2} 10^{\frac{(m-1)\rho_k}{10}} \quad (12)$$

where $m = 1, 2, \ldots, M$.

### 3) MAXIMUM TRANSMISSION POWER CONSTRAINT

Define $p_{k,m}^{\max}$ as the maximum transmission power of the $m$-th MTCD in the $k$-th NOMA cluster. Substituting (8) into $p_{k,m} \leq p_{k,m}^{\max}$, we can get

$$p_{k,1} \leq p_{k,m}^{\max} \frac{l_{k,m}^2}{l_{k,1}^2} 10^{\frac{(m-1)\rho_k}{10}} \quad (13)$$

## III. NON-ORTHOGONAL RANDOM ACCESS AND DATA TRANSMISSION PROTOCOL

In this section, we introduce a new hybrid NORA-DT protocol for M2M communications. As Fig.2 shows, in our proposed protocol, we have five steps executed in sequence in each cycle, which are i) NOMA cluster establishment, ii) Preamble transmission, iii) Random access response (RAR), iv) Power back-off and data packet transmission, and v) SIC reception and acknowledge (ACK). From now on, we explain the above steps one by one.

### A. NOMA CLUSTER ESTABLISHMENT

This step is to discover NOMA clusters and select the cluster center MTCDs. The number of NOMA clusters as well as the device clustering strategy will be discussed in this subsection.

Different from traditional clustering strategy based on the distinct channel coefficients, the proposed device clustering algorithm utilize the information of channel coefficients, minimal data rate requirements and maximum transmission power constraints at the devices simultaneously. With the increase of the number of MTCDs, the available information of different MTCDs are more abundant, which is better for device clustering. We consider the wireless channels are independent and identically distributed (i.i.d.) block Rayleigh

fading, which means the fading channel gain is constant during a frame. Assuming that perfect channel state information (CSI) is available at BS for each frame [46]. BS can use the CSI for NOMA cluster establishmen and power allocation. Then BS informs the configuration to cluster center MTCDs. Without loss of generality, the noise power on the $k$-th subchannel at BS is defined as $\sigma^2$, the power back-off size on he $k$-th subchannel is $\rho$, and the target arrived power of the cluster center MTCD on the $k$-th subchannel is defined as $p_u$. For the sake of notational simplicity, we define X(Y) as the Y-th element of X. The MTCDs clustering strategy is summarized in Algorithm 1, in which the following factors are considered.

a) According to (10), the order of SIC in NOMA is usually based on the descending order of the Rayleigh fading coefficient. The Rayleigh fading coefficients of $q$ MTCDs are first arranged from large to small.

b) As the Phase I shows, the devices can be clustered together if the channel conditions of the MTCDs belong to the condition of (13), which satisfies the transmission power constraints of each multiplexing MTCDs.

c) As the Phase II shows, the devices can be clustered together if the channel conditions of the MTCDs belong to the condition of (12), which satisfies the data rate constraint of each multiplexing MTCD.

In Algorithm 1, the MTCD with the largest Rayleigh fading coefficient is selected as a cluster center MTCD, and Phase I-Phase II are used to judge whether other MTCDs can be assigned to a NOMA cluster with the cluster center MTCD. If these devices can be clustered together, from the remaining MTCDs after removing these devices, the MTCD with the largest Rayleigh fading channel coefficient is selected as the cluster center MTCD, and then judgment is made according to Phase I-Phase II. If clustering fails, from the remaining MTCDs after removing the MTCD with the largest Rayleigh fading channel coefficient, the MTCD with the largest Rayleigh fading channel coefficient is selected as the cluster center user, and the judgment is made according to the above process. By analogy, each NOMA cluster and cluster center MTCD can be identified. We can get the number of NOMA clusters $u = k$ when $\bar{G} = \emptyset$.

### B. PREAMBLE TRANSMISSION

The cluster center MTCD transmits not only preamble but also the identity (ID) information and a cyclic redundancy check (CRC) on PRACH. We assume MTCD ID and CRC are mapped to the same position in the guard band, and preamble sequences are still mapped to the central 839 random access channel (RACH) subcarriers. After that, the cluster center MTCD broadcasts its transmission power, the list of power back-off index and corresponding MTCDs identity to the non-cluster center MTCDs in the same cluster. The index of the transmitted preamble on PRACH is also attached. Non-cluster center MTCD checks if its MTCD identity is included in the list. If included, it goes to next step to receive RAR. Otherwise, it reattempts in the next RA cycle.

---

**Algorithm 1** Device Clustering Algorithm

1: Define $G = \bar{G} = \left\{ |g_1|^2, |g_2|^2, \ldots, |g_q|^2 \right\}$ as the set of candidate devices' Rayleigh fading coefficient, where $|g_1|^2 \geq |g_2|^2 \geq \cdots \geq |g_q|^2$. Define $L = \bar{L} = \left\{ l_1^2, l_2^2, \ldots, l_q^2 \right\}$, $R^\dagger = \bar{R}^\dagger = \left\{ R_1^\dagger, R_2^\dagger, \ldots, R_q^\dagger \right\}$ and $P = \bar{P} = \left\{ p_1^{\max}, p_2^{\max}, \ldots, p_q^{\max} \right\}$ as the set of above candidate devices' large scale fading coefficient, the minimal data rate and the maximum allowed transmission power, respectively. $J = |G|$. Initialize the number of NOMA clusters as $k = 0$.
2: **while** $\bar{G} \neq \emptyset$ **do**
3:   $k = k + 1$. Initialize the number of multiplex devices in the $k$-th NOMA cluster as $i = 0$. Initialize the set of multiplex devices' Rayleigh fading coefficient, large scale fading coefficient and minimal data rate as $\Lambda_k^{(i)} = \emptyset$, $\Gamma_k^{(i)} = \emptyset$, and $\Psi_k^{(i)} = \emptyset$, respectively. Initialize the transmission power of the cluster center MTCD as $p_1$. $s = 0$.
4:   **PhaseI** :
5:   **while** $i < M$ **do**
6:     **for** $j = s + 1$ to $J$ **do**
7:       **if** $p_1 \leq P(j) \frac{L(j)}{L(1)} 10^{\frac{(j-1)\rho}{10}}$ **then**
8:         $\Lambda_k^{(i+1)} = \left\{ \Lambda_k^{(i)}, G(j) \right\}, \Gamma_k^{(i+1)} = \left\{ \Gamma_k^{(i)}, L(j) \right\}$, $\Psi_k^{(i+1)} = \left\{ \Psi_k^{(i)}, R^\dagger(j) \right\}$. $s = j$. $i = i + 1$.
9:         *break*;
10:       **else**
11:         $G = G \backslash G(s)$. $L = L \backslash L(s)$. $R^\dagger = R^\dagger \backslash R^\dagger(s)$. $P = P \backslash P(s)$.
12:       **end if**
13:     **end for**
14:   **end while**
15:   **PhaseII** :
16:   **for** $m = 1$ to $M - 1$ **do**
17:     $a(m) = \Gamma_k^{(i)}(1)/\sigma^2 \sum_{i=m}^{M-1} 10^{\frac{-(i-1)\rho}{10}} \Lambda_k^{(i)}(1)$. $b(m) = \Gamma_k^{(i)}(1)/\sigma^2 \sum_{i=m+1}^{M-1} 10^{\frac{-(i-1)\rho}{10}} \Lambda_k^{(i)}(1)$.
18:     **if** $\frac{2^{\Psi_k^{(i)}(m)}-1}{a(m)-2^{\Psi_k^{(i)}(m)}b(m)} > p_1$ (Based on (12)) **then**
19:       $i = i - 1$. $G = \bar{G} = \bar{G} \backslash G(s)$. $L = \bar{L} = \bar{L} \backslash L(s)$. $R^\dagger = \bar{R}^\dagger = \bar{R}^\dagger \backslash R^\dagger(s)$. Go to 2.
20:     **else**
21:       $i = i - 1$. $G = \bar{G} = \bar{G} \backslash \Lambda_k^{(i)}$, $L = \bar{L} = \bar{L} \backslash \Gamma_k^{(i)}$, $R^\dagger = \bar{R}^\dagger = \bar{R}^\dagger \backslash \Psi_k^{(i)}$. Go to 2.
22:     **end if**
23:   **end for**
24: **end while**

---

The length of MTCD ID and the position of mapping subcarriers are predefined. The number of subcarriers which the spread MTCD ID occupies can be calculated as $n_{sub} = \frac{n_{ID} \cdot F}{M}$, where $n_{ID}$ is the number of bits of MTCD ID, $F$ is the spreading factor, and $M$ is the modulation order. If the
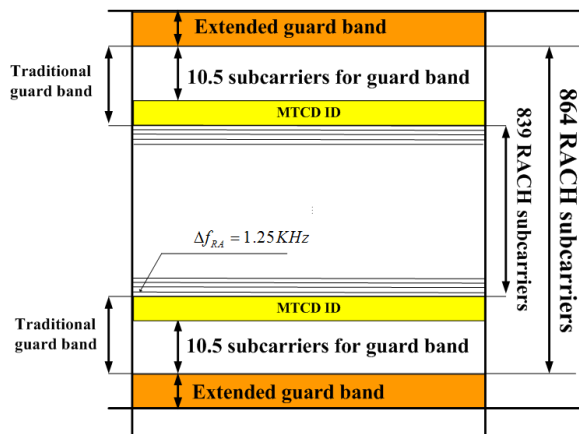
**FIGURE 3.** The extension of guard band for PRACH mapping.

modulation mode is quadrature phase shift keying (QPSK), 4 bits of MTCD ID is spread by 2 bits of orthogonal code, then 4 subcarriers are needed for mapping (yellow part in Fig.3). However, if the modulation mode is 16 Quadrature Amplitude Modulation (16QAM), 8 bits of MTCD ID is spread by 16 bits of orthogonal code, 32 subcarriers for mapping. Given that there are only 25 subcarriers in guard band, the traditional area of guard band needs to be extended. Fig.3 shows the extension of guard band. The extension of guard band may impact the scheduling of PUSCH. However, in LTE systems, the subcarrier spacing for PRACH and PUSCH is 1.25kHz and 15kHz, respectively. In other words, the size of a PUSCH subcarrier can hold 12 PRACH subcarriers. It is expected the impact of extended guard band to PUSCH scheduling is limited since very fewer PUSCH subcarriers are served as PRACH guard band. For example, only 1 PUSCH subcarrier is sufficient to meet the requirement of 8 bits of MTCD ID with 16QAM modulation mode (the spreading factor is 16).

### C. PREAMBLE DETECTION AND RAR TRANSMISSION

On receiving preambles on PRACH, BS first calculates the power delay profile (PDP) of a preamble to detect if this preamble is selected [43]. If this preamble is selected, BS decodes the MTCD ID to check if this selected preamble is collided. Given that every MTCD has a unique ID, when more than one cluster center MTCDs selects the same preamble and transmits their IDs on the same predefined subcarriers, BS cannot correctly decode MTCD ID because of the interferences. Thus, BS thinks that collisions happen and will not schedule PUSCH to this preamble. Therefore, resources wasting can be avoided. Upon correctly decoding MTCD ID, BS schedules PUSCH and sends corresponding RAR through downlink channel. Given all MTCDs in a NOMA cluster know the index of transmit preamble sequence, they are expected to receive the same RAR.

Code division multiple access (CDMA) is used to distinguish each MTCD, namely, the ID and CRC bits are encoded by gold sequence [47]. Denote $a_i(t) = \sum_{n=0}^{\infty} a_n g_a(t - nT_a)$

as the ID of the $i$-th MTCD, where $a_n$ is bipolar information code, $T_a$ is the code element width, and $g_a$ is gate-function. $c_i(t) = \sum_{n=0}^{\infty} c_n g_c(t - nT_c)$ is the pseudo-random sequence of the $i$-th MTCD, where $c_n$ is spread code element, $T_c$ is the chip width, and $g_c$ is gate-function. The information of the $i$-th MTCD after spread spectrum modulation is

$$d_i(t) = a_i(t) c_i(t) = \sum_{n=0}^{\infty} d_n g_c(t - nT_c) \quad (14)$$

where $d_n = \begin{cases} +1, a_n = c_n \\ -1, a_n \neq c_n \end{cases}$. Assuming there are $I$ MTCDs that select the same preamble. Then the information of multi-MTCDs after spread spectrum modulation can be expressed as

$$d(t) = a_i(t) c_i(t) + \sum_{j=1, j \neq i}^{I} a_j(t) c_j(t) \quad (15)$$

Given that the transmission signal is modulated by BPSK, then

$$s(t) = a_i(t) c_i(t) \cos w_0 t + \sum_{i \neq j} a_j(t) c_j(t) \cos w_0 t \quad (16)$$

where $w_0$ is the carrier of the signal. For simplicity, the received signal at receiver through the channel can expressed as follows

$$r(t) = s_I(t) + n_I(t) + s_J(t) \quad (17)$$

where $s_I(t)$ is the desired signal, $n_I(t)$ is the channel noise, and $s_J(t)$ is the interference caused by multi-MTCD. Then, we use pseudo-random sequence $c'_i(t)$ which is same to transmitter to do correlation dispreading.

$$\begin{aligned} r'(t) &= r(t) c'_i(t) = (s_I(t) + n_I(t) + s_J(t)) c'_i(t) \\ &= s_I(t) c'_I(t) + n_I(t) c'_i(t) + s_J(t) c'_i(t) \\ &= s'_I(t) + n'_I(t) + s'_J(t) \end{aligned} \quad (18)$$

where $s'_I(t) = s_I(t) c'_i(t) = a_i(t) c_i(t) c'_i(t) \cos w_i t$. If $c(t) c'(t) = 1$. If $c_i(t) c'_i(t) = 1$, $s'_i(t) = a_i(t) \cos w_i t$. However, BS cannot decode $a_i(t)$ because of the interference from $s'_J(t) = \sum_{j=1, j \neq i}^{I} a_j(t) c_j(t) c'_i(t) \cos w_j t$. Therefore, the MTCD ID can be decoded correctly only when $i = j$.

### D. POWER BACK-OFF AND DATA PACKET TRANSMISSION

When the MTCD in a NOMA cluster receives RAR, it checks if the RAR message matches the preamble sequence sent by cluster center MTCD. Upon receiving a matching RAR, the MTCD adjusts the transmit power based on power back-off index and transmits data packet on the assigned subchannels.

#### 1) ENERGY EFFICIENCY MAXIMIZATION PROBLEM FORMULATION

In this subsection, we formulate the energy-efficient power allocation as an optimization problem for the data transmission process in NORA-DT scheme. We desire to maximize the energy efficiency while providing devices' QoS

guarantees by finding the optimal power allocation. Energy efficiency (EE) is defined as the ratio of the achievable sum rate of the devices to the total power consumption, which is given by $EE = \sum_{k=1}^{K} \frac{R_k}{P_k}$, where $R_k$ is the achievable sum rate of the devices on the $k$-th subchannel, and $P_k$ is the sum transmission power of the devices on the $k$-th subchannel. $R_k$ and $P_k$ can be written as

$$R_k = B \sum_{m=1}^{M} \log_2 \left( 1 + \frac{p_{k,m} |h_{k,m}|^2}{\sum_{i=m+1}^{M} p_{k,i} |h_{k,i}|^2 + \sigma_k^2} \right) \quad (19)$$

$$P_k = \sum_{m=1}^{M} p_{k,m} \quad (20)$$

Since $EE = \sum_{k=1}^{K} \frac{R_k}{P_k}$, we derive the optimal power allocation policy that maximizes the EE per NOMA cluster and in turn maximizes the overall system energy efficiency. Thereby, the EE maximization problem is formulated as

$$\max_{p_{k,1}} \frac{R_k}{P_k}$$

$$subject\ to: C_1 : p_{k,1} \leq p_{k,m}^{\max} \frac{l_{k,m}^2}{l_{k,1}^2} 10^{\frac{(m-1)\rho_k}{10}}, \forall m \in \mathcal{M}$$

$$C_2 : |g_{k,m}|^2 - \varphi_{k,m} \sum_{i=m+1}^{M} \frac{|g_{k,i}|^2}{10^{\frac{(i-m)\rho_k}{10}}}$$

$$\geq \frac{\varphi_{k,m} \sigma_k^2 10^{\frac{(m-1)\rho_k}{10}}}{p_{k,1} l_{k,1}^2}, \forall m \in \mathcal{M}$$

$$C_3 : |g_{k,m}|^2 \geq |g_{k,m+1}|^2, \forall m \in \mathcal{M} \setminus \{M\} \quad (21)$$

Based on the relation of device's transmission power in (8), $\frac{R_k}{P_k}$ is derived as the function of cluster center MTCD's transmission power, i.e., $p_{k,1}$. $\frac{R_k}{P_k}$ is shown at the top of next page. Constraint $C_1$ guarantees the power constraint of each device. Constraint $C_2$ ensures the minimum data rate requirement of each device. In order to use the DC programming approach, we can convert (22), as shown at the bottom of this page, to DC representation that can be simply written by (23), as shown at the bottom of this page, where

$$a(m) = l_{k,1}^2 \Big/ \sigma_k^2 \sum_{i=m}^{M} 10^{\frac{-(i-1)\rho_k}{10}} |g_{k,i}|^2 \quad (24)$$

$$b(m) = l_{k,1}^2 \Big/ \sigma_k^2 \sum_{i=m+1}^{M} 10^{\frac{-(i-1)\rho_k}{10}} |g_{k,i}|^2 \quad (25)$$

$$d = l_{k,1}^2 \sum_{m=1}^{M} \frac{1}{l_{k,m}^2} 10^{\frac{-(m-1)\rho_k}{10}} \quad (26)$$

Note that the optimization problem in (23) is non-convex with respect to $p_{k,1}$. We can convert $\max \frac{R_k}{P_k}$ to $\min \left( -\frac{R_k}{P_k} \right)$. Let $f(p_{k,1}) = -B \sum_{m=1}^{M} \log_2 (1 + a(m) p_{k,1})$, $g(p_{k,1}) = -B \sum_{m=1}^{M} \log_2 (1 + b(m) p_{k,1})$, then $-\frac{R_k}{P_k} = \frac{f(p_{k,1}) - g(p_{k,1})}{dp_{k,1}}$. Thereby, the EE maximization problem is converted as

$$\min_{p_{k,1}} \frac{f(p_{k,1}) - g(p_{k,1})}{dp_{k,1}}$$

$$subject\ to: C_1 : p_{k,1} \leq p_{k,m}^{\max} \frac{l_{k,m}^2}{l_{k,1}^2} 10^{\frac{(m-1)\rho_k}{10}}, \forall m \in \mathcal{M}$$

$$C_2 : p_{k,1} \geq \frac{\varphi_{k,m}}{a(m) - (\varphi_{k,m} + 1) b(m)}, \forall m \in \mathcal{M}$$

$$C_3 : |g_{k,m}|^2 \geq |g_{k,m+1}|^2, \forall m \in \mathcal{M} \setminus \{M\} \quad (27)$$

### 2) RESOLUTION FOR OPTIMAL PROBLEM

The gradient of $f(p_{k,1})$ and $g(p_{k,1})$ is denoted by $\nabla f(p_{k,1})$ and $\nabla g(p_{k,1})$, respectively. $f(p_{k,1})$ and $g(p_{k,1})$ are convex functions with respect to $p_{k,1}$ because $\nabla^2 f(p_{k,1}) > 0$ and $\nabla^2 g(p_{k,1}) > 0$. Proposition 1 proves the quasi-convexity of $\frac{f(p_{k,1})}{dp_{k,1}}$ and $\frac{g(p_{k,1})}{dp_{k,1}}$. Therefore, we can use the DC programming approach to realize energy-efficient power allocation.

*Proposition 1*: If $-f(p_{k,1}) = B \sum_{m=1}^{M} \log_2 (1 + a(m) p_{k,1})$ and $-g(p_{k,1}) = B \sum_{m=1}^{M} \log_2 (1 + b(m) p_{k,1})$ are strictly concave in $p_{k,1}$, $-\frac{f(p_{k,1})}{dp_{k,1}}$ and $-\frac{g(p_{k,1})}{dp_{k,1}}$ are quasi-concave. Inspired by [48], we can prove Proposition 1 as follows.

*Proof*: Denote the $\tau$-sublevel sets of function $-\frac{f(p_{k,1})}{dp_{k,1}}$ as

$$S_\tau = \left\{ p_{k,1} > 0 \left| -\frac{f(p_{k,1})}{dp_{k,1}} \geq \tau \right. \right\} \quad (28)$$

Based on the Proposition 1, $-\frac{f(p_{k,1})}{dp_{k,1}}$ is is strictly quasi-concave if and only if $S_\tau$ is strictly convex for any $\tau$. In this case, when $\tau < 0$, there are no points satisfying $-\frac{f(p_{k,1})}{dp_{k,1}} = \tau$. Therefore, $S_\tau$ is strictly convex when $\tau \leq 0$. When $\tau > 0$, we can rewrite $S_\tau$ as $S_\tau = \left\{ p_{k,1} > 0 \left| d\tau p_{k,1} + f(p_{k,1}) \leq 0 \right. \right\}$. Since $f(p_{k,1})$ is strictly convex in $p_{k,1}$, $S_\tau$ is therefore also strictly convex. Hence,

$$\frac{R_k}{P_k} = \frac{B \sum_{m=1}^{M} \log_2 \left( 1 + \frac{p_{k,1} l_{k,1}^2 10^{\frac{-(m-1)\rho_k}{10}} |g_{k,m}|^2}{p_{k,1} l_{k,1}^2 \sum_{i=m+1}^{M} 10^{\frac{-(i-1)\rho_k}{10}} |g_{k,i}|^2 + \sigma_k^2} \right)}{p_{k,1} l_{k,1}^2 \sum_{m=1}^{M} \frac{1}{l_{k,m}^2} 10^{\frac{-(m-1)\rho_k}{10}}} = \frac{B \sum_{m=1}^{M} \log \left( \frac{p_{k,1} l_{k,1}^2 \Big/ \sigma_k^2 \sum_{i=m}^{M} 10^{\frac{-(i-1)\rho_k}{10}} |g_{k,i}|^2 + 1}{p_{k,1} l_{k,1}^2 \Big/ \sigma_k^2 \sum_{i=m+1}^{M} 10^{\frac{-(i-1)\rho_k}{10}} |g_{k,i}|^2 + 1} \right)}{p_{k,1} l_{k,1}^2 \sum_{m=1}^{M} \frac{1}{l_{k,m}^2} 10^{\frac{-(m-1)\rho_k}{10}}} \quad (22)$$

$$\frac{R_k}{P_k} = \frac{B \sum_{m=1}^{M} \log_2 (1 + a(m) p_{k,1}) - B \sum_{m=1}^{M} \log_2 (1 + b(m) p_{k,1})}{dp_{k,1}} \quad (23)$$

$-\frac{f(p_{k,1})}{dp_{k,1}}$ and $-\frac{g(p_{k,1})}{dp_{k,1}}$ are strictly quasi-concave. Therefore, $\frac{f(p_{k,1})}{dp_{k,1}}$ and $\frac{g(p_{k,1})}{dp_{k,1}}$ are quasi-convex.

The proposed power allocation algorithm is presented in Algorithm 2. The function $\frac{g(p_{k,1})}{dp_{k,1}}$ can be approximated by its first-order Taylor expansion at $p_{k,1}^{(l)}$. i.e., $\frac{g(p_{k,1}^{(l)})}{dp_{k,1}^{(l)}} + \nabla \frac{g(p_{k,1}^{(l)})}{dp_{k,1}^{(l)}}\left(p_{k,1} - p_{k,1}^{(l)}\right)$ in each iteration, where $\nabla \frac{g(p_{k,1}^{(l)})}{dp_{k,1}^{(l)}}$ denotes the gradient of $\frac{g(p_{k,1}^{(l)})}{dp_{k,1}^{(l)}}$ at $p_{k,1}^{(l)}$.

---

**Algorithm 2** DC Programming Algorithm for the $k$-th Subchannel Device's Power Allocation

---

1: Initialize $p_{k,1}^{(0)}$, set iteration number $l = 0$. The convex functions $\frac{f(p_{k,1})}{dp_{k,1}}$ and $\frac{g(p_{k,1})}{dp_{k,1}}$.

2: **while** $\left| q\left(p_{k,1}^{(l+1)}\right) - q\left(p_{k,1}^{(l)}\right) \right| > \varepsilon$ **do**

3: Define convex approximation of $q^{(l)}\left(p_{k,1}\right)$ as
$$q^{(l)}\left(p_{k,1}\right) = \frac{f(p_{k,1})}{dp_{k,1}} - \frac{g(p_{k,1}^{(l)})}{dp_{k,1}^{(l)}} - \frac{\nabla g(p_{k,1}^{(l)})\left(p_{k,1}-p_{k,1}^{(l)}\right)}{dp_{k,1}^{(l)}}$$

4: Solve $\min_{p_{k,1}} q^{(l)}\left(p_{k,1}\right)$ by convex programming solvers such as CVX [49] to obtain the optimal solution $p_{k,1}^{opt}$.

5: Set $p_{k,1}^{(l+1)} = p_{k,1}^{opt}$. $l = l + 1$.

6: **end while**

---

### 3) CONVERGENCE AND COMPLEXITY ANALYSIS

In each iteration of Algorithm 2, the solution $p_{k,1}^{(l+1)}$ is generated as the optimal solution at the last iteration. Thus, we have $\frac{f(p_{k,1}^{(l)})}{dp_{k,1}^{(l)}} - \frac{g(p_{k,1}^{(l)})}{dp_{k,1}^{(l)}} \geq \frac{f(p_{k,1}^{(l+1)})}{dp_{k,1}^{(l+1)}} - \frac{g(p_{k,1}^{(l)})}{dp_{k,1}^{(l)}} - \nabla \frac{g(p_{k,1}^{(l)})}{dp_{k,1}^{(l)}}\left(p_{k,1}^{(l+1)} - p_{k,1}^{(l)}\right)$. As function $\frac{g(p_{k,1})}{dp_{k,1}}$ is convex, we have that $\frac{g(p_{k,1})}{dp_{k,1}} \geq \frac{g(p_{k,1}^{(l)})}{dp_{k,1}^{(l)}} + \nabla \frac{g(p_{k,1}^{(l)})}{dp_{k,1}^{(l)}}\left(p_{k,1} - p_{k,1}^{(l)}\right)$ at any $p_{k,1}$. Therefore, we derive that $\frac{f(p_{k,1}^{(l)})}{dp_{k,1}^{(l)}} - \frac{g(p_{k,1}^{(l)})}{dp_{k,1}^{(l)}} \geq \frac{f(p_{k,1}^{(l+1)})}{dp_{k,1}^{(l+1)}} - \frac{g(p_{k,1}^{(l+1)})}{dp_{k,1}^{(l+1)}}$. This implies that the objective value $\frac{f(p_{k,1}) - g(p_{k,1})}{dp_{k,1}}$ is reduced after each iteration. In the algorithm, the iterative process is terminated when $\left| q\left(p_{k,1}^{(l+1)}\right) - q\left(p_{k,1}^{(l)}\right) \right| \leq \varepsilon$. Because $p_{k,1}^{(l+1)} = p_{k,1}^{opt}$ in Algorithm 2, we can obtain $\left| p_{k,1}^{(l+1)} - p_{k,1}^{(l)} \right| \leq \varepsilon^*$ if $\left| q\left(p_{k,1}^{(l+1)}\right) - q\left(p_{k,1}^{(l)}\right) \right| \leq \varepsilon^*$, where $\varepsilon^*$ is a error tolerance parameter. Thus, the sequence of $\left\{p_{k,1}^{(l)}\right\}$ generated by Algorithm 2 is a Cauchy sequence. In addition, since the constraint set is compact, the sequence of always converges by Cauchy theorem [50].

In the proposed energy efficient power allocation scheme, sophisticated calculations can be avoided. Since the closed-form solution for the relation of MTCD's transmission power is formulated in (8), the transmission power of other MTCDs in a cluster can be achieved by the transmission power of the cluster center MTCD. Assuming the energy efficient power allocation for the cluster center MTCD is derived at $L$-th iteration. Then the energy efficient target arrived power of the cluster center MTCD is $p_{k,u}^{ee} = p_{k,1}^{L}\left/\left(\varsigma\left(1/l_{k,1}^2\right)^w\right)\right.$. The energy efficient target arrived power of the non-cluster center MTCDs are $p_{k,u}^{ee} - (k-1)\rho$.

### E. SIC RECEPTION AND ACK

BS already gets the number of multiplex MTCDs by distinguishing the preamble sub-subsets in step 3. As all subchannels are pairwise orthogonal, BS could decode messages from each subchannel independently. It performs SIC and decodes the data packets one by one. After that, BS sends the acknowledge message and corresponding MTCD identity via control channel.

## IV. THROUGHPUT ANALYSIS AND ADAPTIVE RESOURCE ALLOCATION BASED ON DEVICE NUMBER INTERVALS

### A. THROUGHPUT ANALYSIS

In this subsection, we analyze the throughput of the proposed scheme within a cycle. In a RA procedure, two uplink channels are required, i.e., PRACH for preamble transmission and PUSCH for data transmission. Since the uplink available resources are limited. The more the PRACH resources allocated to alleviate preamble collision problem, the less the radio resources available for uplink data transmission. While if more RBs are allocated to PUSCH, preamble collision may be increased. Therefore, the number of RBs allocated to PRACH and PUSCH will be discussed in this subsection. In LTE systems, a PRACH consists of 6 RBs in a subframe, which occupies 864 subcarriers [51]. As shown in Fig.4, a periodic sequence of time-frequency resources called random access slots (RA slots) are reserved in the PRACH for preamble transmission. In time domain, the duration of RA slot depends on the preamble format. There are at least one RA slot per two LTE frames and at most ten RA slots per LTE frame [51].

Assuming the total number of uplink RBs allocated to M2M communications is $Q$. We assume $N$ RBs allocated to PRACH, and $\kappa$ preambles constructed from per 6 RBs. Therefore, $N_p$ preambles can be constructed by $N_p = \kappa N/6$. Note that $N$ is integral multiple of 6. Suppose BS schedules one PUSCH for fixed-size data packets transmission and the number of RBs constituting one PUSCH is $\varsigma$, then the number of available PUSCH is $(Q-N)/\varsigma$.

According to the hybrid NORA-DT protocol, each cluster center MTCD randomly selects one preamble from the $N_p$ preambles and transmits it on PRACH. Let $B_i$ denote the number of cluster center MTCDs which select the same preamble $i$. Assuming the number of cluster center MTCDs
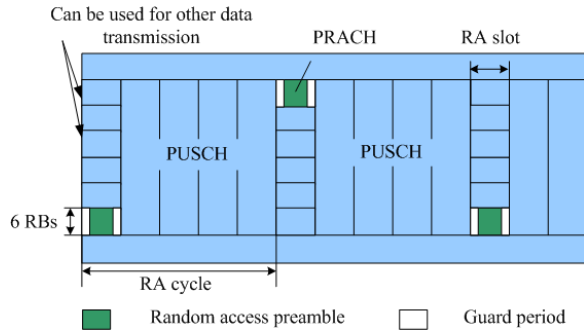
**FIGURE 4.** Periodic time-frequency resources for RA and data transmission.

selecting a preamble follows the Poisson distribution[1] with mean $u/N_p$, where $u$ is the number of NOMA clusters derived by device clustering algorithm. The probability that preamble $i$ is selected by $l$ cluster center MTCDs is

$$\Pr[B_i = l] = \left(\frac{6u}{\kappa N}\right)^l \exp\left(-\frac{6u}{\kappa N}\right) \Big/ l! \qquad (29)$$

In NORA-DT scheme, due to the enhanced preamble transmission, preamble $i$ is detected by BS if and only if one cluster center MTCD select the preamble $i$. Since the detection of different preambles is independent, the number of preambles that BS can detect follows a binomial distribution $V^{one} \sim Binom\left(\frac{\kappa N}{6}, Pr[B_i = 1]\right)$, where $V^{one}$ is the number of non-collision preambles, each of which is selected by only one cluster center MTCD in cycle. The expected number of non-collision preambles, denoted by $E[V^{one}]$, is given by

$$E[V^{one}] = \frac{\kappa N}{6} \Pr[B_i = 1] \qquad (30)$$

If BS schedules subchannel for the non-collision preamble $i$, the NOMA cluster conduct the non-orthogonal RA and data transmission successfully. According to the hybrid NORA-DT protocol, BS schedules enough data channels for the preambles that each is selected by only one cluster center MTCD. Upon receiving a matching RAR, all the MTCDs in the same NOMA cluster adjust the transmit power based on power back-off index and transmit data packets on the assigned subchannels. Therefore, the system throughput is given by $B_{succ} = M \min\left\{E[V^{one}], (Q-N)/\varsigma\right\}$, where $M$ is the number of multiplexing MTCDs in each NOMA cluster. If $E[V^{one}] \leq (Q-N)/\varsigma$, BS schedules respective data channels for the preambles that each is selected by only one cluster center MTCD. The number of MTCDs that successfully transmit data packets is $M$ times to the number of the preambles. On the other hand, if $E[V^{one}] > (Q-N)/\varsigma$,

BS schedules data channels only for $(Q-N)/\varsigma$ preambles. The number of MTCDs that successfully transmit data packets is $M$ times to the number of data channels. Let $N^*$ denote the number of PRACH RBs to minimize the gap between $E[V^{one}]$ and $(Q-N)/\varsigma$. $N^*$ is obtained by exhaustive attack algorithm as

$$N^* = \underset{6 \leq N < Q}{\arg\min} \left\{\left|E[V^{one}] - (Q-N)/\varsigma\right|\right\}$$

$$= \underset{6 \leq N < Q}{\arg\min} \left\{\left|u \exp\left(-\frac{6u}{\kappa N}\right) - (Q-N)/\varsigma\right|\right\} \quad (31)$$

Let $d = u \exp\left(-\frac{6\mu}{\kappa N^*}\right) - \frac{Q-N^*}{\varsigma}$.

1) If $d = 0$, we can get $E[V^{one}] = (Q-N)/\varsigma$, where $N = N^*$. Assuming $u$ is fixed, then $E[V^{one}]$ is convex function of $N$, and with the increment of $N$, $(Q-N)/\varsigma$ decreases. Therefore, when $N < N^*$, $B_{succ} = E[V^{one}]$. When $N > N^*$, $B_{succ} = (Q-N)/\varsigma$. Therefore, the throughput is given as

$$B_{succ} = \begin{cases} Mu \exp\left(-\frac{6u}{\kappa N}\right) & \text{if } N \leq N^* \\ \frac{M(Q-N)}{\varsigma} & \text{if } N \geq N^* \end{cases} \quad (32)$$

2) If $d > 0$, the throughput is given as

$$B_{succ} = \begin{cases} Mu \exp\left(-\frac{6u}{\kappa N}\right) & \text{if } N < N^* \\ \frac{M(Q-N)}{\varsigma} & \text{if } N \geq N^* \end{cases} \quad (33)$$

In this case, the optimal PRACH RBs is $N^* - 6$ or $N^*$ by comparing $u \exp\left(-\frac{6\mu}{\kappa(N^*-6)}\right)$ with $\frac{Q-N^*}{\varsigma}$, respectively.

3) If $d < 0$, the throughput is given as

$$B_{succ} = \begin{cases} Mu \exp\left(-\frac{6u}{\kappa N}\right) & \text{if } N \leq N^* \\ \frac{M(Q-N)}{\varsigma} & \text{if } N > N^* \end{cases} \quad (34)$$

In this case, the optimal PRACH RBs is $N^*$ or $N^* + 6$ by comparing $u \exp\left(-\frac{6\mu}{\kappa N^*}\right)$ with $\frac{Q-(N^*+6)}{\varsigma}$, respectively.

### B. ADAPTIVE RESOURCE ALLOCATION BASED ON DEVICE NUMBER INTERVALS

From above, it can be seen that with change of $u$, $N^*$ is recalculated. If $|d| > 0$, $u \exp\left(-\frac{6\mu}{\kappa N}\right)$ and $\frac{Q-N}{\varsigma}$ is further computed and compared when $N$ is $N^* - 6$ and $N^* + 6$, respectively. To distinguish this solution from the proposed resource allocation scheme in this paper, we call this solution as reference scheme. Although the reference scheme realizes a resource trade-off between PRACH and PUSCH, the procedure is complicated. In order to improve the system throughput and reduce the computation of resource allocation, a computationally efficient adaptive resource allocation scheme is proposed in this subsection.

First, the value of $u$ which satisfies $d = 0$ is derived as follows. Take both $u = x_m$ and $N = N^* = 6m$ into $u \exp\left(-\frac{6\mu}{\kappa N}\right) = \frac{Q-N}{\varsigma}$, then we have

$$x_m \exp\left(-\frac{x_m}{\kappa m}\right) = \frac{Q - 6m}{\varsigma} \qquad (35)$$

---

[1]If $\lim_{n \to \infty} np_n = \lambda > 0$, where $n$ is the total number of MTCD, $p_n$ is the communication probability of MTCD, $\lambda$ is the mean arrival rate, then $\lim_{n \to \infty} \binom{n}{k} p_n^k (1-p_n)^{n-k} = \frac{\lambda^k e^{-\lambda}}{k!}$, where $\binom{n}{k} = \frac{n!}{(k!) \times (n-k)!}$, and $\frac{\lambda^k e^{-\lambda}}{k!}$ is the probability density function of Poisson distribution. Basically, $n$ has to be very large and $p_n$ has to be appropriately small. This is called the rare events limit because it can be interpreted as applying to the number of occurrences of a rare event in a very large population of individual or trails.

$x_m$ is formulated as

$$x_m = -\kappa m \left\lfloor W_0\left(\frac{Q - 6m}{-\varsigma\kappa m}\right)\right\rfloor \tag{36}$$

where $W_0$ is the principle branch of the lambert W-function, $\lfloor\rfloor$ denotes the bottom integer function. When $\frac{Q}{\varsigma\kappa e^{-1}+6} \leq m < \frac{Q}{6}$, $W_0\left(\frac{Q-6m}{-\varsigma\kappa m}\right) < 0$ and $x_m > 0$ (see Appendix A for detail). Since $m$ is an integer, $\left\lfloor\frac{Q}{\varsigma\kappa e^{-1}+6}\right\rfloor \leq m < \left\lfloor\frac{Q}{6}\right\rfloor$. Therefore, when $n = x_m$, $d = 0$, $6m$ RBs are allocated to PRACH. Then the RBs allocation between PRACH and PUSCH for any number of active MTCDs can be derived in Algorithm 3.

1) If $u \in (x_{m+1}, x_m)$, where $\left\lfloor\frac{Q}{\varsigma\kappa e^{-1}+6}\right\rfloor \leq m < \left\lfloor\frac{Q}{6}\right\rfloor - 1$, $6m$ RBs are allocated to PRACH for the following reasons.

If $u \in \left(x_{m+1}, \frac{x_{m+1}+x_m}{2}\right)$, we can get $d > 0$, $N^* = 6(m+1)$ (see Appendix B for detail). According to the reference scheme, the optimal number of PRACH RBs is $6m$ or $6(m+1)$ by comparing $u \exp\left(-\frac{\mu}{\kappa m}\right)$ with $\frac{Q-6(m+1)}{\varsigma}$, respectively. If $u \in \left[\frac{x_{m+1}+x_m}{2}, x_m\right)$, we can get $d < 0$, $N^* = 6m$ (see Appendix B for detail). According to the reference scheme, the optimal number of PRACH RBs is $6m$ or $6(m+1)$ by comparing $u \exp\left(-\frac{\mu}{\kappa m}\right)$ with $\frac{Q-6(m+1)}{\varsigma}$, respectively. Therefore, if $u \in (x_{m+1}, x_m)$, the optimal number of PRACH RBs is $6m$ or $6(m+1)$ by comparing $u \exp\left(-\frac{\mu}{\kappa m}\right)$ with $\frac{Q-6(m+1)}{\varsigma}$, respectively. Note that $x_{m+1} \exp\left(-\frac{x_{m+1}}{\kappa(m+1)}\right) = \frac{Q-6(m+1)}{\varsigma}$, $x_m \exp\left(-\frac{x_m}{\kappa m}\right) = \frac{Q-6m}{\varsigma}$, then $x_{m+1} \exp\left(-\frac{x_{m+1}}{\kappa(m+1)}\right) < x_m \exp\left(-\frac{x_m}{\kappa m}\right)$. Since $x_m \exp\left(-\frac{x_m}{\kappa m}\right) < u \exp\left(-\frac{\mu}{\kappa m}\right)$, then $\frac{Q-6(m+1)}{\varsigma} < u \exp\left(-\frac{\mu}{\kappa m}\right)$. Therefore, if $u \in (x_{m+1}, x_m)$, where $\left\lfloor\frac{Q}{\varsigma\kappa e^{-1}+6}\right\rfloor \leq m < \left\lfloor\frac{Q}{6}\right\rfloor - 1$, the optimal number of PRACH RBs is $6m$.

2) If $u \in (0, x_{m+1})$, where $m = \left\lfloor\frac{Q}{6}\right\rfloor - 2$, $6(m+1)$ RBs would be allocated to PRACH. Since when $u = x_{m+1}$, $N^* = Q - 6$. With the decrease of $u$, $N^*$ is non-decreasing. Thereby, we have $N^* = Q-6$ and $d < 0$ when $u \in (0, x_{m+1})$. According to the reference scheme, the optimal number of PRACH RBs is $Q - 6$.

3) If $u > x_m$, where $m = \left\lfloor\frac{Q}{\varsigma\kappa e^{-1}+6}\right\rfloor$. First, $N^*$ is calculated by (31), $d = u \exp\left(-\frac{6\mu}{\kappa N^*}\right) - \frac{Q-N^*}{\varsigma}$. If $d < 0$, the optimal number of PRACH RBs is $N^*+6$ (see Appendix C for detail). If $d > 0$, the optimal number of PRACH RBs is $N^* - 6$ (see Appendix C for detail).

In the proposed device number intervals based resource allocation, frequently calculations can be avoided. BS first learns the resource allocation for some device number intervals, then the resource allocation for any number of MTCDs is derived from those intervals.

## V. PERFORMANCE EVALUATION

In this section, we present the performance of the proposed NORA-DT scheme in terms of the system throughput, overall

---

**Algorithm 3** Resource Allocation based on Device Number Intervals in NORA-DT

1: Initialize $Q$, $\varsigma$, $\kappa$, $u$, $m_a = \left\lfloor\frac{Q}{\varsigma\kappa e^{-1}+6}\right\rfloor$, $m_b = \left\lfloor\frac{Q}{6}\right\rfloor - 1$.
2: **for** $m = m_a$ to $m_b$ **do**
3:      $x_m = -\kappa m \left\lfloor W_0\left(\frac{Q-6m}{-\varsigma\kappa m}\right)\right\rfloor$.
4: **end for**
5: **if** $u \in (x_{m+1}, x_m]$, $\forall m \in [m_a, m_b]$ **then**
6:      $6m$ RBs are allocated to PRACH.
7:      $Q - 6m$ RBs are allocated to PUSCH.
8: **end if**
9: **if** $u \in \left(0, x_{m_b}\right]$ **then**
10:      $6m_b$ RBs are allocated to PRACH. $Q - 6m_b$ RBs are allocated to PUSCH.
11: **end if**
12: **if** $u \in \left(x_{m_a}, +\infty\right)$ **then**
13:      Calculate $N^*$ by (31). $d = u \exp\left(-\frac{6\mu}{\kappa N^*}\right) - \frac{Q-N^*}{\varsigma}$.
14:      **if** $d < 0$ **then**
15:          $N^*+6$ RBs are allocated to PRACH. $Q - (N^* + 6)$ RBs are allocated to PUSCH.
16:      **end if**
17:      **if** $d > 0$ **then**
18:          $N^*-6$ RBs are allocated to PRACH. $Q - (N^* - 6)$ RBs are allocated to PUSCH.
19:      **end if**
20: **end if**

**TABLE 1.** Simulation parameters.

| Parameter | Value |
|---|---|
| System bandwidth | 20MHz |
| Bandwidth of a RB, $B$ | 180KHz |
| Number of preambles mapped to one PRACH, $\kappa$ | 24 |
| Number of RBs constituting one PUSCH, $\varsigma$ | 1 |
| Maximum transmission power | $20 \sim 30$dBm |
| Minimum data rate | $1 \sim 2$bps |
| Distance between MTCDs and BS | $0.5 \sim 1$km |
| Number of multiplexing devices | 2 |
| Variance of additive noise | 0dB |

resource efficiency, PUSCH resource efficiency and average EE. The simulation parameters are given in Table I. Average EE can be quantitatively measured by the bits of information reliably transferred to a receiver per unit consumed energy per unit bandwidth at the transmitter. The overall resource efficiency is defined as the number of RBs for successful RA and data transmission over the total number of RBs for M2M communications. The overall resource efficiency can be written as:

$$r_{total} = \frac{B_{succ}\left(\frac{6}{\kappa} + \varsigma\right)}{Q} \tag{37}$$

The PUSCH resource efficiency is defined as ratio between the number of PUSCH for successful accesses and the total number of available PUSCH. The PUSCH resource efficiency can be written as:

$$r_{PUSCH} = \frac{\varsigma B_{succ}}{Q - N} \tag{38}$$

**TABLE 2.** Device number intervals for different evaluation cycle size.

| Evaluation cycle | $1ms$ | $1ms$ | $2ms$ | $5ms$ | $10ms$ |
|---|---|---|---|---|---|
| $Q$ | 60 | 80 | 160 | 400 | 800 |
| $m$ | [5,9] | [4,12] | [6,26] | [14,66] | [28,132] |
| $x_m$ | [6,43] | [8,84] | [4,284] | [4,948] | [8,1896] |

**TABLE 3.** Calculation of $m$, $x_m$, and $6m$ for $Q = 60$.

| $m$ | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|
| $x_m$ | 43 | 29 | 20 | 13 | 6 |
| $6m$ | 30 | 36 | 42 | 48 | 54 |

Table II shows the device number intervals corresponding to different evaluation cycle size. The evaluation cycle size is integer multiples of the duration of RA cycle. Assuming the duration of RA cycle is $1ms$. As the evaluation cycle size increases, the total number of RBs for M2M communications increases. $m$ is obtained from $\left\lfloor \frac{Q}{\varsigma \kappa e^{-1} + 6} \right\rfloor \leq m < \left\lfloor \frac{Q}{6} \right\rfloor$ and $x_m$ is obtained from (36). It can be seen that with the increase of $Q$, $m$ and $x_m$ increase. We can also observe that the number of device number intervals, i.e., $\max(m) - \min(m) + 1$, is 5, 9, 21, 53 and 105, respectively. As $Q$ increases, the number of device number intervals increases. Furthermore, if $Q$ is a fixed value in each evaluation cycle, BS does not have to calculate and reserve the RBs allocation in every cycle, which realizes simpler hardware implementation. While in the reference scheme and scheme in ORA and ORA-DT, the RBs allocation for any number of MTCDs needs to be frequently calculated by BS in every cycle.

Table III shows the device number intervals when $Q = 60$ for the evaluation cycle of 1 $ms$. It can be seen that for $m = 9, 8, 7, 6, 5$, $x_m$ is 6, 13, 20, 29 and 43, respectively. The device number intervals are (0, 6], (6, 13], (13, 20], (20, 29], (29, 43]. The optimal number of RBs allocated to PRACH for these device number intervals are 54, 48, 42, 36, 30, respectively. Then the RBs allocation for any number of MTCDs can be known from the RBs allocation for these device number intervals. For example, if the number of active MTCDs in current cycle is 56, then the number of RBs allocated to PRACH and PUSCH is 30 and 30, respectively. While if the number of active MTCDs in next cycle is 34, then the number of RBs allocated to PRACH and PUSCH is 36 and 24, respectively.

Fig.5 shows the number of PRACH RBs (i.e., $N^*$) to minimize the gap between the number of successful preamble transmission and the number of available PUSCH for different values of active MTCDs. The number of active MTCDs is denoted as $u$. We have $5 \leq m \leq 9$ when $Q = 60$ for the evaluation cycle of 1 $ms$. It can be seen that with the increment of $u$, $N^*$ is non-increasing. While with the increment of $x_m$, $N^*$ is decreased by 6 RBs. We can also observe that when $x_{m+1} < u < \frac{x_{m+1}+x_m}{2}$, $N^* = 6(m+1)$. While when $\frac{x_{m+1}+x_m}{2} \leq u < x_m$, $N^* = 6m$.

Fig.6 shows the comparison of the number of successful preamble transmission and the number of data channels
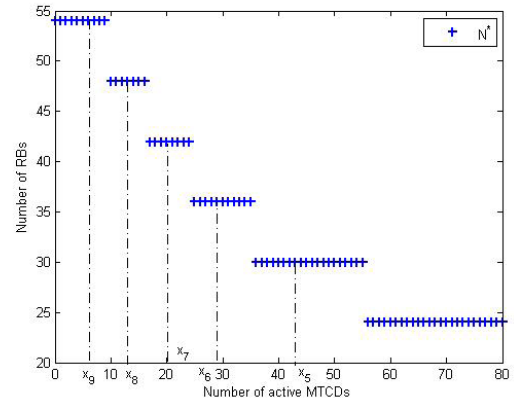


**FIGURE 5.** Comparison of the number of PRACH RBs (i.e., $N^*$) for different number of active MTCDs, $Q = 60$.
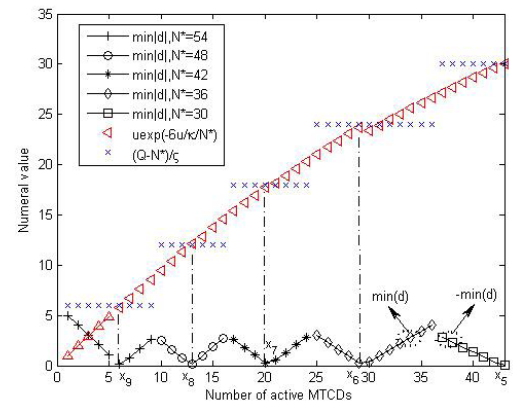


**FIGURE 6.** Comparison of preamble transmission and the number of data channels when the number of PRACH RBs is $N^*$.

when the number of PRACH RBs is $N^*$. $Q = 60$. In Fig.6, $u \exp\left(-\frac{6\mu}{\kappa N^*}\right)$ is the number of successful preamble transmission, and $\frac{Q-N^*}{\varsigma}$ is the number of data channels. $d = u \exp\left(-\frac{6\mu}{\kappa N^*}\right) - \frac{Q-N^*}{\varsigma}$. It can be seen that if $u \in \left(x_{m+1}, \frac{x_{m+1}+x_m}{2}\right)$, $\min\{d\} > 0$. With the increment of $u$, $|d|$ increases. If $u \in \left[\frac{x_{m+1}+x_m}{2}, x_m\right)$, $\min\{d\} < 0$. With the decrement of $u$, $|d|$ increases. We can also observe that $N^* = 6(m+1)$ when $u \in \left(x_{m+1}, \frac{x_{m+1}+x_m}{2}\right)$, $N^* = 6m$ when $u \in \left[\frac{x_{m+1}+x_m}{2}, x_m\right)$.

Fig.7 shows the comparison of the system throughput in reference scheme, the number of successful preamble transmission when $N^* = 6m$, and the number of data channels when $N^* = 6(m+1)$. $Q = 60$. In Fig.7, $u \exp\left(-\frac{\mu}{\kappa m}\right)$ is the number of successful preamble transmission when the number of RBs allocated to PRACH is $6m$. $\frac{Q-6(m+1)}{\varsigma}$ is the number of data channels when the number of RBs allocated to PRACH is $6(m+1)$. According to the analysis in reference scheme, when $u \in (x_{m+1}, x_m)$, the optimal number of PRACH RBs is $6m$ or $6(m+1)$ by comparing $u \exp\left(-\frac{\mu}{\kappa m}\right)$ with $\frac{Q-6(m+1)}{\varsigma}$, respectively. As Fig.7 shows, the value of
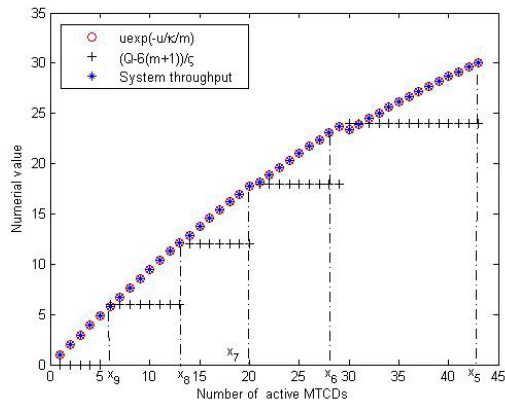
**FIGURE 7.** Comparison of system throughput in reference scheme, the number of successful preamble transmission when $N^* = 6m$, and the number of data channels when $N^* = 6(m+1)$.



**FIGURE 8.** The expected number of active MTCDs over 50 cycles.



**FIGURE 9.** Comparison of system throughput among ORA, ORA-DT and NORA-DT for different number of active MTCDs, $Q = 60$.

$u \exp\left(-\frac{\mu}{\kappa m}\right)$ is larger than $\frac{Q-6(m+1)}{\varsigma}$ when $u \in (x_{m+1}, x_m)$. It means that the system throughput is maximum when $6m$ RBs are allocated to PRACH, which is conformed with the analysis in the proposed resource allocation scheme.

Fig.8 shows the expected number of active MTCDs over 50 cycles. Let $Q = 60$. As Table III shows, the device number intervals are (0, 6], (6, 13], (13, 20], (20, 29], (29, 43]. The optimal number of RBs allocated to PRACH for these device number intervals are 54, 48, 42, 36, 30, respectively. Other than reattempting MTCDs which failed in previous RA cycles, we assume new active MTCDs arrive at the system for transmitting data packets, and the number of new arriving MTCDs follows a Poisson distribution with mean $\lambda$. We define the number of active MTCDs in cycle $t = 1$ is $u$, and the new active MTCDs arrive in cycle $t \geq 2$. 1) Let $u = 40$, $\lambda=20$. As Fig.8 shows, when $t \geq 3$, the optimal number of PRACH RBs is 36, while the optimal number of PRACH RBs before cycle $t = 3$ is 30. 2) Let $u = 60$, $\lambda=30$. As Fig.8 shows, when $t \geq 7$, the optimal number of PRACH RBs is 30, while the optimal number of PRACH RBs before cycle $t = 7$ is obtained by (31). 3) Let $u = 60$, $\lambda=10$. As Fig.8 shows, when $t \geq 5$, the optimal number of PRACH RBs is 48. When $t = 4$, the optimal number of PRACH RBs is 42. When $t = 3$, the optimal number of PRACH RBs is 36. While the optimal number of PRACH RBs before cycle $t = 3$ is obtained by (31). From 1)-3), it can be seen that the optimal number of PRACH RBs is relative to the device number intervals but not to $u$, $\lambda$ and cycle $t$. Therefore, upon the RBs allocation is determined and reserved by BS, BS does not have to calculate the RBs allocation in every cycle, the RBs allocation for any cycle can be fast known from BS.

Fig.9 shows the system throughput of ORA, ORA-DT and NORA-DT for different number of active MTCDs. Theoretical analysis and simulation results are also compared. It can be seen that the system throughput of ORA first increases to its maximum and then drops greatly due to serious preamble collisions and unreasonable resource allocation between PRACH and PUSCH. The system throughput in ORA-DT first increases to its maximum and then achieve a stable level
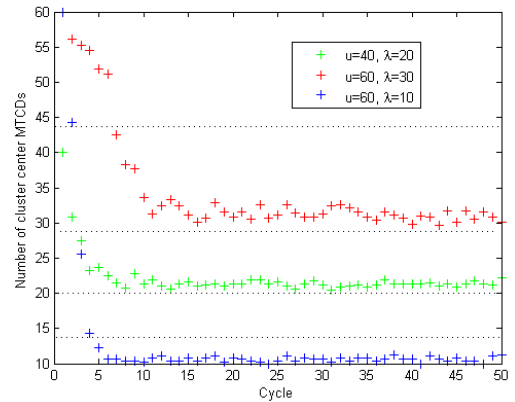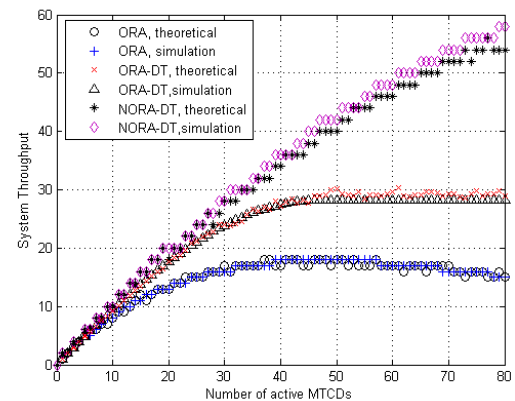
as the number of active MTCD increases. This is because the ACB scheme is incorporated to control the number of active MTCDs when the number of active MTCDs is higher than the optimal number of participating MTCDs, and the number of RBs allocated to PRACH and PUSCH is set to the same. In contrast, for the proposed NORA-DT procedure, the system throughput is obviously increased even when the number of active MTCDs is much higher. This is because with the increase of the number of MTCDs, the minimum data rate requirements and channel conditions of different MTCDs are more abundant. NOMA can effectively improve the connections by supporting more and more NOMA clusters.

In subsection IV-A, we assume the number of cluster center MTCDs selecting a preamble follows the Poisson distribution, the throughput function are derived on this assumption. To show that the Poisson assumption is acceptable, the system throughput of NORA-DT for different number of active MTCDs by analysis (i.e., the throughput function in section VI) and by simulation is presented in Fig.9. Similarly, the system throughput of ORA and ORA-DT by analysis and by simulation is also provided. We present an example to illustrate how to obtain the simulation results. If there are $\kappa N/6$ preambles are assigned for M2M communications, and $u$ MTCDs randomly chooses a preamble out of these
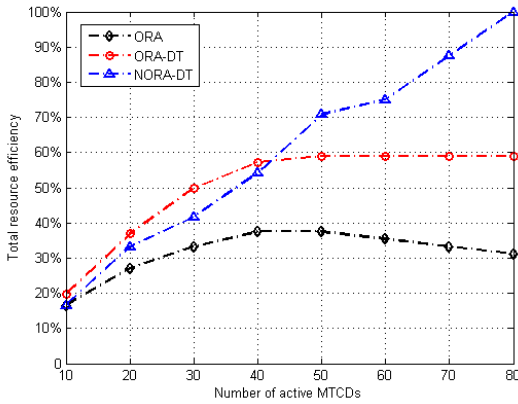
**FIGURE 10.** Comparison of total resource efficiency among ORA, ORA-DT and NORA-DT, $Q = 60$.



**FIGURE 11.** Comparison of PUSCH resource efficiency among ORA, ORA-DT and NORA-DT, $Q = 60$.

preambles. Denote $i$ as index for the $i$-th preamble, where $i \in \left\{ 1, \cdots, \kappa N / 6 \right\}$, and $B_i$ as the number of MTCDs that select the preamble $i$. We can categorize each preamble $i$ into the following three cases. 1) Idle preamble: $B_i = 0$, preamble $i$ is not selected by any MTCD. 2) Collision preamble: $B_i \geq 2$, preamble $i$ is selected by more than one MTCD. 3) Non-collision preamble: $B_i = 1$, preamble $i$ is selected by only one MTCD. For each preamble $i$, if $B_i = 1$, let the number of non-collision preambles increases by one until all preambles are exhausted. Then the number of MTCDs that successfully transmit data packets is obtained by $\min \left\{ n_{non-coll}, (Q-N)/\varsigma \right\}$, where $n_{non-coll}$ is the number of all non-collision preambles, and $(Q-N)/\varsigma$ is the number of data channels. $N$ is the number of RBs for the PRACH and the uplink data channels which can be obtained in advance by the proposed resource allocation scheme. It can be seen the system throughput by analysis agrees with the actual throughput obtained by simulation.

Fig.10 compares the total resource efficiency in (37) among ORA, ORA-DT and NORA-DT for different number of active MTCDs. The overall resource efficiency for the traditional ORA is very low. When $u = 50$, the overall resource efficiency of ORA-DT and NORA-DT can reach 59% and 71%, respectively. The overall resource efficiency of NORA-DT is about 20% higher than that of ORA-DT. When $u = 70$, the overall resource efficiency of ORA-DT is still 59%. However, the overall resource efficiency of NORA-DT can reach 89%. The overall resource efficiency of NORA-DT is about 50% higher than that of ORA-DT. It can be concluded that as the number of active MTCDs increases, the overall resource efficiency is obviously improved in NORA-DT.

Fig.11 compares the PUSCH resource efficiency in (38) among ORA, ORA-DT and NORA-DT for different number of active MTCDs. It is worth observing that in traditional ORA, the PUSCH resource efficiency is the same as overall resource efficiency. This is because the number of available preambles is the same with the number of available PUSCH, i.e., $\frac{\kappa N}{6} = \frac{Q-N}{\varsigma}$, then $N = \left\lfloor \frac{Q}{1+\varsigma\kappa/6} \right\rfloor$. It can be seen that $N$ is a fixed value according to different number of active
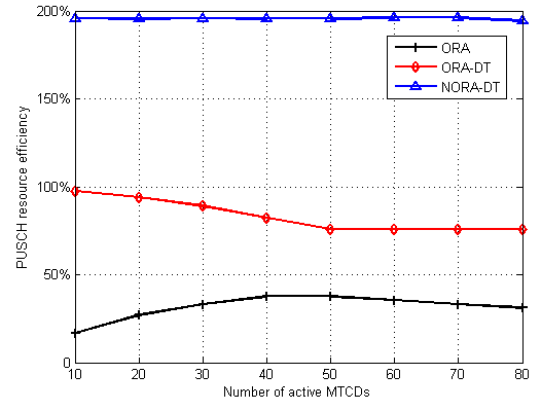
MTCDs. The PUSCH resource efficiency in ORA-DT first decreases to its minimum. This is because as the number of active MTCDs increases, the probability of preamble collisions become larger. Since BS cannot detect the preamble collisions before scheduling PUSCH, more and more RBs are wasted on the collided PUSCH. And then the PUSCH resource efficiency in ORA-DT achieve a stable level as the number of active MTCD increases. This is because when the number of active MTCDs is higher than the optimal number of participating MTCDs, the system throughput in ORA-DT is constant, as is shown in Fig.9. Besides, the RBs allocation between PRACH and PUSCH is the same as that when the number of active MTCDs equals the optimal number of participating MTCDs. In contrast, the proposed NORA-DT maintains a value of almost 200% even when the probability of preamble collisions increases. There are two reasons, one is due to the BS can perfectly detect the preamble collisions in advance and schedules PUSCH only to the NOMA clusters without collision, the other is that compared with the ORA and ORA-DT, two multiplexing MTCDs can transmit data packets simultaneously on the same PUSCH in NORA-DT.

Fig. 12 shows the average EE among ORA, ORA-DT and NORA-DT corresponding to different cluster center MTCD's target arrived SNR for $\rho = 5$ dB, $u = 80$. As seen in the figures, the cluster center MTCD's target arrived SNR ranges from 0dB to 20dB, with a interval of 2dB. It can be found that as the cluster center MTCD's target arrived SNR increases, NORA-DT scheme has a significant improvement in energy efficiency compared with ORA and ORA-DT. When the cluster center MTCD' target arrived SNR is relatively small, ORA and ORA-DT has a higher energy efficiency than that of NORA-DT scheme. This is because the value of target arrived SNR is directly affected by the cluster center MTCD's target arrived power (i.e., $p_{k,1}$). The target arrived SNR decreases as $p_{k,1}$ decreases. When $p_{k,1}$ is small, ORA and ORA-DT with only one MTCD sending data on the same channel resource is more likely to meet the devices's minimum data rate requirement and maximum transmission power. For NORA-DT with multiple devices sending data on the same channel resource,
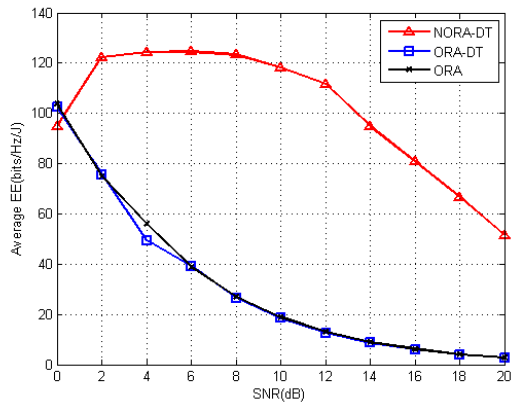
**FIGURE 12.** Comparison of average EE among ORA, ORA-DT and NORA-DT, $\rho = 5$ dB, $u = 80$.
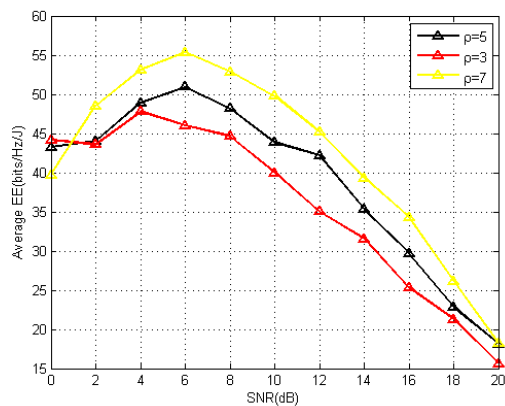


**FIGURE 14.** Comparison of average EE in NORA-DT versus cluster center MTCD's target arrived SNR for different number of active MTCDs, $\rho = 5$ dB.



**FIGURE 13.** Comparison of average EE in NORA-DT versus cluster center MTCD's target arrived SNR for $\rho = 3$, 5 and 7dB, $u = 20$.

$\rho = 5$ dB. As seen in this figure, for different number of active MTCDs, the average EE first increases to its maximum and then obviously decreases. The reason for the increase of average EE is that with the increase of the target arrived SNR of cluster center MTCD, the accuracy of SIC in decoding multiple MTCDs increases. The reason for the decrease of average EE is that with the increase of the target arrived SNR of cluster center MTCD, devices clustering and power allocation strategy are more difficult to meet the minimum data rate requirements and the maximum transmission power of all multiplexing MTCDs.

## VI. CONCLUSION

In this work, we have proposed a new hybrid NORA-DT scheme for M2M communications to resolve the excessive signalling overhead and resource allocation problems. A power back-off scheme based uplink NOMA transmission is introduced to adjust device's target arrived power, and the difference of transmission power among MTCDs is formulated. Based on the range of transmission power derived under the maximum transmission power constraints and minimum rate requirements at the MTCDs, the MTCDs are clustered into a set of NOMA clusters. A new hybrid NORA-DT protocol has been proposed to reduce excessive signalling overhead and improve resource efficiency, in which the cluster cluster center MTCD transmits a extended preamble on behalf of the MTCDs in a cluster on the PRACH for connection request. The BS can perfectly detect the preamble collisions in advance and schedules PUSCH only to the NOMA clusters without collision. Then the MTCDs in the same NOMA clusters transmit data packets right after preamble transmission on the PUSCH. An energy-efficient power allocation as an optimization problem of the cluster center MTCD's transmission power is formulated while providing devices' QoS guarantees. DC programming is used to transform the original non-convex problem for energy efficiency maximization to convex optimization problem. Then the transmission power of the cluster center MTCD is obtained by an iterative algorithm. By the relation with

it is difficult to meet the minimum data rate requirements and maximum transmission power of all multiplexing devices.

Fig.13 shows the average EE in NORA-DT corresponding to different cluster center MTCD's target arrived SNR for $\rho = 3$, 5 and 7dB. $u = 20$. As seen in this figure, the average EE first increases to its maximum and then obviously decreases as the cluster center MTCD's target arrived SNR increases. And we also observe that when the value of cluster center MTCD's target arrived SNR is less than 1dB, by increasing $\rho$ the performance of the average EE decreases. This is because for NORA-DT with multiple devices sending data on the same channel resource, it is difficult to meet the minimum data rate requirements and maximum transmission power of all multiplexing devices. When the value of cluster center MTCD's target arrived SNR is more than 1dB, by increasing $\rho$ the performance of the average EE increases. The reason is that for uplink, the power allocation strategy needs to ensure the difference of received power at BS among multiplex devices. As the difference (i.e., $\rho$) increases, the accuracy of SIC in decoding multiple MTCDs increases.

Fig.14 shows the average EE in NORA-DT versus different values of cluster center MTCD's target arrived SNR when the number of active MTCDs is 20, 40, 60 and 80, respectively.
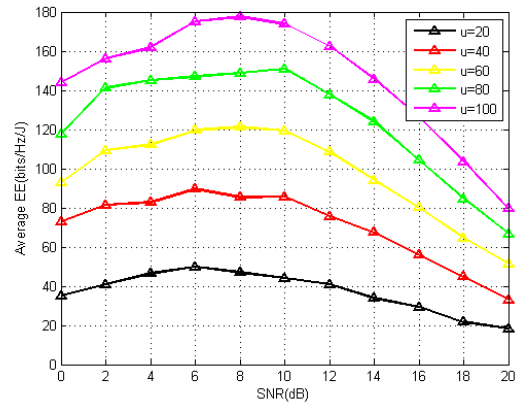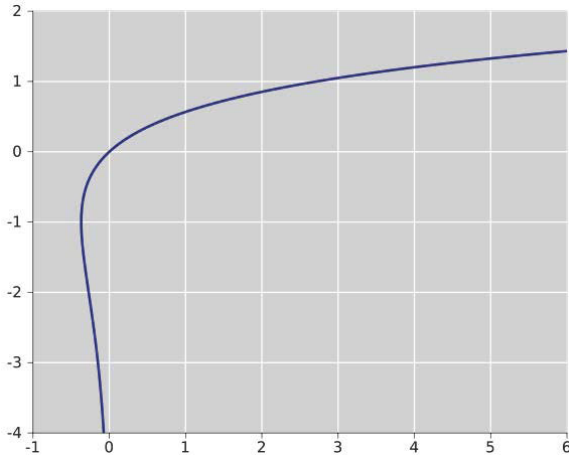
**FIGURE 15.** The principle branch of the lambertW-function.

the cluster center MTCD, the transmission power of other MTCDs in a cluster can be also obtained. A computationally efficient adaptive resource allocation scheme has been proposed to maximize the M2M throughput and resource efficiency to resolve the congestion problem both in the PRACH and PUSCH. We have derived a closed-form analytic expression for the expected throughput and have obtained the optimal number of RBs allocated to PRACH and PUSCH for some device number intervals in advance, which avoids frequent computation. Simulation results show that the proposed NORA-DT scheme outperforms the conventional ORA and ORA-DT schemes in terms of the number of successful data packet transmissions, while guaranteeing good system's resource efficiency and energy efficiency performance.

## APPENDIXES
### APPENDIX A

Let $f(z) = ze^z$, $z = f^{-1}(ze^z) = W(ze^z)$. As Fig.15 shows, the relation $W$ is multivalued (except at 0) due to the function $f$ is not injective. If we restrict attention to real-valued $W$, the complex variable $z$ is then replaced by the real variable $x$, and the relation is defined only for $x \geq -1/e$, and is double-valued on $(-1/e, 0)$. The additional constraint $W \geq -1$ defines a single-valued function $W_0(x)$. We have $W_0(0) = 0$ and $W_0(-1/e) = -1$. Therefore, $-e^{-1} \leq \frac{Q(t)-6m}{-\varsigma\kappa m} < 0$, then $\frac{Q}{\varsigma\kappa e^{-1}+6} \leq m < \frac{Q}{6}$.

### APPENDIX B

1) Assuming $6(m+1)$ RBs are allocated to PRACH when $u \in (x_{m+1}, x_m)$, $\forall m \in \left[\left\lfloor\frac{Q}{\varsigma\kappa e^{-1}+6}\right\rfloor, \left\lfloor\frac{Q}{6}\right\rfloor - 1\right)$, we have $d = u\exp\left(-\frac{\mu}{\kappa(m+1)}\right) - \frac{Q-6(m+1)}{\varsigma}$. Due to $u\exp\left(-\frac{\mu}{\kappa(m+1)}\right) > x_{m+1}\exp\left(-\frac{x_{m+1}}{\kappa(m+1)}\right)$, and $x_{m+1}\exp\left(-\frac{x_{m+1}}{\kappa(m+1)}\right) = \frac{Q-6(m+1)}{\varsigma}$, we have $d > 0$. With the increment of $u$, $|d|$ increases.

2) Assuming $6m$ RBs are allocated to PRACH when $u \in (x_{m+1}, x_m)$, $\forall m \in \left[\left\lfloor\frac{Q}{\varsigma\kappa e^{-1}+6}\right\rfloor, \left\lfloor\frac{Q}{6}\right\rfloor - 1\right)$, we have $d = u\exp\left(-\frac{\mu}{\kappa m}\right) - \frac{Q-6m}{\varsigma}$. Due to $u\exp\left(-\frac{\mu}{\kappa m}\right) < x_m\exp\left(-\frac{x_m}{\kappa m}\right)$,

and $x_m\exp\left(-\frac{x_m}{\kappa m}\right) = \frac{Q-6m}{\varsigma}$, we have $d < 0$. With the decrement of $u$, $|d|$ increases.

From above, in order to make $|d|$ as small as possible, the number of RBs allocated to PRACH when $u \in \left(x_{m+1}, \frac{x_{m+1}+x_m}{2}\right)$ is $6(m+1)$, while the number of RBs allocated to PRACH when $u \in \left[\frac{x_{m+1}+x_m}{2}, x_m\right)$ is $6m$.

### APPENDIX C

In subsection IV-B, we have derived only if $u = x_m$, $|d| = 0$, where $\left\lfloor\frac{Q}{\varsigma\kappa e^{-1}+6}\right\rfloor \leq m \leq \left\lfloor\frac{Q}{6}\right\rfloor - 1$. Therefore, if $u > x_\tau$, where $\tau = \left\lfloor\frac{Q}{\varsigma\kappa e^{-1}+6}\right\rfloor$, we have $|d| > 0$. $6(\tau-n)$ RBs are allocated to PRACH when $u > x_\tau$, where $n \geq 0$. This is because with the increment of $u$, the number of PRACH RBs in (31) is non-increasing. Consider the following cases: $d < 0$ and $d > 0$.

1) If $d < 0$, $u\exp\left(-\frac{\mu}{\kappa(\tau-n)}\right) < \frac{Q-6(\tau-n)}{\varsigma}$. According to subsection IV-A, $u\exp\left(-\frac{\mu}{\kappa(\tau-n)}\right)$ and $\frac{Q-6(\tau-n+1)}{\varsigma}$ is compared.

Since $u\exp\left(-\frac{\mu}{\kappa(\tau-n)}\right) < u\exp\left(-\frac{\mu}{\kappa(\tau-n+1)}\right)$, and $u\exp\left(-\frac{\mu}{\kappa(\tau-n+1)}\right) < \frac{Q-6(\tau-n+1)}{\varsigma}$, we have $u\exp\left(-\frac{\mu}{\kappa(\tau-n)}\right) < \frac{Q-6(\tau-n+1)}{\varsigma}$. Therefore, if $d < 0$, $6(\tau-n+1)$ RBs are allocated to PRACH.

2) If $d > 0$, $u\exp\left(-\frac{\mu}{\kappa(\tau-n)}\right) > \frac{Q-6(\tau-n)}{\varsigma}$. According to subsection IV-A, $u\exp\left(-\frac{\mu}{\kappa(\tau-n-1)}\right)$ and $\frac{Q-6(\tau-n)}{\varsigma}$ is compared.

Since $u\exp\left(-\frac{\mu}{\kappa(\tau-n-1)}\right) > \frac{Q-6(\tau-n-1)}{\varsigma}$, and $\frac{Q-6(\tau-n-1)}{\varsigma} > \frac{Q-6(\tau-n)}{\varsigma}$, we have $u\exp\left(-\frac{\mu}{\kappa(\tau-n-1)}\right) > \frac{Q-6(\tau-n)}{\varsigma}$. Therefore, if $d > 0$, $6(\tau-n-1)$ RBs are allocated to PRACH.

## REFERENCES

[1] Z. Wu, K. Lu, C. Jiang, and X. Shao, "Comprehensive study and comparison on 5G NOMA schemes," *IEEE Access*, vol. 6, pp. 18511–18519, 2018.

[2] N. Xia, H.-H. Chen, and C.-S. Yang, "Radio resource management in machine-to-machine communications—A survey," *IEEE Commun. Surveys Tuts.*, vol. 20, no. 1, pp. 791–828, 1st Quart., 2018.

[3] J. W. Raymond, T. O. Olwal, and A. M. Kurien, "Cooperative communications in machine to machine (M2M): Solutions, challenges and future work," *IEEE Access*, vol. 6, pp. 9750–9766, 2018.

[4] C. X. Mavromoustakis, G. Mastorakis, and J. M. Batalla, "A mobile edge computing model enabling efficient computation offload-aware energy conservation," *IEEE Access*, vol. 7, pp. 102295–102303, 2019.

[5] F. Ding, R. Su, E. Tong, D. Zhang, H. Zhu, and M. W. M. Ismail, "Toward a M2M-based Internet of vehicles framework for wireless monitoring applications," *IEEE Access*, vol. 6, pp. 67699–67708, 2018.

[6] Z. Zhou, Y. Guo, Y. He, X. Zhao, and W. M. Bazzi, "Access control and resource allocation for m2m communications in industrial automation," *IEEE Trans. Ind. Informat.*, vol. 15, no. 5, pp. 3093–3103, May 2019.

[7] *Cisco Visual Networking Index (VNI) Complete Forecast for 2017–2022*, Cisco, San Jose, CA, USA, 2018.

[8] A. Laya, L. Alonso, and J. Alonso-Zarate, "Is the random access channel of LTE and LTE—A suitable for M2M communications? A survey of alternatives," *IEEE Commun. Surveys Tuts.*, vol. 16, no. 1, pp. 4–16, 1st Quart., 2014.

[9] H. S. Jang, S. M. Kim, H.-S. Park, and D. K. Sung, "A preamble collision resolution scheme via tagged preambles for cellular IoT/M2M communications," *IEEE Trans. Veh. Technol.*, vol. 67, no. 2, pp. 1825–1829, Feb. 2018.

[10] C. Wan and J. Sun, "Access class barring parameter adaptation based on load estimation model for mMTC in LTE-A," in *Proc. Int. Conf. Commun., Inf. Syst. Comput. Eng. (CISCE)*, 2014, vol. 16, no. 1, pp. 4–16.

[11] K. Lee and J. W. Jang, "An efficient contention resolution scheme for massive IoT devices in random access to LTE—A networks," *IEEE Access*, vol. 6, pp. 67118–67130, 2018.

[12] Z. Alavikia and A. Ghasemi, "Collision-aware resource access scheme for LTE-based machine-to-machine communications," *IEEE Trans. Veh. Technol.*, vol. 67, no. 5, pp. 4683–4688, May 2018.

[13] L. Zhen, H. Qin, Q. Zhang, Z. Chu, G. Lu, J. Jiang, and M. Guizani, "Optimal preamble design in spatial group-based random access for satellite-M2M communications," *IEEE Wireless Commun. Lett.*, vol. 8, no. 3, pp. 953–956, Jun. 2019.

[14] H. Althumali and M. Othman, "A survey of random access control techniques for machine-to-machine communications in LTE/LTE—A networks," *IEEE Access*, vol. 6, pp. 74961–74983, 2018.

[15] H. Li and H. Liu, "An analysis of uplink OFDMA optimality," *IEEE Trans. Wireless Commun.*, vol. 6, no. 8, pp. 2972–2983, Aug. 2007.

[16] S. M. R. Islam, N. Avazov, O. A. Dobre, and K.-S. Kwak, "Power-domain non-orthogonal multiple access (NOMA) in 5G systems: Potentials and challenges," *IEEE Commun. Surveys Tuts.*, vol. 19, no. 2, pp. 721–742, 2nd Quart., 2017.

[17] Y. Gao, B. Xia, K. Xiao, Z. Chen, X. Li, and S. Zhang, "Theoretical analysis of the dynamic decode ordering sic receiver for uplink NOMA systems," *IEEE Commun. Lett.*, vol. 21, no. 10, pp. 2246–2249, Oct. 2017.

[18] Z. Ding, Y. Liu, J. Choi, Q. Sun, M. Elkashlan, 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.

[19] M. Shirvanimoghaddam, M. Dohler, and S. J. Johnson, "Massive non-orthogonal multiple access for cellular IoT: Potentials and limitations," *IEEE Commun. Mag.*, vol. 55, no. 9, pp. 55–61, Sep. 2017.

[20] 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.

[21] J. Choi, "On power and rate allocation for coded uplink NOMA in a multicarrier system," *IEEE Trans. Commun.*, vol. 66, no. 6, pp. 2762–2772, Jun. 2018.

[22] A. E. Mostafa, Y. Zhou, and V. W. S. Wong, "Connectivity maximization for narrowband IoT systems with NOMA," in *Proc. Int. Conf. Commun. (ICC)*, May 2017, pp. 1–6.

[23] Z. Tan, H. Qu, J. Zhao, G. Ren, and W. Wang, "Self-sustainable dense cellular M2M system with hybrid energy harvesting and high sensitivity rectenna," *IEEE Access*, vol. 7, pp. 19447–19460, 2019.

[24] J. Liu, G. Wu, S. Xiao, X. Zhou, G. Y. Li, S. Guo, and S. Li, "Joint power allocation and user scheduling for device-to-device-enabled heterogeneous networks with non-orthogonal multiple access," *IEEE Access*, vol. 7, pp. 62657–62671, 2019.

[25] D. Tweed, M. Derakhshani, S. Parsaeefard, and T. Le-Ngoc, "Outage-constrained resource allocation in uplink NOMA for critical applications," *IEEE Access*, vol. 5, pp. 27636–27648, 2017.

[26] N. Zhang, J. Wang, G. Kang, and Y. Liu, "Uplink nonorthogonal multiple access in 5G systems," *IEEE Commun. Lett.*, vol. 20, no. 3, pp. 458–461, Mar. 2016.

[27] M. R. G. Aghdam, R. Abdolee, F. A. Azhiri, and B. M. Tazehkand, "Random user pairing in massive-MIMO-NOMA transmission systems based on mmWave," in *Proc. IEEE 88th Veh. Technol. Conf. (VTC-Fall)*, Aug. 2018, pp. 1–6, doi: 10.1109/vtcfall.2018.8690578.

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

[29] M. S. Ali, H. Tabassum, and E. Hossain, "Dynamic user clustering and power allocation for uplink and downlink non-orthogonal multiple access (NOMA) systems," *IEEE Access*, vol. 4, pp. 6325–6343, 2016.

[30] E. Cabrera and R. Vesilo, "An enhanced k-means clustering algorithm with non-orthogonal multiple access (NOMA) for MMC networks," in *Proc. 28th Int. Telecommun. Netw. Appl. Conf. (ITNAC)*, Nov. 2018, pp. 1–8, doi: 10.1109/atnac.2018.8615298.

[31] M. R. G. Aghdam, S. M. Pishvaei, R. Abdolee, B. M. Tazehkand, and F. T. Miandoab, "User grouping and optimal random beamforming in mmWave MIMO-NOMA transmission systems," in *Proc. IEEE 20th Int. Symp. 'World Wireless, Mobile Multimedia Netw.' (WoWMoM)*, Jun. 2019, pp. 1–6, doi: 10.1109/wowmom.2019.8793042.

[32] M. Pischella and D. Le Ruyet, "NOMA-relevant clustering and resource allocation for proportional fair uplink communications," *IEEE Wireless Commun. Lett.*, vol. 8, no. 3, pp. 873–876, Jun. 2019.

[33] M.-R. Hojeij, C. Abdel Nour, J. Farah, and C. Douillard, "Waterfilling-based proportional fairness scheduler for downlink non-orthogonal multiple access," *IEEE Wireless Commun. Lett.*, vol. 6, no. 2, pp. 230–233, Apr. 2017.

[34] Z. Yang, W. Xu, H. Xu, J. Shi, and M. Chen, "Energy efficient non-orthogonal multiple access for machine-to-machine communications," *IEEE Commun. Lett.*, vol. 21, no. 4, pp. 817–820, Apr. 2017.

[35] M. Zeng, W. Hao, O. A. Dobre, and H. V. Poor, "Energy-efficient power allocation in uplink mmWave massive MIMO with NOMA," *IEEE Trans. Veh. Technol.*, vol. 68, no. 3, pp. 3000–3004, Mar. 2019.

[36] B. Gu, C. Zhang, H. Wang, Y. Yao, and X. Tan, "Power control for cognitive m2m communications underlaying cellular with fairness concerns," *IEEE Access*, vol. 7, pp. 80789–80799, 2019.

[37] Z. Yang, W. Xu, Y. Pan, C. Pan, and M. Chen, "Energy efficient resource allocation in machine-to-machine communications with multiple access and energy harvesting for IoT," *IEEE Internet Things J.*, vol. 5, no. 1, pp. 229–245, Feb. 2018.

[38] Z. Li and J. Gui, "Energy-efficient resource allocation with hybrid TDMA–NOMA for cellular-enabled machine-to-machine communications," *IEEE Access*, vol. 7, pp. 105800–105815, 2019.

[39] A. B. Rozario and M. F. Hossain, "An architecture for M2M communications over cellular networks using clustering and hybrid TDMA-NOMA," in *Proc. 6th Int. Conf. Inf. Commun. Technol. (ICoICT)*, May 2018, pp. 18–23, doi: 10.1109/icoict.2018.8528767.

[40] S. Alemaishat, O. A. Saraereh, I. Khan, and B. J. Choi, "An efficient resource allocation algorithm for D2D communications based on NOMA," *IEEE Access*, vol. 7, pp. 120238–120247, 2019.

[41] Z. Na, M. Zhang, M. Jia, M. Xiong, and Z. Gao, "Joint uplink and downlink resource allocation for the Internet of Things," *IEEE Access*, vol. 7, pp. 15758–15766, 2019.

[42] Y. Liang, X. Li, J. Zhang, and Z. Ding, "Non-orthogonal random access for 5G networks," *IEEE Trans. Wireless Commun.*, vol. 16, no. 7, pp. 4817–4831, Jul. 2017.

[43] D. T. Wiriaatmadja and K. W. Choi, "Hybrid random access and data transmission protocol for machine-to-machine communications in cellular networks," *IEEE Trans. Wireless Commun.*, vol. 14, no. 1, pp. 33–46, Jan. 2015.

[44] G. Roche, A. Glazunov, and B. Allen, *LTE-Advanced and Next Generation Wireless Networks Channel Modelling and Propagation*. London, U.K.: Wiley, 2013.

[45] *Physical Layer Procedures, V10.12.0*, document TR 36.213, 3GPP, 2015.

[46] W. Han, Y. Zhang, X. Wang, J. Li, M. Sheng, and X. Ma, "Orthogonal power division multiple access: A green communication perspective," *IEEE J. Sel. Areas Commun.*, vol. 34, no. 12, pp. 3828–3842, Dec. 2016.

[47] D. Axehill, F. Gunnarsson, and A. Hansson, "A low-complexity high-performance preprocessing algorithm for multiuser detection using gold sequences," *IEEE Trans. Signal Process.*, vol. 56, no. 9, pp. 4377–4385, Sep. 2008.

[48] 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.

[49] M. Grant and S. Boyd. (Jun. 2015). *CVX: MATLAB Software for Disciplined Convex Programming, Version 2.1.* [Online]. Available: http://cvxr.com/cvx/

[50] H. Kha, H. D. Tuan, and H. H. Nguyen, "Fast global optimal power allocation in wireless networks by local D.C. programming," *IEEE Trans. Wireless Commun.*, vol. 11, no. 2, pp. 510–515, Feb. 2012.

[51] *Evolved Universal Terrestrial Radio Access (E-UTRA); Physical Channels Modulation, V10.4.0*, document TS 36.211, 3GPP, 2011.
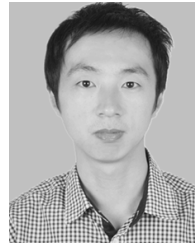
**YALI WU** received the Ph.D. degree in communication and information systems from the Beijing University of Posts and Telecommunications (BUPT), in 2017. She is currently a Lecturer with the Hebei University of Economics and Business. Her research interests include wireless communications and machine to machine communications.

**NINGBO ZHANG** received the Ph.D. degree from the Beijing University of Posts and Telecommunications (BUPT), in 2010. He is currently an Associate Professor with BUPT. From 2010 to 2014, he was a Senior Engineer with the Research and Development Wireless Department, Huawei Technologies. Since 2014, he has been the Project Manager of several national projects such as the National Natural Science Foundation of China, the National Science and Technology Major Project of China, and the National 863 Project. He has authored or coauthored over 40 journal and conference papers. He has expertise in the physical layer of 5G wireless systems and the machine-to-machine communications in the IoT.

**KAIXUAN RONG** received the Ph.D. degree in electronic circuit and systems from Xidian University, in 2016. He is currently a Senior Engineer with The 54th Research Institute of China Electronics Technology Group Corporation. His major research interests include joint channel estimation and joint detection, multicarrier transmission, channel coding, and multiple access.

• • •