

Design and Evaluation of Buffer-Aided Cooperative NOMA with Direct Transmission in IoT

Peng Xu, *Member, IEEE*, Yunwu Wang, Gaojie Chen, *Senior Member, IEEE*,
Gaofeng, Pan, *Senior Member, IEEE* and Zhiguo Ding, *Fellow, IEEE*

Abstract—The high spectrum efficiency of non-orthogonal multiple access (NOMA) is attractive to solve the massive number of connections in the Internet of Things (IoT). This paper investigates a buffer-aided cooperative NOMA (C-NOMA) system in the IoT, where the intended users are equipped with buffers for cooperation. The direct transmission from the access point to the users and the buffer-aided cooperative transmission between the intended users are coordinated. In particular, a novel buffer-aided C-NOMA scheme is proposed to adaptively select a direct or cooperative transmission mode, based on the instantaneous channel state information and the buffer state. Then, the performance of the proposed scheme, in terms of the system outage probability and average delay, is theoretically derived with closed-form expressions. Furthermore, the full diversity order of three is demonstrated to be achieved for each user pair if the buffer size is not less than three, which is larger than conventional non-buffer-aided C-NOMA schemes whose diversity order is only two in the considered C-NOMA system in the IoT.

Index Terms—Cooperative non-orthogonal multiple access, buffer-aided cooperation, outage probability, diversity order, average delay.

I. INTRODUCTION

Due to the connection of the massive number of devices in the Internet of Things (IoT), the design of multiple access (MA) technique faces significant challenges in wireless networks [1], [2]. Non-orthogonal multiple access (NOMA) is a critical promising MA technique for the future wireless networks [3]–[5], which has a great potential to solve the problem of massive connectivity in the IoT [6]–[9]. Using power-domain NOMA, a base station (BS) or access point (AP) serves more than one users at the same time, code

and frequency, but with different power levels. To deal with interference signals introduced by multiple users, intended receivers decode their own messages by applying successive interference cancellation (SIC). Compared to orthogonal multiple access (OMA), NOMA can achieve higher spectrum efficiency, lower transmission delay, and more user connections [4], [5].

Cooperative communication can further improve the efficiency and reliability of NOMA systems. Relay cooperation in NOMA systems were investigated in many existing works (e.g., [10]–[13]). In these existing works, dedicated relays between the BS and multiple users are utilized to assist the users to receive messages. On the other hand, cooperative NOMA (C-NOMA) with cooperation between intended users was first introduced in [14]. Motivated by the fact that the performance of NOMA is limited by the user with poor channel conditions, the work in [14] arranges the stronger users to help the weaker users for data transmission, so that all the users achieve the diversity order which is equal to the number of users. The basic idea of user cooperation in [14] was also extended to downlink NOMA systems in some other existing works (e.g., [15]–[18]).

Buffer-aided cooperative communication technique can provide an additional degree of freedom for the wireless cooperative communication networks, which helps to overcome the bottleneck effect of traditional cooperative communication technology (e.g., [19]–[27]). For buffer-aided relaying, the key issue is to design an efficient decision scheme at each buffer to receive or transmit packets, according to the instantaneous channel state information (CSI) and the buffer state. Existing buffer-aided relaying schemes have been shown to achieve full diversity orders of the corresponding relay networks. In addition, buffer-aided relaying for multiuser NOMA has also been investigated in recent years [9], [28]–[34]. In [28]–[32], downlink or uplink NOMA systems with a single buffer-aided relay were considered, where the BS transmits packets to the relay first, and then the relay uses NOMA to forward these packets to the users. Moreover, the work in [30] investigated an uplink NOMA system with a single buffer-aided relay. Recently, the works in [9], [33], [34] investigated downlink NOMA systems with a pair of users and multiple relays, and several relay selection schemes were proposed to enhance the outage performance.

Different from the existing buffer-aided relay NOMA systems [9], [28]–[34], this paper investigates a downlink buffer-aided C-NOMA system with an AP and multiple users in the IoT, where buffer-aided cooperation between a pair of intended users are considered. Note that conventional buffer-aided cooperation with dedicated relays is significantly dif-

The work of P. Xu was supported in part by the National Natural Science Foundation of China under Grant 61701066 and Grant 61971080, and in part by the Chongqing Natural Science Foundation Project under Grant cstc2019jcyj-msxmX0032. The work of G. Chen was supported by EPSRC grant number EP/R006377/1 (M3NETs). This paper was presented in part at IEEE/CIC International Conference on Communications in China (ICCC), Chongqing, China, 2020.

P. Xu and Y. Wang are with Chongqing Key Laboratory of Mobile Communications Technology, School of Communication and Information Engineering, Chongqing University of Posts and Telecommunications, Chongqing, 400065, P. R. China. (Email: xupeng@cqupt.edu.cn, wangyunwu@163.com).

G. Chen is with School of Engineering, University of Leicester, Leicester LE1 7RH, U.K. (Email: gaojie.chen@leicester.ac.uk).

G. Pan is with the School of Information and Electronics Engineering, Beijing Institute of Technology, Beijing 100081, China, and he is also with Computer, Electrical and Mathematical Sciences and Engineering Division, King Abdullah University of Science and Technology (KAUST), Thuwal 23955-6900, Saudi Arabia. (Email: penggaofeng@qq.com).

Z. Ding is with School of Electrical and Electronic Engineering, The University of Manchester, Manchester M13 9PL, U.K. (Email: zhiguo.ding@manchester.ac.uk).

Copyright (c) 2020 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

ferent from the considered buffer-aided scheme employing the cooperation between intended users. For example, the dedicated relays do not have their own packets to receive, and they just receive and retransmit the packets for those users; whereas, for buffer-aided cooperation between intended users in the considered system model, the intended users respectively have their own packets to receive, and also perform mutual buffer-aided cooperation with each other for packet transmission. Moreover, compared to buffer-aided dedicated relays, the considered buffer-aided intended users experience more complicated buffer state transition processes since each user receives part of its packets directly from the AP, and also receives the other part of packets from the other user.

Specifically, for the considered buffer-aided C-NOMA system, transmission modes in terms of the direct transmission (from the AP to the users) and the buffer-aided cooperative transmission (between a pair of users) are coordinated. It is highlighted that the issue of exploiting direct links in buffer-aided relaying systems is important to guarantee the transmission efficiency, which has been investigated in several existing works in [20], [32], [35]. However, these works only investigated the coordination between direct transmission and buffer-aided relaying transmission in cooperative systems with dedicated buffer-aided relays. This motivates us to investigate buffer-aided C-NOMA with direct transmission and cooperative transmission between intended users in the IoT. Although the recent work in [36] investigated C-NOMA with buffer-aided user cooperation, only unidirectional cooperation from a near user to a far user was considered, and the performance evaluation was focused only on the far user. To the best of the authors' knowledge, this is the first work to investigate buffer-aided C-NOMA with coordinated direct transmission and bidirectional buffer-aided cooperative transmission between the intended users.

Now, we summarize the contributions as follows:

- Proposing a novel buffer-aided C-NOMA scheme in the IoT, which adaptively selects a direct or cooperative transmission mode, based on the instantaneous CSI and the buffer state.
- Analyzing the performance of the proposed scheme with closed-form expressions for system outage probability (SOP) and average delay, by formulating the Markov chain (MC) of the proposed buffer-aided C-NOMA scheme and analyzing the corresponding state transition probabilities.
- Deriving the diversity order of the proposed scheme, by simplifying the original MC and approximating the SOP at high signal-to-noise ratio (SNR). The proposed buffer-aided C-NOMA scheme is proved to achieve the full diversity order of three for each user pair when the buffer size is larger than or equal to three.

It should be noted that the NOMA-based transmission mode is not simply injected into the considered cooperative system as an additional mode. In fact, the use of NOMA-based transmission mode leads to a global and deep impact on the design of buffer-aided C-NOMA schemes as well as the corresponding performance analysis. One can also consult

many related existing works (e.g., [9], [28]–[31], [33], [34]) considering the hybrid NOMA and OMA transmission, which clearly shows the transmission scheme design and performance analysis are significantly different from those in classical OMA based buffer-aided cooperative systems.

In this paper, Section II describes the system model and preliminary. In Section III, the buffer-aided C-NOMA scheme is proposed to adaptively select transmission modes. Section IV analyzes the SOP, average delay and diversity order. Section V discusses the proposed scheme with adaptive power allocation (PA). Simulation results are provided in Section VI. We conclude this paper in Section VII.

Throughout this paper, we define $[x]^+ = \max\{0, x\}$; $\log(\cdot)$ denotes the 2-based logarithm function; $[0 : K]$ denotes the set $\{0, 1, \dots, K\}$; “ \wedge ” and “ \vee ” denote “and” and “or”, respectively; $\bar{\mathcal{A}}$ denotes the complementary set of \mathcal{A} ; $\Pr\{\cdot\}$ denotes the probability of an event; $\mathbb{E}[\cdot]$ denotes the expectation of a random variable.

II. SYSTEM MODEL AND PRELIMINARIES

A. System Model

Consider a buffer-aided downlink C-NOMA system in the IoT, which includes a single-antenna AP and multiple single-antenna users. Since the user pairing technique is commonly used to reduce the complexity of power-domain NOMA [37], [38], these multiple users are partitioned into different groups, and each group includes a pair of users. Different user pairs are served on different orthogonal frequency resource blocks (RBs). Therefore, as shown in Fig. 1, we focus on a pair of users sharing in the same frequency RB. The details of how to design the user pairing strategy for the buffer-aided C-NOMA systems are out of the scope of this paper, which could be an interesting future topic. For each user pair in the considered C-NOMA system, the two users cooperate with each other with the help of buffers. It is assumed that time is divided into slots of equal length and each packet spans one time slot. In particular, each user i is equipped with a buffer, i.e., B_i , $i = 1, 2$. Each buffer consists of $L \geq 2$ storage units, where a storage unit at buffer B_i is used to store a packet intended to the other user. In addition, all devices are assumed to work under half-duplex modes, i.e., they cannot transmit and receive packets simultaneously. In each time slot, the AP or a user will be selected to transmit packets. Both the direct transmission from the AP to the two users and the buffer-aided cooperative transmission between the two users are considered. When the NOMA-based direct transmission cannot be performed, we adopt the cooperative transmission: the AP first transmits a mixed packet formed by the two users' packets to a certain user, and this user decodes both the two packets and stores the other user's information as a packet in one buffer unit; then, this stored packet will be eventually sent to the intended user in another time slot. The target transmission rate for user i is denoted by R_i in bits per slot.

The channel gain from the AP to user i is denoted as $h_i(t)$ in time slot t , and the channel gain between the two users is denoted as $h_3(t)$. The channels are assumed to be flat Rayleigh block fading channels which remain constant during one time

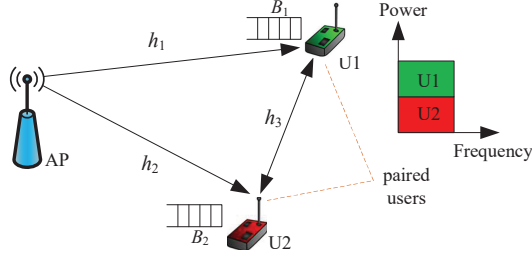


Fig. 1. System model of a buffer-aided C-NOMA system in the IoT.

slot and change randomly from one time slot to another. Each channel is modeled as $h_j(t) = d_j^{-\gamma/2} g_j(t)$ in time slot t , where the small scale fading gain is Rayleigh distributed, i.e., $g_j(t) \sim \mathcal{CN}(0, 1)$, γ is the path loss exponent and d_j denotes the corresponding distance, $j \in [1 : 3]$. Specifically, d_1 (d_2) denote the distances between the AP and user 1 (user 2), and d_3 is the distance between the two users. These distances remain constant for different time slots. Furthermore, $|h_j(t)|^2$ is exponentially distributed with mean $\Omega_j = d_j^{-\gamma}$. In addition, it is assumed that the AP and each user are constrained by the maximum transmit powers P_b and P_u , respectively. For the sake of brevity, the time slot index t will be omitted in the rest of this paper.

B. Transmission Modes and CSI Requirements

For the proposed buffer-aided C-NOMA system, six possible transmission modes are considered, which are denoted by $\mathcal{M}_0, \mathcal{M}_1, \dots, \mathcal{M}_5$, and the required CSI region for mode \mathcal{M}_k is defined as \mathcal{R}_k , where $k \in [0 : 5]$. Note that \mathcal{M}_0 is the direct transmission mode, $\{\mathcal{M}_k\}_{k \in [1:4]}$ are the cooperative transmission modes, and \mathcal{M}_5 is the silent mode. These transmission modes as well as their CSI requirements are shown in Table I, which are also described in details as follows.

1) *Mode \mathcal{M}_0* : Mode \mathcal{M}_0 denotes the direct transmission mode using the NOMA scheme, and both the two users can decode their packets correctly. Specifically, the AP sends a superposition signal ($x \triangleq \sqrt{\alpha_1} s_{\pi_1} + \sqrt{\alpha_2} s_{\pi_2}$) to the two users, where s_{π_i} ($\mathbb{E}[|s_{\pi_i}|^2] = 1$) carries information bits for user π_i , where $(\pi_1, \pi_2) = (1, 2)$ or $(2, 1)$, which denotes the ordering of the two users. In addition, α_i denotes the PA factor for user π_i satisfying $0 < \alpha_2 < \alpha_1$ and $\alpha_1 + \alpha_2 = 1$.¹ Therefore, user π_i receives

$$y_{\pi_i} = \sqrt{P_b} h_{\pi_i} (\sqrt{\alpha_1} s_{\pi_1} + \sqrt{\alpha_2} s_{\pi_2}) + n_{\pi_i}, \quad (1)$$

where n_{π_i} is the complex Gaussian noise at user π_i with mean zero and variance σ^2 .

Based on the decoding principle of SIC, user π_1 decodes its packet by treating user π_2 's signal as the noise, and the signal-to-interference-plus-noise ratio (SINR) is denoted by

$$\Gamma_1 \triangleq \frac{\alpha_1 \rho_b |h_{\pi_1}|^2}{\alpha_2 \rho_b |h_{\pi_1}|^2 + 1}, \quad (2)$$

¹The fixed PA strategy will be mainly considered for performance analysis in Section IV, and the adaptive PA strategy will also be discussed in Section V. Note that NOMA with fixed PA enjoys a low cost of system overhead since α_i remains unchanged and does not depend on the instantaneous CSI.

where $\rho_b \triangleq \frac{P_b}{\sigma^2}$. User π_2 first tries to decode the packet of user π_1 by treating its own signal as the noise, and the corresponding SINR is denoted by Γ_{21} ; then it cancels user π_1 's signal from its observation, and the receive SNR for decoding the packet of user π_2 is denoted by Γ_{22} , where

$$\Gamma_{21} \triangleq \frac{\alpha_1 \rho_b |h_{\pi_2}|^2}{\alpha_2 \rho_b |h_{\pi_2}|^2 + 1}, \quad \Gamma_{22} \triangleq \alpha_2 \rho_b |h_{\pi_2}|^2. \quad (3)$$

One can observe that the following conditions are necessary to perform \mathcal{M}_0 :

$$\log(1 + \min\{\Gamma_1, \Gamma_{21}\}) \geq R_{\pi_1}, \quad \log(1 + \Gamma_{22}) \geq R_{\pi_2}. \quad (4)$$

By considering the two possible ordering cases $(\pi_1, \pi_2) = (1, 2)$ and $(\pi_1, \pi_2) = (2, 1)$, the CSI region \mathcal{R}_0 is given by

$$\mathcal{R}_0 \triangleq \bigcup_{(|h_1|^2, |h_2|^2)} \left\{ \left\{ |h_1|^2 \geq \frac{\epsilon_{b1}}{\hat{\alpha}_1}, |h_2|^2 \geq \max \left\{ \frac{\epsilon_{b1}}{\hat{\alpha}_1}, \frac{\epsilon_{b2}}{\hat{\alpha}_2} \right\} \right\} \right. \\ \left. \vee \left\{ |h_1|^2 \geq \max \left\{ \frac{\epsilon_{b2}}{\hat{\alpha}_2}, \frac{\epsilon_{b1}}{\hat{\alpha}_1} \right\}, |h_2|^2 \geq \frac{\epsilon_{b2}}{\hat{\alpha}_2} \right\} \right\}, \quad (5)$$

where $\epsilon_{b1} \triangleq \frac{2^{R_1}-1}{\rho_b}$, $\epsilon_{b2} \triangleq \frac{2^{R_2}-1}{\rho_b}$, $\hat{\alpha}_1 \triangleq [1 - 2^{R_1} \alpha_2]^+$, and $\hat{\alpha}_2 \triangleq [1 - 2^{R_2} \alpha_2]^+$.

2) *Mode \mathcal{M}_1* : Mode \mathcal{M}_1 denotes the cooperative mode that the AP transmits both the two users' packets to user 1. In particular, the AP maps the two users' packets into one mixed packet [13], [29], from which user 1 can decode both the two users' packets correctly. Then, user 1 will store the packet of user 2 in one of its buffer units, and wait for another time slot to cooperatively transmit this packet to user 2. In this mode, the channel of user 1 is required to be so strong that user 1 can decode both the two users' packets correctly, i.e.,

$$\frac{1}{2} \log(1 + \rho_b |h_1|^2) \geq R_1 + R_2, \quad (6)$$

where the factor $\frac{1}{2}$ exists since two time slots are required to perform each cooperative transmission.

Accordingly, the CSI region \mathcal{R}_1 is given by

$$\mathcal{R}_1 \triangleq \bigcup_{|h_1|^2} \{ |h_1|^2 \geq \xi_{b1} \}, \quad (7)$$

where $\xi_{b1} \triangleq \frac{2^{2(R_1+R_2)}-1}{\rho_b}$.

3) *Mode \mathcal{M}_2* : Similar to \mathcal{M}_1 , \mathcal{M}_2 denotes the cooperative mode that the AP successfully transmits both the two users' packets to user 2, and user 2 will store user 1's packet in one of its buffer units. In this mode, h_2 is required to be so strong that user 2 can decode both the two users' packets correctly. Accordingly, the CSI region \mathcal{R}_2 is given by

$$\mathcal{R}_2 \triangleq \bigcup_{|h_2|^2} \{ |h_2|^2 \geq \xi_{b1} \}. \quad (8)$$

4) *Mode \mathcal{M}_3* : Mode \mathcal{M}_3 denotes that user 1 sends a packet stored in its buffer to user 2, and user 2 receives the following signal:

$$y_2^{(u)} = \sqrt{P_u} h_3 x_1^u + n_2^u, \quad (9)$$

TABLE I
NECESSARY REQUIREMENTS FOR EACH TRANSMISSION MODE (T: TRANSMIT; R: RECEIVE; S: SILENT)

Mode	AP	User 1	User 2	CSI Requirement	Buffer Requirement
\mathcal{M}_0	T	R	R	$\mathcal{R}_0 \triangleq \bigcup_{(h_1 ^2, h_2 ^2)} \left\{ \begin{aligned} & \left\{ h_1 ^2 \geq \frac{\epsilon_{b1}}{\alpha_1}, h_2 ^2 \geq \max \left\{ \frac{\epsilon_{b1}}{\alpha_1}, \frac{\epsilon_{b2}}{\alpha_2} \right\} \right\} \\ & \vee \left\{ h_1 ^2 \geq \max \left\{ \frac{\epsilon_{b2}}{\alpha_2}, \frac{\epsilon_{b1}}{\alpha_1} \right\}, h_2 ^2 \geq \frac{\epsilon_{b2}}{\alpha_2} \right\} \end{aligned} \right\}$	$\forall l_1, l_2$
\mathcal{M}_1	T	R	S	$\mathcal{R}_1 \triangleq \bigcup_{ h_1 ^2} \left\{ h_1 ^2 \geq \xi_{b1} \right\}$	$l_1 < L$
\mathcal{M}_2	T	S	R	$\mathcal{R}_2 \triangleq \bigcup_{ h_2 ^2} \left\{ h_2 ^2 \geq \xi_{b1} \right\}$	$l_2 < L$
\mathcal{M}_3	S	T	R	$\mathcal{R}_3 \triangleq \bigcup_{ h_3 ^2} \left\{ h_3 ^2 \geq \xi_{u2} \right\}$	$l_1 > 0$
\mathcal{M}_4	S	R	T	$\mathcal{R}_4 \triangleq \bigcup_{ h_3 ^2} \left\{ h_3 ^2 \geq \xi_{u1} \right\}$	$l_2 > 0$
\mathcal{M}_5	S	S	S	$\forall h_1 ^2, h_2 ^2, h_3 ^2$	$\forall l_1, l_2$

where x_1^u denotes the transmitted signal by user 1, which carries the information bits for user 2, and n_2^u is the complex Gaussian noise with zero mean and variance σ^2 . Thus, the necessary CSI requirement for \mathcal{M}_3 can be expressed as follows:

$$\frac{1}{2} \log(1 + \rho_u |h_3|^2) \geq R_2, \quad (10)$$

where $\rho_u \triangleq \frac{P_u}{\sigma^2}$ and the factor $\frac{1}{2}$ exists since two time slots are required to perform each cooperative transmission. Note that, although mode \mathcal{M}_3 only requires one time slot for user 1 to send a packet to user 2, another time slot has been previously consumed for user 1 to receive this packet from the AP via mode \mathcal{M}_1 . Accordingly, the CSI region \mathcal{R}_3 is given by

$$\mathcal{R}_3 \triangleq \bigcup_{|h_3|^2} \left\{ |h_3|^2 \geq \xi_{u2} \right\}, \quad (11)$$

where $\xi_{u2} \triangleq \frac{2^{2R_2}-1}{\rho_u}$.

5) *Mode \mathcal{M}_4* : Similar to \mathcal{M}_3 , \mathcal{M}_4 denotes that user 2 sends a packet stored in its buffer to user 1, and the CSI region \mathcal{R}_4 can be easily given by

$$\mathcal{R}_4 \triangleq \bigcup_{|h_3|^2} \left\{ |h_3|^2 \geq \xi_{u1} \right\}, \quad (12)$$

where $\xi_{u1} \triangleq \frac{2^{2R_1}-1}{\rho_u}$.

6) *Mode \mathcal{M}_5* : Mode \mathcal{M}_5 denotes that all nodes keep silent and do not transmit any packets, whose CSI requirement is arbitrary.

C. The Buffer Requirements

At each user, the number of storage units at each user is denoted as l_i with $0 \leq l_i \leq L$. \mathcal{M}_0 and \mathcal{M}_5 denote the direct transmission and silent modes, respectively, and hence they can be performed for any buffer state; \mathcal{M}_1 (\mathcal{M}_2) arranges user 1 (user 2) to store a packet intended for the other user, so it requires the corresponding buffers to be not full, i.e., $l_1 < L$ ($l_2 < L$); \mathcal{M}_3 (\mathcal{M}_4) arranges user 1 (user 2) to send a packet to the other user, so it requires the corresponding buffers to be not empty, $l_1 > 0$ ($l_2 > 0$).

III. PROPOSED BUFFER-AIDED C-NOMA SCHEME

In this section, we propose an efficient buffer-aided C-NOMA scheme to naturally select a transmission mode among $\{\mathcal{M}_k\}_{k \in [0:5]}$. For the basic idea of the proposed buffer-aided C-NOMA scheme, each mode \mathcal{M}_k is allocated with a buffer

state based selection function (SF), denoted by Λ_k . Note that a larger SF represents a higher priority to be selected. For example, if Λ_k has the largest value among $\{\Lambda_k\}_{k \in [0:5]}$, mode \mathcal{M}_k has the highest priority, which will be selected as long as its CSI requirement (shown in Table I) can be satisfied. If the CSI requirement for \mathcal{M}_k cannot be satisfied, we will keep on to consider the mode with the second highest priority, and so on.

Each SF Λ_k for mode \mathcal{M}_k , $k \in [0 : 5]$, is defined as follows:

$$\Lambda_0 \triangleq \beta_4, \quad (13a)$$

$$\Lambda_k \triangleq \beta_2(L - l_k), \quad k = 1, 2, \quad (13b)$$

$$\Lambda_k \triangleq \omega_k(\beta_3[l_{k-2} - l^{\text{thr}}]^+ + \beta_1), \quad k = 3, 4, \quad (13c)$$

$$\Lambda_5 \triangleq \beta_0, \quad (13d)$$

where $\omega_k = 0$ if $l_{k-2} = 0$, and $\omega_k = 1$, otherwise, for $k = 3, 4$; l^{thr} is a predefined threshold to control the value of Λ_k for $k = 3, 4$; $\{\beta_i\}_{i \in [0:4]}$ are used to denote five different priority levels of the transmission modes, $\beta_0 > 0$, $\beta_i \ll \beta_j$ for $\forall i, j \in [0 : 4]$ and $i < j$. The motivation of using ω_k in (13c) lies in setting $\Lambda_k = 0$ if $l_{k-2} = 0$, for $k = 3, 4$, which ensures that mode \mathcal{M}_k cannot be selected to transmit any packet if the corresponding buffer is empty. Specifically, direct transmission mode \mathcal{M}_0 has the highest priority by setting $\Lambda_0 \triangleq \beta_4$, which means that \mathcal{M}_0 will be selected as long as its required CSI is satisfied. In other words, any other transmission mode can only be selected when $(|h_1|^2, |h_2|^2) \notin \mathcal{R}_0$.

Here we briefly explain the motivation of the setting of Λ_k in (13b) and (13c). In particular, when $l_{k-2} \geq l^{\text{thr}} + 1$, $k = 3, 4$, mode \mathcal{M}_k lies on level β_3 , which enjoys the second highest priority level. However, when $0 < l_{k-2} \leq l^{\text{thr}}$, \mathcal{M}_k falls down to level β_1 , and this level is lower than level β_2 corresponding to \mathcal{M}_1 and \mathcal{M}_2 . This means that each buffer prioritizes to store rather than send a packet if its length does not exceed l^{thr} . Finally, \mathcal{M}_5 lies on the lowest priority level β_0 , which will be selected only if the SF values of any other modes do not exceed β_0 or their CSI requirements cannot be satisfied.

With the allocated SF for each mode, the buffer-aided C-NOMA scheme can be mathematically expressed as follows. In particular, mode \mathcal{M}_{k^*} is selected in each time slot, where

$$k^* = \arg \max_{k \in [0:5], \mathcal{H} \in \mathcal{R}_k} \Lambda_k, \quad (14)$$

and $\mathcal{H} \triangleq \{|h_i|^2\}_{i \in [1:3]}$ denotes the global instantaneous CSI. Note that there may exist more than one modes with the

same SF value as shown in (13). For example, $\Lambda_1 = \Lambda_2$ and $\Lambda_3 = \Lambda_4$ occur when $l_1 = l_2$. Thus, there exists a positive possibility that more than one modes are selected using (14). In this case, we randomly select a mode among these modes, for the sake of fairness. Besides, k^* in (14) can be easily obtained by jointly considering the CSI and SF values. For example, the AP can first identify the feasible transmission modes by checking their CSI requirements, and then select the one with the largest SF value among these modes whose CSI requirements are satisfied. We have the following remarks for the proposed mode selection scheme in (14).

Remark 1: At the beginning of each time slot, the AP can be employed to select a transmission mode, based on channel estimation and keeping track of buffer length at each user. The AP can obtain the CSI of h_1 and h_2 based on pilot symbols transmitted by the users, and the AP can know whether h_3 lies in \mathcal{R}_3 or \mathcal{R}_4 based on one-bit feedback from each user. However, it should be noted that the AP does not always need to receive feedback concerning h_3 . For example, if $(|h_1|^2, |h_2|^2) \in \mathcal{R}_0$, mode \mathcal{M}_0 will be selected regardless of the strength of h_3 .

Remark 2: For the proposed scheme, we set $l^{\text{thr}} = 1$, which means that the each buffer will give priority to transmit rather than to receive if its buffer length is not less than 2. This setting of l^{thr} is to balance the performance between the SOP and average delay [25], [29], so that each buffer is likely to remain at length 1 or 2 especially at high SNR. Therefore, each buffer is neither full nor empty in most time slots as long as $L \geq 3$. Note that enlarging l^{thr} may enhance the SOP performance, but will lead to a higher delay.

Remark 3: The proposed scheme ensures that the direct transmission mode \mathcal{M}_0 has the highest priority, so that the AP can timely serve both the two users. However, cooperative modes are also important to enhance the diversity order, which will be selected if there exists a user whose channel condition is so poor that the CSI requirement of \mathcal{M}_0 cannot be satisfied.

IV. PERFORMANCE ANALYSIS

In this section, the performance of the proposed buffer-aided C-NOMA scheme will be analyzed by formulating a MC as well as its transition matrix to model the evolution of the buffer states. To formulate the transition matrix, we first provide some preliminary results in the following subsection.

A. Preliminary Results

Base on (5) and the fact that $|h_1|^2$ and $|h_2|^2$ are exponentially distributed, the probability that the CSI requirement of

\mathcal{M}_0 is satisfied can be derived as follows:

$$\begin{aligned} \phi_0 &\triangleq \Pr\{|h_1|^2, |h_2|^2 \in \mathcal{R}_0\} \\ &= \int_{\zeta_1}^{\infty} \int_{\frac{\epsilon_{b1}}{\hat{\alpha}_1}}^{\infty} \frac{1}{\Omega_1} \exp\left(-\frac{x}{\Omega_1}\right) \frac{1}{\Omega_2} \exp\left(-\frac{y}{\Omega_2}\right) dx dy \\ &\quad + \int_{\frac{\epsilon_{b2}}{\hat{\alpha}_2}}^{\infty} \int_{\zeta_2}^{\infty} \frac{1}{\Omega_1} \exp\left(-\frac{x}{\Omega_1}\right) \frac{1}{\Omega_2} \exp\left(-\frac{y}{\Omega_2}\right) dx dy \\ &\quad - \int_{\zeta_1}^{\infty} \int_{\zeta_2}^{\infty} \frac{1}{\Omega_1} \exp\left(-\frac{x}{\Omega_1}\right) \frac{1}{\Omega_2} \exp\left(-\frac{y}{\Omega_2}\right) dx dy \\ &= \exp\left(-\frac{\epsilon_{b1}}{\Omega_1 \hat{\alpha}_1} - \frac{\zeta_1}{\Omega_2}\right) + \exp\left(-\frac{\zeta_2}{\Omega_1} - \frac{\epsilon_{b2}}{\Omega_2 \hat{\alpha}_2}\right) \\ &\quad - \exp\left(-\frac{\zeta_2}{\Omega_1} - \frac{\zeta_1}{\Omega_2}\right), \end{aligned} \quad (15)$$

where $\zeta_1 \triangleq \max\left\{\frac{\epsilon_{b1}}{\hat{\alpha}_1}, \frac{\epsilon_{b2}}{\hat{\alpha}_2}\right\}$, and $\zeta_2 \triangleq \max\left\{\frac{\epsilon_{b2}}{\hat{\alpha}_2}, \frac{\epsilon_{b1}}{\hat{\alpha}_1}\right\}$.

Since a necessary condition to select \mathcal{M}_1 is that \mathcal{M}_0 's CSI requirement is not satisfied (Section III), we need to calculate the probability that \mathcal{M}_0 's CSI requirement is not satisfied but \mathcal{M}_1 's CSI requirement is satisfied. From (5) and (7), we have

$$\begin{aligned} \phi_1 &\triangleq \Pr\{|h_1|^2, |h_2|^2 \notin \mathcal{R}_0 \wedge |h_1|^2 \in \mathcal{R}_1\} \\ &= \int_0^{\frac{\epsilon_{b2}}{\hat{\alpha}_2}} \int_{\xi_{b1}}^{\infty} \frac{1}{\Omega_1} \exp\left(-\frac{x}{\Omega_1}\right) \frac{1}{\Omega_2} \exp\left(-\frac{y}{\Omega_2}\right) dx dy \\ &= \exp\left(-\frac{\xi_{b1}}{\Omega_1}\right) \left[1 - \exp\left(-\frac{\epsilon_{b2}}{\Omega_2 \hat{\alpha}_2}\right)\right]. \end{aligned} \quad (16)$$

Similarly, the probability that \mathcal{M}_0 's CSI requirement is not satisfied but \mathcal{M}_2 's CSI requirement is satisfied can be expressed as follows:

$$\begin{aligned} \phi_2 &\triangleq \Pr\{|h_1|^2, |h_2|^2 \notin \mathcal{R}_0 \wedge |h_2|^2 \in \mathcal{R}_2\} \\ &= \int_{\xi_{b1}}^{\infty} \int_0^{\frac{\epsilon_{b1}}{\hat{\alpha}_1}} \frac{1}{\Omega_1} \exp\left(-\frac{x}{\Omega_1}\right) \frac{1}{\Omega_2} \exp\left(-\frac{y}{\Omega_2}\right) dx dy \\ &= \exp\left(-\frac{\xi_{b1}}{\Omega_2}\right) \left[1 - \exp\left(-\frac{\epsilon_{b1}}{\Omega_1 \hat{\alpha}_1}\right)\right]. \end{aligned} \quad (17)$$

On the other hand, we can define and derive four probabilities with respect to the channel gain $|h_3|^2$ and the regions \mathcal{R}_3 and \mathcal{R}_4 as follows: specifically, ψ_0 denotes the probability that both the CSI requirements of modes \mathcal{M}_3 and \mathcal{M}_4 can be satisfied; ψ_1 denotes the probability that \mathcal{M}_3 's CSI requirement is satisfied but \mathcal{M}_4 's CSI requirement is not satisfied; ψ_2 denotes the probability that \mathcal{M}_4 's CSI requirement is satisfied but \mathcal{M}_3 's CSI requirement is not satisfied; ψ_3 denotes the probability that neither of the CSI requirements of these two modes is satisfied. Thus, $\{\psi_i\}_{i \in [0:3]}$ can be easily derived as

follows:

$$\begin{aligned}\psi_0 &\triangleq \Pr\{|h_3|^2 \in \mathcal{R}_3 \cap \mathcal{R}_4\} \\ &= \exp\left(-\frac{\max(\xi_{u1}, \xi_{u2})}{\Omega_3}\right),\end{aligned}\quad (18)$$

$$\begin{aligned}\psi_1 &\triangleq \Pr\{|h_3|^2 \in \mathcal{R}_3 \cap \bar{\mathcal{R}}_4\} \\ &= \left[\exp\left(-\frac{\xi_{u2}}{\Omega_3}\right) - \exp\left(-\frac{\xi_{u1}}{\Omega_3}\right)\right]^+, \quad (19)\end{aligned}$$

$$\begin{aligned}\psi_2 &\triangleq \Pr\{|h_3|^2 \in \bar{\mathcal{R}}_3 \cap \mathcal{R}_4\} \\ &= \left[\exp\left(-\frac{\xi_{u1}}{\Omega_3}\right) - \exp\left(-\frac{\xi_{u2}}{\Omega_3}\right)\right]^+, \quad (20)\end{aligned}$$

$$\begin{aligned}\psi_3 &\triangleq \Pr\{|h_3|^2 \in \bar{\mathcal{R}}_3 \cap \bar{\mathcal{R}}_4\} \\ &= 1 - \exp\left(-\frac{\min(\xi_{u1}, \xi_{u2})}{\Omega_3}\right).\end{aligned}\quad (21)$$

B. Transition Matrix and Stationary State Probabilities

Let $s \triangleq (l_1, l_2)$ denote the buffer state, which represents the queue lengths of the two buffers at the users. These $(L+1)^2$ states form the MC, and let \mathbf{A} denote the corresponding $(L+1) \times (L+1)$ state transition matrix, whose entry $\mathbf{A}_{i,j} \triangleq p(s_j \rightarrow s_i) = \Pr\{X_{t+1} = s_i | X_t = s_j\}$ is the transition probability to move from state s_j at time t to s_i at time $t+1$.

Based on the proposed buffer-aided C-NOMA scheme and the preliminary results in the previous subsection, the transition probabilities can be given in the following proposition.

Proposition 1: The transition probabilities of the states of the MC for the proposed buffer-aided C-NOMA scheme are expressed in (22)-(27), as shown at the top of the this page, where ϕ_i and ψ_j are given in (15)-(21), $i \in [0:2]$, $j \in [0:3]$.

Proof: Please refer to Appendix A. ■

The transition matrix in Proposition 1 is generally complicated. When $L = 2$, the transition matrix in Proposition 1 is given in (28) at the top of the next page, which provides transition probabilities among states $(0,0)$, $(0,1)$, $(0,2)$, $(1,0)$, $(1,1)$, $(1,2)$, $(2,0)$, $(2,1)$ and $(2,2)$, where $\phi_3 = 1 - \phi_0 - \phi_1 - \phi_2$, $\phi_0 = 1 - \phi_0$, $\phi_{ij} = \phi_i + \phi_j$ and $\psi_{ij} = \psi_i + \psi_j$, $\forall i, j \in [0:3]$.

Remark 4: In Proposition 1, $P_{(l_1, l_2)}^{(l_1, l_2)}$ corresponds to the probability that either mode \mathcal{M}_0 or \mathcal{M}_5 is selected. This is because the buffer state keeps unchanged for both the direct transmission mode and the silent mode. On the other hand, the buffer state changes only when any cooperative transmission mode \mathcal{M}_k , $k \in [1:4]$, is selected.

One can verify that the transition matrix \mathbf{A} is column stochastic, irreducible and aperiodic², so the stationary state probability vector can be obtained as follows [19]:

$$\boldsymbol{\pi} = (\mathbf{A} - \mathbf{I} + \mathbf{B})^{-1} \mathbf{b}, \quad (29)$$

where $\boldsymbol{\pi} = [\pi_{(0,0)}, \dots, \pi_{(L,L)}]^T$, $\mathbf{b} = [1, 1, \dots, 1]^T$ and $\mathbf{B}_{i,j} = 1$, $\forall i, j$.

In the following, the performance of the SOP, the diversity order and the average delay will be analyzed.

²Column stochastic means that all entries in any column sum up to one; irreducible means that it is possible to move from any state to any state; aperiodic means that it is possible to return to the same state at any time [39].

C. System Outage Probability

The SOP of the considered C-NOMA system is defined as the probability that both the AP and the two users remain silence, i.e., neither the AP nor the users transmit signals [29], [34]. When an outage event happens, neither the direct transmission mode nor the cooperative transmission modes can be performed. Thus, the SOP can be expressed as follows:

$$P_{out} = \sum_{n=1}^{(L+1)^2} \pi_n (\mathbf{A}_{nn} - \phi_0) = (\text{diag}(\mathbf{A}) - \phi_0) \boldsymbol{\pi}, \quad (30)$$

where \mathbf{A}_{nn} is the n th diagonal value in the transition matrix \mathbf{A} , which denotes the probability that the buffer state keeps unchanged at the n th state. In addition, ϕ_0 is subtracted from \mathbf{A}_{nn} , which can be explained as follows: although the buffer state keeps unchanged for both the direct transmission and silent modes (as discussed in Remark 4), only the silent mode corresponds to an outage event.

D. Diversity Order

To derive the diversity order, without loss of generality, we assume that $\rho_b = \rho_u = \rho$ for simplicity. The diversity order is defined as follows:

$$d \triangleq - \lim_{\rho \rightarrow \infty} \frac{\log P_{out}}{\log \rho}. \quad (31)$$

Proposition 2: The full diversity order of 3 can be achieved by the proposed buffer-aided C-NOMA scheme, as long as $L \geq 3$.

Proof: Please refer to Appendix B. ■

Compared to the traditional non-buffer-aided C-NOMA scheme which can only achieve a diversity order of 2 for the considered two-user C-NOMA system [14], [18], the proposed buffer-aided C-NOMA scheme achieves a diversity order of 3. Interestingly, the achieved diversity order is the same as the number of wireless links (i.e., 3) in the considered system, which means that the proposed buffer-aided C-NOMA scheme can fully exploit the wireless channels.

E. Average Delay

In general, the delay of the system includes two parts, i.e., the delay at the AP and the delay at the users. It is also assumed that the propagation delay in a wireless link is negligible. In addition, the target rates of the two users are assumed to be the same, i.e., $R_1 = R_2 = R_0$, for convenience. With these assumptions, we can summarize the derived delay of the proposed buffer-aided C-NOMA scheme in the following proposition.

Proposition 3: The average delay of the proposed buffer-aided C-NOMA scheme can be expressed as follows:

$$\bar{D} = \frac{1 + P_{out} - \phi_0 + 2 \sum_{l_1=0}^L \sum_{l_2=0}^L (l_1 + l_2) \pi_{(l_1, l_2)}}{2(1 - P_{out})}. \quad (32)$$

Proof: Please refer to Appendix C. ■

$$P_{(l_1, l_2)}^{(l'_1, l'_2)} = 0, \text{ if } |l'_1 - l_1| \geq 2 \vee |l'_2 - l_2| \geq 2 \vee \{|l'_1 = l_1 \pm 1| \wedge |l'_2 = l_2 \pm 1|\}. \quad (22)$$

$$P_{(l_1, l_2)}^{(l_1, l_2)} = \begin{cases} 1 - \phi_1 - \phi_2 & \text{if } l_1 = l_2 = 0, \\ (1 - \phi_0 - \phi_1 - \phi_2)(1 - \psi_0 - \psi_2) + \phi_0 & \text{if } l_1 = 0 \wedge 0 < l_2 < L, \\ (1 - \phi_0 - \phi_1 - \phi_2)(1 - \psi_0 - \psi_1) + \phi_0 & \text{if } l_2 = 0 \wedge 0 < l_1 < L, \\ (1 - \phi_0 - \phi_1)(1 - \psi_0 - \psi_2) + \phi_0 & \text{if } l_1 = 0 \wedge l_2 = L, \\ (1 - \phi_0 - \phi_2)(1 - \psi_0 - \psi_1) + \phi_0 & \text{if } l_2 = 0 \wedge l_1 = L, \\ (1 - \phi_0 - \phi_1 - \phi_2)\psi_3 + \phi_0 & \text{if } 0 < l_1 < L \wedge 0 < l_2 < L, \\ (1 - \phi_0 - \phi_2)\psi_3 + \phi_0 & \text{if } l_1 = L \wedge 0 < l_2 < L, \\ (1 - \phi_0 - \phi_1)\psi_3 + \phi_0 & \text{if } l_2 = L \wedge 0 < l_1 < L, \\ (1 - \phi_0)\psi_3 + \phi_0 & \text{if } l_1 = l_2 = L, \\ 0 & \text{otherwise.} \end{cases} \quad (23)$$

$$P_{(l_1, l_2)}^{(l_1+1, l_2)} = \begin{cases} \phi_1 & \text{if } 0 \leq l_1 \leq 1 \wedge 0 \leq l_2 \leq 1, \\ \phi_1(1 - \psi_0 - \psi_2) & \text{if } 0 \leq l_1 \leq 1 \wedge 2 \leq l_2 \leq L, \\ \phi_1(1 - \psi_0 - \psi_1) & \text{if } 2 \leq l_1 < L \wedge 0 \leq l_2 \leq 1, \\ \phi_1\psi_3 & \text{if } 2 \leq l_1 < L \wedge 2 \leq l_2 \leq L, \\ 0 & \text{otherwise.} \end{cases} \quad (24)$$

$$P_{(l_1, l_2)}^{(l_1, l_2+1)} = \begin{cases} \phi_2 & \text{if } 0 \leq l_2 \leq 1 \wedge 0 \leq l_1 \leq 1, \\ \phi_2(1 - \psi_0 - \psi_1) & \text{if } 0 \leq l_2 \leq 1 \wedge 2 \leq l_1 \leq L, \\ \phi_2(1 - \psi_0 - \psi_2) & \text{if } 2 \leq l_2 < L \wedge 0 \leq l_1 \leq 1, \\ \phi_2\psi_3 & \text{if } 2 \leq l_2 < L \wedge 2 \leq l_1 \leq L, \\ 0 & \text{otherwise.} \end{cases} \quad (25)$$

$$P_{(l_1, l_2)}^{(l_1-1, l_2)} = \begin{cases} (1 - \phi_0 - \phi_1 - \phi_2)(\psi_1 + \frac{1}{2}\psi_0) & \text{if } l_1 = 1 \wedge l_2 = 1, \\ (1 - \phi_0)(\psi_1 + \frac{1}{2}\psi_0) & \text{if } 2 \leq l_1 = l_2 \leq L, \\ (1 - \phi_0 - \phi_1 - \phi_2)(\psi_1 + \psi_0) & \text{if } l_1 = 1 \wedge l_2 = 0, \\ (1 - \phi_0 - \phi_1 - \phi_2)\psi_1 & \text{if } l_1 = 1 \wedge 2 \leq l_2 < L, \\ (1 - \phi_0 - \phi_1)\psi_1 & \text{if } l_1 = 1 \wedge l_2 = L, \\ (1 - \phi_0)\psi_1 & \text{if } 2 \leq l_1 < l_2 \leq L, \\ (1 - \phi_0)(\psi_1 + \psi_0) & \text{if } 2 \leq l_1 \leq L \wedge l_2 < l_1, \\ 0 & \text{otherwise.} \end{cases} \quad (26)$$

$$P_{(l_1, l_2)}^{(l_1, l_2-1)} = \begin{cases} (1 - \phi_0 - \phi_1 - \phi_2)(\psi_2 + \frac{1}{2}\psi_0) & \text{if } l_2 = 1 \wedge l_1 = 1, \\ (1 - \phi_0)(\psi_2 + \frac{1}{2}\psi_0) & \text{if } 2 \leq l_2 = l_1 \leq L, \\ (1 - \phi_0 - \phi_1 - \phi_2)(\psi_2 + \psi_0) & \text{if } l_2 = 1 \wedge l_1 = 0, \\ (1 - \phi_0 - \phi_1 - \phi_2)\psi_2 & \text{if } l_2 = 1 \wedge 2 \leq l_1 < L, \\ (1 - \phi_0 - \phi_2)\psi_2 & \text{if } l_2 = 1 \wedge l_1 = L, \\ (1 - \phi_0)\psi_2 & \text{if } 2 \leq l_2 < l_1 \leq L, \\ (1 - \phi_0)(\psi_2 + \psi_0) & \text{if } 2 \leq l_2 \leq L \wedge l_1 < l_2, \\ 0 & \text{otherwise.} \end{cases} \quad (27)$$

V. DISCUSSION OF ADAPTIVE POWER ALLOCATION

The performance of the proposed buffer-aided NOMA scheme can be further enhanced if the AP adopts adaptive PA to perform NOMA, i.e., the PA factors in Section II-B are adapted according to the instantaneous CSI. Using adaptive PA, the CSI requirement for the direct transmission mode \mathcal{M}_0 (i.e., the AP can use NOMA to transmit both users' packets successfully) can be written as follows [13], [29]:

$$\mathcal{R}_0^A \triangleq \bigcup_{(|h_1|^2, |h_2|^2)} \left\{ \left\{ \frac{\epsilon_{b1}}{|h_1|^2} + \frac{2^{R_1} \epsilon_{b2}}{|h_2|^2} \leq 1 \right\} \vee \left\{ \frac{\epsilon_{b2}}{|h_2|^2} + \frac{2^{R_2} \epsilon_{b1}}{|h_1|^2} \leq 1 \right\} \right\}. \quad (33)$$

Following similar derivation steps in Section IV and [13], [29], the stationary steady probabilities for the proposed buffer-

aided C-NOMA scheme with adaptive PA can be easily obtained. The details of deriving the performance of the proposed scheme with adaptive PA are omitted here for simplicity. We will rely on computer simulations later in Section VI to compare the performance of the proposed scheme with fixed PA and adaptive PA.

Remark 5: Despite the use of adaptive PA can achieve a better outage performance in comparison with fixed PA, the performance enhancement may be limited since the proposed scheme with fixed PA has achieved the full diversity order of three. On the other hand, it should be noted that NOMA with adaptive PA suffers from a higher cost of overhead since the AP in mode \mathcal{M}_0 needs to inform both users the PA factors at the beginning of each fading block.

$$\mathbf{A} = \begin{pmatrix} \phi_{03} & \phi_3\psi_{02} & 0 & \phi_3\psi_{01} & 0 & 0 & 0 & 0 & 0 & 0 \\ \phi_2 & \phi_3\psi_{13}+\phi_0 & \bar{\phi}_0\psi_{02} & 0 & \phi_3(\psi_1+\frac{1}{2}\psi_0) & 0 & 0 & 0 & 0 & 0 \\ 0 & \phi_2 & \phi_{23}\psi_{13}+\phi_0 & 0 & 0 & \phi_{23}\psi_1 & 0 & 0 & 0 & 0 \\ \phi_1 & 0 & 0 & \phi_3\psi_{23}+\phi_0 & \phi_3(\psi_2+\frac{1}{2}\psi_0) & 0 & \bar{\phi}_0\psi_{01} & 0 & 0 & 0 \\ 0 & \phi_1 & 0 & \phi_2 & \phi_3\psi_3+\phi_0 & \bar{\phi}_0\psi_{02} & 0 & \bar{\phi}_0\psi_{01} & 0 & 0 \\ 0 & 0 & \phi_1\psi_{13} & 0 & \phi_2 & \phi_{23}\psi_3+\phi_0 & 0 & 0 & 0 & \bar{\phi}_0(\psi_1+\frac{1}{2}\psi_0) \\ 0 & 0 & 0 & \phi_1 & 0 & 0 & \phi_{13}\psi_{23}+\phi_0 & \phi_{13}\psi_2 & 0 & 0 \\ 0 & 0 & 0 & 0 & \phi_1 & 0 & \phi_2\psi_{23} & \phi_{13}\psi_3+\phi_0 & \bar{\phi}_0(\psi_2+\frac{1}{2}\psi_0) & 0 \\ 0 & 0 & 0 & 0 & 0 & \phi_1\psi_{13} & 0 & \phi_2\psi_{23} & \phi_{13}\psi_3+\phi_0 & \bar{\phi}_0\psi_3+\phi_0 \end{pmatrix} \quad (28)$$

VI. NUMERICAL RESULTS

In this section, the performance of the proposed buffer-aided C-NOMA scheme is evaluated by using Monte Carlo simulations. Each channel gain is modeled as $h_j = d_j^{-\gamma/2} g_j$, where $g_j \sim \mathcal{CN}(0, 1)$, $j \in [1 : 3]$. Furthermore, we consider the network topology where the AP, user 1 and user 2 are located at the point $(0, 0)$, $(0, 1)$ and $(x_0, 0)$, respectively, which means that $d_1 = 1$, $d_2 = x_0$ and $d_3 = \sqrt{1 + x_0^2}$. The path loss exponent is chosen as $\gamma = 2$ to reflect a favorable propagation condition. For the other parameters, we also set $L = 5$, $\alpha_1 = 0.8$, $\alpha_2 = 0.2$, $\rho_b = \rho_u$, and $x_0 = 2$ meters for simplicity, unless stated otherwise.

A. Benchmark Schemes

Several benchmark schemes are used for performance comparison, which are listed as follows:

1) *Non-cooperative NOMA scheme*: The non-cooperative NOMA scheme means that the system only adopts NOMA to perform the direct transmission, which does not consider the cooperation between the two the users.

2) *Non-buffer-aided C-NOMA scheme*: The considered non-buffer-aided C-NOMA first tries to perform NOMA at the AP. If the CSI requirement for NOMA cannot be satisfied, the fading block will be divided into two phases, where the AP transmits both users' packets to the strong user in the first phase and the strong user repeat to transmit the packet to the weak user in the second phase.

3) *The scheme with priority to transmit*: The scheme with priority to transmit is a variant scheme of the proposed scheme, where the threshold in Section III is set to be $l^{\text{thr}} = 0$, which means that a buffer has a higher priority to transmit than to receive as long as it is not empty.

4) *SOP minimization scheme*: The SOP minimization scheme is another variant scheme of the proposed scheme, where the SOP is minimized by exhaustive searching an optimal threshold l^{thr} for each buffer, and the corresponding user gives priority to transmit or receive when the buffer length does or does not exceed this threshold, respectively.

B. Fixed Power Allocation

In this subsection, the performance of the proposed scheme and the comparative schemes will be displayed, where all schemes adopts fixed PA strategy.

In Fig. 2, the SOPs of the proposed buffer-aided C-NOMA scheme and several comparative schemes are depicted versus the transmit SNR in dB, where the target rates for the two

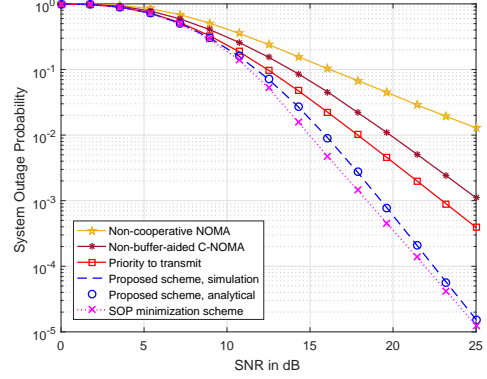


Fig. 2. SOPs of several schemes with fixed PA versus the transmit SNR in dB, where $R_1 = 1$ bit/slot, $R_2 = 0.5$ bits/slot.

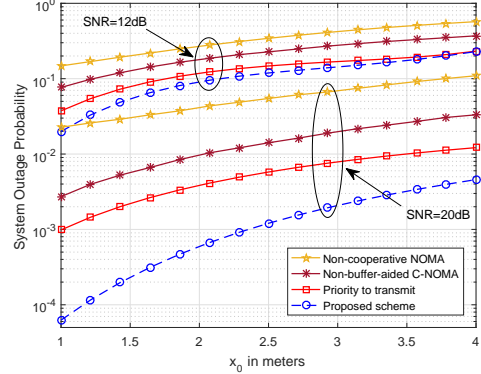


Fig. 3. SOPs of several schemes with fixed PA versus distance x_0 in meters, where $R_1 = 1$ bit/slot, $R_2 = 0.5$ bits/slot.

users are set as $R_1 = 1$ bit/slot and $R_2 = 0.5$ bits/slot. As shown in this figure, the proposed scheme significantly outperforms the non-cooperative NOMA scheme, the non-buffer-aided C-NOMA scheme and the scheme with priority to transmit. The analytical results in (30) match well with the simulation results. As discussed in Section IV-D, the proposed scheme can achieve the full diversity order of three, while the achieved diversity order of the non-buffer-aided C-NOMA scheme is only two. In addition, a gap exists between the proposed scheme and the SOP minimization scheme. However, this gap is negligible when the SNR approaches to 25 dB.

In Fig. 3, the SOP performance of the proposed scheme as well as the comparative schemes is shown as a function of x_0 in meters, where SNR= 12 or 20 dB, $R_1 = 1$ bit/slot and $R_2 = 0.5$ bits/slot. Since d_2 and d_3 increase with x_0

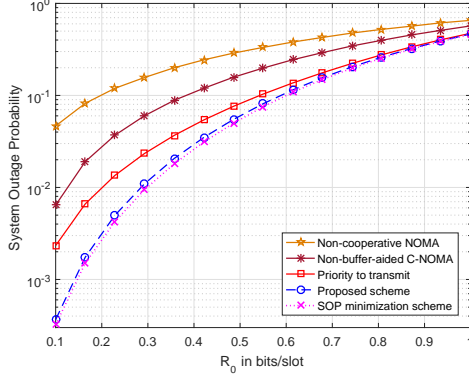


Fig. 4. SOPs of several schemes with fixed PA versus the target rate R_0 in bits/slot, where $\text{SNR}=10$ dB.

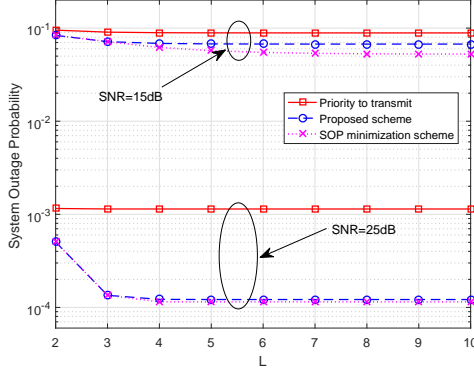


Fig. 5. SOPs of several schemes with fixed PA versus buffer size L , where $R_1 = R_2 = 1$ bit/slot.

for the considered network topology, the SOP of each scheme increases with x_0 . The proposed scheme achieves the lowest SOP in comparison with the non-cooperative NOMA scheme, the non-buffer-aided C-NOMA scheme and the scheme with priority to transmit. In addition, the gap between the curves of the proposed scheme and the comparative schemes becomes larger when the SNR changes from 12 dB to 20 dB, as shown in this figure.

In Fig. 4, we set $R_0 = R_1 = R_2$ (i.e., the same target rate for the two users are assumed) for simplicity, and the SOP performance of the proposed scheme is compared with the comparative schemes versus R_0 in bits/slot, where $\text{SNR}=10$ dB. The proposed scheme outperforms the non-cooperative NOMA scheme, the non-buffer-aided C-NOMA scheme and the scheme with priority to transmit when R_0 ranges from 0.1 to 1 bits/slot. Moreover, there is only a slight gap between the proposed scheme and the SOP minimization scheme especially when R_0 is larger than 0.5 bits/slot.

Fig. 5 presents the comparison of SOPs among the proposed scheme, the scheme with priority to transmit and the SOP minimization scheme, as functions of the buffer size L , where we set $R_1 = R_2 = 1$ bit/slot, and $\text{SNR}=15$ or 25 dB. The scheme with priority to transmit does not benefit from increasing L , whereas the SOP of the proposed scheme decreases with L when L changes from 2 to 4. When L exceeds 5, the SOP of the proposed scheme almost keeps unchanged since the user gives priority to transmit as long

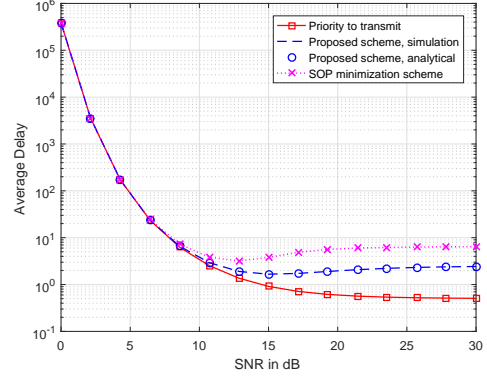


Fig. 6. Average delays of several schemes with fixed PA versus the transmit SNR in dB, where $R_1 = R_2 = 1$ bit/slot.

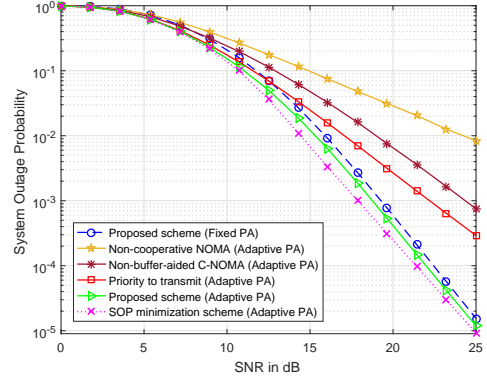


Fig. 7. SOPs of several schemes versus the transmit SNR in dB, where $R_1 = 1$ bit/slot, $R_2 = 0.5$ bits/slot.

as the buffer length exceeds the threshold of 2 as shown in Section III. In addition, a little gap exists between the curves of the proposed scheme and the SOP minimization scheme when $\text{SNR}=15$ dB, but the two schemes achieve almost the same SOP when $\text{SNR}=25$ dB.

Fig. 6 presents the performance of the average delay of the proposed buffer-aided C-NOMA scheme versus the transmit SNR in dB, where $R_1 = R_2 = 1$ bit/slot. The analytical results in Proposition 3 match well with the simulation results. The proposed scheme achieves a higher average delay compared to the scheme with priority to transmit, but a lower average delay compared to the SOP minimization scheme, which means that it achieves a tradeoff between the SOP and average delay. For example, when the SNR is large, Fig. 6 shows that the average delay is about 2.5 slots for the proposed scheme, whereas the SOP minimization scheme suffers a larger delay which is around 6.4 slots. This means that, compared to the SOP minimization scheme, the proposed scheme reduces the average delay by about 3.9 slots, at the price of only a slight loss of outage performance (shown in Fig. 2).

C. Adaptive Power Allocation

Fig. 7 shows the SOP performance for the proposed scheme with fixed PA and adaptive PA as well as the comparative schemes with adaptive PA. As can be seen from this figure, the proposed scheme also significantly outperforms the non-cooperative NOMA scheme, the non-buffer-aided C-NOMA

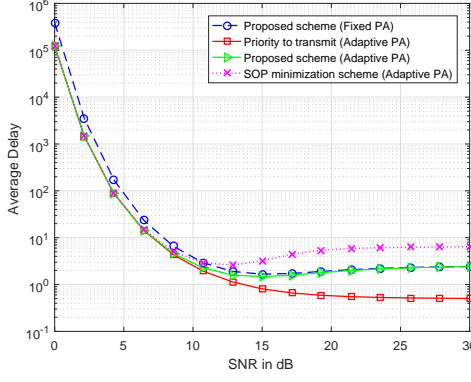


Fig. 8. Average delays of several schemes versus the transmit SNR in dB, where $R_1 = R_2 = 1$ bit/slot.

scheme and the scheme with priority to transmit. In addition, there only exists a little gap between the curves with respect to the proposed scheme with fixed PA and adaptive PA, which means that the SOP improvement for adaptive PA is limited for the proposed scheme as the full diversity order of three has been achieved by the fixed PA strategy.

Fig. 8 shows the average delay performance of the proposed scheme with fixed PA and adaptive PA as well as the scheme with priority to transmit (adaptive PA) and the SOP minimization scheme (adaptive PA). From this figure, one can also observe that the proposed scheme achieve a moderate average delay compared to the other two schemes. Moreover, there is only a little gap between the curves with respect to the proposed scheme with fixed PA and adaptive PA when the SNR is less than 15dB, and these two curves overlap each other for high SNRs.

VII. CONCLUSION

In this paper, a buffer-aided C-NOMA system in the IoT was investigated, where the AP communicates with multiple users, and these users are equipped with buffers. The direct transmission from the AP to the users and the buffer-aided cooperative transmission between a pair of users were coordinated. In particular, we proposed a buffer-aided C-NOMA scheme, which adaptively selects the direct and buffer-aided cooperative transmission modes. Then, the SOP and average delay of the proposed scheme were theoretically derived with closed-form expressions. Furthermore, we proved that the full diversity order of three can be achieved by the proposed scheme as long as the buffer size is not less than three, whereas the conventional non-buffer-aided C-NOMA schemes can only achieve a diversity order of two. An future interest is to design the optimal user pairing strategy as well as the optimal resource allocation strategy in the considered buffer-aided C-NOMA system in the IoT.

APPENDIX A PROOF OF PROPOSITION 1

To prove this proposition, we first define $\{F_k\}_{k \in [0:5]}$ as the probability that mode \mathcal{M}_k is selected, where $\sum_{k=0}^5 F_k = 1$ and $F_0 = \phi_0$ as shown in (15). Based on the proposed buffer-aided C-NOMA scheme in Section III and the preliminary

results in Section IV-A, we can derive different transition probabilities as follows:

- 1) Since each buffer at most receives or transmits only one packet in one time slot, $P_{(l_1, l_2)}^{(l'_1, l'_2)} = 0$ if $|l'_i - l_i| \geq 2$, $i = 1, 2$. Moreover, according to the proposed buffer-aided C-NOMA scheme, the two buffer lengths cannot change in the same time slot, i.e., one buffer length should keep unchanged if the other buffer length increases or decreases. Thus, (22) can be obtained.
 - 2) $P_{(l_1, l_2)}^{(l_1, l_2)}$ corresponds to the case that the buffer state keeps unchanged, which means that the direct mode \mathcal{M}_0 or the silent mode \mathcal{M}_5 is selected, i.e., $P_{(l_1, l_2)}^{(l_1, l_2)} = F_0 + F_5$. Note that \mathcal{M}_0 has the highest priority to be selected, but \mathcal{M}_5 has the lowest priority to be selected. $P_{(l_1, l_2)}^{(l_1, l_2)}$ is derived according to different buffer states as follows:
 - When $l_1 = l_2 = 0$, the buffer requirement of modes \mathcal{M}_3 and \mathcal{M}_4 cannot be satisfied. Thus, \mathcal{M}_5 can be selected only if none of the CSI requirements of modes $\{\mathcal{M}_k\}_{k \in [0:2]}$ can be satisfied. Based on Section IV-A, it is easy to obtain that $F_5 = 1 - \phi_0 - \phi_1 - \phi_2$, and $P_{(l_1, l_2)}^{(l_1, l_2)} = F_0 + F_5 = 1 - \phi_1 - \phi_2$.
 - When $l_1 = 0 \wedge 0 < l_2 < L$, the buffer requirement of \mathcal{M}_3 cannot be satisfied. \mathcal{M}_5 can only be selected if none of the CSI requirements of modes $\{\mathcal{M}_k\}_{k \in \{0, 1, 2, 4\}}$ can be satisfied. In this case, based on Section IV-A, $F_5 = (1 - \phi_0 - \phi_1 - \phi_2)(1 - \psi_0 - \psi_2)$, thus $P_{(l_1, l_2)}^{(l_1, l_2)} = (1 - \phi_0 - \phi_1 - \phi_2)(1 - \psi_0 - \psi_2) + \phi_0$.
 - Similarly, when $l_2 = 0 \wedge 0 < l_1 < L$, we can obtain $P_{(l_1, l_2)}^{(l_1, l_2)} = (1 - \phi_0 - \phi_1 - \phi_2)(1 - \psi_0 - \psi_1) + \phi_0$.
 - When $l_1 = 0 \wedge l_2 = L$, the buffer requirements of \mathcal{M}_2 and \mathcal{M}_3 cannot be satisfied. \mathcal{M}_5 can only be selected if none of the CSI requirements of $\{\mathcal{M}_k\}_{k \in \{0, 1, 4\}}$ can be satisfied. In this case, $F_5 = (1 - \phi_0 - \phi_1)(1 - \psi_0 - \psi_2)$, and thus $P_{(l_1, l_2)}^{(l_1, l_2)} = (1 - \phi_0 - \phi_1)(1 - \psi_0 - \psi_2) + \phi_0$.
 - Similarly, when $l_2 = 0 \wedge l_1 = L$, we obtain $P_{(l_1, l_2)}^{(l_1, l_2)} = (1 - \phi_0 - \phi_2)(1 - \psi_0 - \psi_1) + \phi_0$.
 - When $0 < l_1 < L \wedge 0 < l_2 < L$, \mathcal{M}_5 can only be selected if none of the CSI requirements of $\{\mathcal{M}_k\}_{k \in [0:4]}$ can be satisfied. Thus, $F_5 = (1 - \phi_0 - \phi_1 - \phi_2)\psi_3$, and $P_{(l_1, l_2)}^{(l_1, l_2)} = (1 - \phi_0 - \phi_1 - \phi_2)\psi_3 + \phi_0$.
 - When $0 < l_1 < L \wedge l_2 = L$, the buffer requirement of mode \mathcal{M}_2 cannot be satisfied. \mathcal{M}_5 can only be selected if none of the CSI requirements of $\{\mathcal{M}_k\}_{k \in \{0, 1, 3, 4\}}$ can be satisfied. Thus, $F_5 = (1 - \phi_0 - \phi_1)\psi_3$, and $P_{(l_1, l_2)}^{(l_1, l_2)} = (1 - \phi_0 - \phi_1)\psi_3 + \phi_0$.
 - Similarly, When $l_1 = L \wedge 0 < l_2 < L$, we obtain $P_{(l_1, l_2)}^{(l_1, l_2)} = (1 - \phi_0 - \phi_2)\psi_3 + \phi_0$.
 - When $l_1 = l_2 = L$, the buffer requirements of \mathcal{M}_1 and \mathcal{M}_2 cannot be satisfied. \mathcal{M}_5 can only be selected if none of the CSI requirements of $\{\mathcal{M}_k\}_{k \in \{0, 3, 4\}}$ can be satisfied. Thus, $F_5 = (1 - \phi_0)\psi_3$, and $P_{(l_1, l_2)}^{(l_1, l_2)} = (1 - \phi_0)\psi_3 + \phi_0$.
- In summary, (23) can be obtained.
- 3) $P_{(l_1, l_2)}^{(l_1+1, l_2)}$ corresponds to the case that the mode \mathcal{M}_1 is selected, which is derived according to different buffer states as follows:

- When $0 \leq l_1 \leq 1 \wedge 0 \leq l_2 \leq 1$, according to the proposed buffer-aided C-NOMA scheme in Section III, \mathcal{M}_1 can be selected only if the CSI requirement of \mathcal{M}_1 can be satisfied but the CSI requirement of \mathcal{M}_0 cannot be satisfied, i.e., $(|h_1|^2, |h_2|^2) \notin \mathcal{R}_0 \wedge |h_1|^2 \in \mathcal{R}_1$. In this case, $P_{(l_1, l_2)}^{(l_1+1, l_2)} = \phi_1$ as shown in (16).
- When $0 \leq l_1 \leq 1 \wedge 2 \leq l_2 \leq L$, \mathcal{M}_1 can be selected only if the CSI requirement of \mathcal{M}_1 can be satisfied, but none of the CSI requirements of \mathcal{M}_0 and \mathcal{M}_4 can be satisfied. Thus, $P_{(l_1, l_2)}^{(l_1+1, l_2)} = \phi_1(1 - \psi_0 - \psi_2)$.
- Similarly, When $2 \leq l_1 < L \wedge 0 \leq l_2 \leq 1$, we can derive that $P_{(l_1, l_2)}^{(l_1+1, l_2)} = \phi_1(1 - \psi_0 - \psi_1)$.
- When $2 \leq l_1 < L \wedge 2 \leq l_2 \leq L$, $\Lambda_0 > \{\Lambda_i\}_{i \in [3:4]} > \{\Lambda_j\}_{j \in [1:2]}$, \mathcal{M}_1 can be selected only if the CSI requirement of \mathcal{M}_1 can be satisfied but none of the CSI requirements of $\{\mathcal{M}_k\}_{k \in \{0, 3, 4\}}$ can be satisfied. In this case, we can obtain $P_{(l_1, l_2)}^{(l_1+1, l_2)} = \phi_1 \psi_3$.

In summary, (24) can be obtained.

- 4) Following similar derivation steps of $P_{(l_1, l_2)}^{(l_1+1, l_2)}$, $P_{(l_1, l_2)}^{(l_1, l_2+1)}$ can be obtained shown in (25).
- 5) $P_{(l_1, l_2)}^{(l_1-1, l_2)}$ corresponds to the case that the mode \mathcal{M}_3 is selected. In the following, we only derive $P_{(l_1, l_2)}^{(l_1-1, l_2)}$ for two buffer states in (26) as follows:
 - When $l_1 = l_2 = 1$, there are two events to select \mathcal{M}_3 : (i) The CSI requirements of mode \mathcal{M}_3 can be satisfied but none of the CSI requirements of modes $\{\mathcal{M}_k\}_{k \in \{0, 1, 2, 4\}}$ can be satisfied, whose probability is $(1 - \phi_0 - \phi_1 - \phi_2)\psi_1$. (ii) Both the CSI requirements of mode \mathcal{M}_3 and \mathcal{M}_4 can be satisfied, but none of the CSI requirements of modes $\{\mathcal{M}_k\}_{k \in \{0, 1, 2\}}$ can be satisfied; in addition, \mathcal{M}_3 is randomly selected among $\{\mathcal{M}_k\}_{k \in [3:4]}$. This event has a probability of $\frac{1}{2}\psi_0(1 - \phi_0 - \phi_1 - \phi_2)$. Thus, $P_{(l_1, l_2)}^{(l_1-1, l_2)} = (1 - \phi_0 - \phi_1 - \phi_2)(\psi_1 + \frac{1}{2}\psi_0)$.
 - When $2 \leq l_1 = l_2 \leq L$, the derivation steps for $P_{(l_1, l_2)}^{(l_1-1, l_2)}$ are almost the same as those in the previous case, except that $(1 - \phi_0 - \phi_1 - \phi_2)$ should be replaced by $(1 - \phi_0)$ since \mathcal{M}_3 and \mathcal{M}_4 have higher priority than \mathcal{M}_1 and \mathcal{M}_2 in this case. Therefore, we can obtain $P_{(l_1, l_2)}^{(l_1-1, l_2)} = (1 - \phi_0)(\psi_1 + \frac{1}{2}\psi_0)$.
 - When $l_1 = 1 \wedge l_2 = 0$, the buffer requirement of mode \mathcal{M}_4 cannot be satisfied. \mathcal{M}_3 can only be selected if none of the CSI requirements of $\{\mathcal{M}_k\}_{k \in \{0, 1, 2\}}$ can be satisfied. Thus, $P_{(l_1, l_2)}^{(l_1-1, l_2)} = (1 - \phi_0 - \phi_1 - \phi_2)(\psi_1 + \psi_0)$.
 - When $l_1 = 1 \wedge 2 \leq l_2 < L$, \mathcal{M}_3 can be selected only if the CSI requirement of \mathcal{M}_3 can be satisfied but none of the CSI requirements of $\{\mathcal{M}_k\}_{k \in \{0, 1, 2, 4\}}$ can be satisfied. In this case, we can obtain $P_{(l_1, l_2)}^{(l_1-1, l_2)} = (1 - \phi_0 - \phi_1 - \phi_2)\psi_1$.
 - When $l_1 = 1 \wedge l_2 = L$, the buffer requirement of mode \mathcal{M}_2 cannot be satisfied. \mathcal{M}_3 can only be selected if none of the CSI requirements of $\{\mathcal{M}_k\}_{k \in \{0, 1, 4\}}$ can be satisfied. Thus, $P_{(l_1, l_2)}^{(l_1-1, l_2)} = (1 - \phi_0 - \phi_1)\psi_1$.
 - When $2 \leq l_1 < l_2 \leq L$, \mathcal{M}_3 can be selected only if the

CSI requirement of \mathcal{M}_3 can be satisfied, but none of the CSI requirements of $\{\mathcal{M}_k\}_{k \in \{0, 4\}}$ can be satisfied.

Thus, $P_{(l_1, l_2)}^{(l_1-1, l_2)} = (1 - \phi_0)\psi_1$.

- When $2 \leq l_1 \leq L \wedge l_2 < l_1$, \mathcal{M}_3 can be selected only if the CSI requirement of \mathcal{M}_3 can be satisfied but the CSI requirement of \mathcal{M}_0 cannot be satisfied. In this case, we can obtain $P_{(l_1, l_2)}^{(l_1-1, l_2)} = (1 - \phi_0)(\psi_1 + \psi_0)$.

- 6) Following similar derivation steps of $P_{(l_1, l_2)}^{(l_1-1, l_2)}$, $P_{(l_1, l_2)}^{(l_1, l_2-1)}$ can be obtained shown in (27).

APPENDIX B

PROOF OF PROPOSITION 2

To derive the diversity order of the proposed buffer-aided C-NOMA scheme, $L = 3$ is considered and all the $(3+1)^2 = 16$ buffer states are divided into three state sets for simplicity. As shown in Fig. 9, each set includes several states with the same number of available³ links with respect to the buffers, where $\hat{S}_1 \triangleq \{(0, 1), (1, 0), (1, 1), (0, 2), (1, 2), (2, 0), (2, 1), (2, 2)\}$, $\hat{S}_2 \triangleq \{(0, 0), (0, 3), (1, 3), (2, 3), (3, 0), (3, 1), (3, 2)\}$, and $\hat{S}_3 \triangleq \{(3, 3)\}$. One can verify that every state in \hat{S}_i has $4 - i$ available links. For example, state $(0, 1)$ has 2 available receive links for the two buffers (i.e., h_1 and h_2) and 1 available transmit link for buffer B_2 (i.e., h_3), so it has totally 3 available links.

As shown in Fig. 9, the three state sets $\{\hat{S}_i\}_{i \in [1:3]}$ form a simplified MC, where \hat{S}_1 and \hat{S}_3 can only be connected through \hat{S}_2 . Each transition probability from the state set \hat{S}_i to the state set \hat{S}_j , denoted by $\hat{P}_i^j \triangleq \Pr(X_{t+1} \in \hat{S}_j | X_t \in \hat{S}_i)$, $\forall i, j \in [1 : 3]$. Obviously $\hat{P}_i^j > 0$ if $|i - j| \leq 1$, and $\hat{P}_i^j = 0$, otherwise. To derive \hat{P}_i^j , we are interested in the edge states⁴ of each set. As shown in Fig. 9, for \hat{S}_1 , except $(1, 1)$, all the other 7 states are edge states; for \hat{S}_2 and \hat{S}_3 , all the corresponding states are edge states.

A. Transition Probabilities

Closed-form expressions for \hat{P}_i^j are hard to be obtained. Alternatively, we will derive the exponential orders of \hat{P}_i^j as well as the stationary probabilities with respect to the MC formed by these three state sets. Note that $f(\rho)$ is termed to be exponentially equal to ρ^b , denoted by $f(\rho) \doteq \rho^b$, if $\lim_{\rho \rightarrow \infty} \frac{\log f(\rho)}{\log \rho} = b$; $f(\rho) \gtrsim \rho^b$ and $f(\rho) \lesssim \rho^b$ can also be defined similarly [40].

To obtain exponential order of \hat{P}_1^2 , we need to analyze the transition probabilities of the corresponding edge states. From (26), $P_{(1, 0)}^{(0, 0)}$ can be approximated as follows:

$$\begin{aligned}
 P_{(1, 0)}^{(0, 0)} &= (1 - \phi_0 - \phi_1 - \phi_2)(\psi_1 + \psi_0) \approx (1 - \phi_0 - \phi_1 - \phi_2) \\
 &\approx \left[1 - \exp\left(-\frac{\zeta_2}{\Omega_1}\right)\right] \times \left[1 - \exp\left(-\frac{\zeta_1}{\Omega_2}\right)\right] \\
 &\approx \frac{\zeta_1 \zeta_2}{\Omega_1 \Omega_2}.
 \end{aligned} \tag{34}$$

³Note that a receive (or transmit) link with respect to a buffer is available when this buffer is not full (or not empty).

⁴If a state in set \hat{S}_i in time slot t has a positive probability to transit to another state in set \hat{S}_j in time slot $t + 1$, $i \neq j$, it is termed as an edge state.

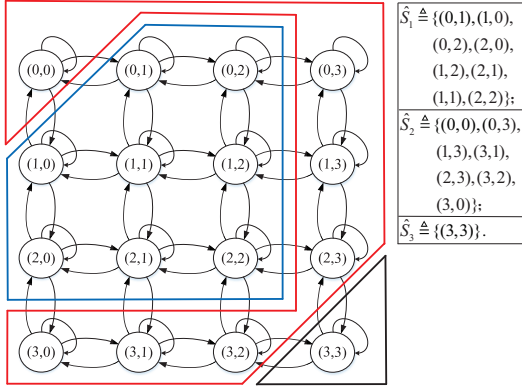


Fig. 9. Diagram of the formed state sets.

Thus, we have $P_{(1,0)}^{(0,0)} \doteq \rho^{-2}$. Similarly, $P_{(0,1)}^{(0,0)} \doteq \rho^{-2}$, $P_{(0,2)}^{(0,3)} \doteq \rho^{-2}$, $P_{(1,2)}^{(1,3)} \doteq \rho^{-2}$, $P_{(2,2)}^{(2,3)} \doteq \rho^{-2}$, $P_{(2,2)}^{(3,2)} \doteq \rho^{-2}$, $P_{(2,1)}^{(3,1)} \doteq \rho^{-2}$ and $P_{(2,0)}^{(3,0)} \doteq \rho^{-2}$ can be obtained, and thus

$$\begin{aligned} \hat{P}_1^2 &= \Pr(X_{t+1} \in \hat{S}_2 | X_t \in \hat{S}_1) \\ &= \sum_{s \in \hat{S}_1} \Pr(X_{t+1} \in \hat{S}_2 | X_t = s) \Pr(X_t = s | X_t \in \hat{S}_1) \\ &= \sum_{i=1}^7 P_{u_i}^{v_i} \Pr(X_t = u_i | X_t \in \hat{S}_1) \\ &\doteq \rho^{-2} (1 - \Pr(X_t = (1,1) | X_t \in \hat{S}_1)) \leq \rho^{-2}, \end{aligned} \quad (35)$$

where we denote $u_1 = (1,0)$, $u_2 = (0,1)$, $u_3 = (0,2)$, $u_4 = (1,2)$, $u_5 = (2,2)$, $u_6 = (2,1)$, $u_7 = (2,0)$; and $v_1 = v_2 = (0,0)$, $v_3 = (0,3)$, $v_4 = (1,3)$, $v_5 = \{(2,3), (3,2)\}$, $v_6 = (3,1)$, $v_7 = (3,0)$.

Following similar steps for \hat{P}_1^2 , we have

$$\hat{P}_1^1 \doteq \rho^{-1}, \quad \hat{P}_2^1 \doteq \rho^{-2}, \quad \hat{P}_3^1 \doteq \rho^{-1}, \quad (36)$$

where we use “ \doteq ” rather than “ \leq ” in (35), since all the states in \hat{S}_2 and \hat{S}_3 are the edge states.

B. Stationary Probabilities

Based on (35), (36) and the fact that $\hat{P}_1^3 = \hat{P}_3^1 = 0$, it is easy to obtain the stationary probabilities of \hat{S}_i , denoted by $\hat{\pi}_i$, $i \in [1:3]$, as follows:

$$\hat{\pi}_1 = \frac{\hat{P}_2^1 \hat{P}_3^2}{\theta}, \quad \hat{\pi}_2 = \frac{\hat{P}_1^2 \hat{P}_3^2}{\theta}, \quad \hat{\pi}_3 = \frac{\hat{P}_1^1 \hat{P}_2^3}{\theta}, \quad (37)$$

where $\theta \triangleq \hat{P}_2^1 \hat{P}_3^2 + \hat{P}_1^2 \hat{P}_3^2 + \hat{P}_1^1 \hat{P}_2^3$.

Since $\hat{P}_1^2 \hat{P}_3^2 \doteq \rho^{-2}$, $\hat{P}_1^2 \hat{P}_3^2 \leq \rho^{-3}$, and $\hat{P}_1^1 \hat{P}_2^3 \leq \rho^{-4}$ as shown in (35) and (36), $\theta \doteq \rho^{-2}$ holds. Thus,

$$\hat{\pi}_1 \doteq \rho^0, \quad \hat{\pi}_2 \leq \rho^{-1}, \quad \hat{\pi}_3 \leq \rho^{-2}. \quad (38)$$

C. Diversity Order

Firstly, from (23), $(P_{(1,0)}^{(1,0)} - \phi_0)$ can be approximated as follows:

$$\begin{aligned} P_{(1,0)}^{(1,0)} - \phi_0 &= (1 - \phi_0 - \phi_1 - \phi_2)(1 - \psi_0 - \psi_1) \\ &\approx \frac{\zeta_1 \zeta_2 \max(\xi_{u1}, \xi_{u2})}{\Omega_1 \Omega_2 \Omega_3}. \end{aligned} \quad (39)$$

Thus, we have $P_{(1,0)}^{(1,0)} - \phi_0 \doteq \rho^{-3}$. Similarly, it can be proved that

$$P_s^s - \phi_0 \leq \rho^{i-4}, \quad \forall i \in [1:3], \forall s \in \hat{S}_i. \quad (40)$$

Secondly, from (30), (38) and (40), the SOP of the proposed buffer-aided C-NOMA scheme can be expressed as follows:

$$\begin{aligned} P_{out} &= \sum_{i \in [1:3]} \sum_{s \in \hat{S}_i} (P_s^s - \phi_0) \pi_s \\ &\doteq \sum_{i \in [1:3]} \rho^{i-4} \sum_{s \in \hat{S}_i} \pi_s = \sum_{i \in [1:3]} \rho^{i-4} \hat{\pi}_i \\ &\leq \rho^{-3} \rho^0 + \rho^{-2} \rho^{-1} + \rho^{-1} \rho^{-2} \doteq \rho^{-3}. \end{aligned} \quad (41)$$

From the definition of diversity order in (31), the diversity order of 3 can be achieved by the proposed buffer-aided C-NOMA scheme.

APPENDIX C PROOF OF PROPOSITION 3

The system delay includes the delay at the AP and the delay at the users. Thus, we first need to analyze the transmit throughputs at the AP and the users, respectively. Based on the law of the conservation of flow, the sum bits received by the users should be equal to the sum bits transmitted by the users over a long period, so

$$F_1 + F_2 = F_3 + F_4. \quad (42)$$

Moreover,

$$F_1 + F_2 + F_3 + F_4 + P_{out} + \phi_0 = 1. \quad (43)$$

Thus, we have

$$F_1 + F_2 = F_3 + F_4 = \frac{1 - P_{out} - \phi_0}{2}. \quad (44)$$

Denote η_U as the sum transmit throughput at the two users, which can be expressed as

$$\eta_U = 2R_0(F_3 + F_4) = R_0(1 - P_{out} - \phi_0). \quad (45)$$

Denote η_S as the sum transmit throughput at the AP. The first part of η_S corresponds to mode \mathcal{M}_0 , denoted by $\eta_S^{(1)}$. Obviously, $\eta_S^{(1)}$ is equal to $2R_0\phi_0$. The second part of η_S corresponds to modes \mathcal{M}_1 and \mathcal{M}_2 , denoted by $\eta_S^{(2)}$. Note that $\eta_S^{(2)}$ should be twice of η_U in (45), since both the two users' packets are transmitted by the AP in modes \mathcal{M}_1 and \mathcal{M}_2 , while only one user's packets are transmitted in modes \mathcal{M}_3 and \mathcal{M}_4 . In summary, η_S can be expressed as

$$\eta_S = 2R_0\phi_0 + 2\eta_U = 2R_0(1 - P_{out}). \quad (46)$$

A. Average Delay at the AP

Denote Q_S as the average sum queuing length at the AP, which is determined by how fast the packets are delivered by the AP. Thus, Q_S is given by

$$Q_S = 2R_0[1 - (F_1 + F_2) - \phi_0] = R_0(1 + P_{out} - \phi_0). \quad (47)$$

According to Little's law [41], the average delay at the AP is

$$D_S = \frac{Q_S}{\eta_S} = \frac{1 + P_{out} - \phi_0}{2(1 - P_{out})}. \quad (48)$$

B. Average Delay at the Users

On the other hand, the average sum queuing length at the users is give by

$$Q_U = 2R_0 \sum_{l_1=0}^L \sum_{l_2=0}^L (l_1 + l_2) \pi_{(l_1, l_2)}. \quad (49)$$

Thus, the average delay at the users is

$$D_U = \frac{Q_U}{\eta_U} = \frac{2 \sum_{l_1=0}^L \sum_{l_2=0}^L (l_1 + l_2) \pi_{(l_1, l_2)}}{1 - P_{out} - \phi_0}. \quad (50)$$

C. Total Average Delay of the System

The total average delay of the system can be defined as follows:

$$\bar{D} = (1 - r_c)D_S + r_c(D_S + D_U) = D_S + r_c D_U, \quad (51)$$

where r_c denotes the percentage between the cooperative throughput (i.e., the throughput corresponding to the cooperative modes $\{\mathcal{M}_k\}_{k \in [1:4]}$) and the total throughput, which is given by

$$r_c = \frac{\eta_U}{2R_0(1 - P_{out})} = \frac{1 - P_{out} - \phi_0}{2(1 - P_{out})}. \quad (52)$$

Accordingly, the percentage between the direct throughput (i.e., the throughput corresponding to the direct mode \mathcal{M}_0) and the total throughput is $1 - r_c$.

In summary, from (48), (50), (51) and (52), the total average delay of the system can be obtained shown in (32).

REFERENCES

- [1] G. Chen, J. Tang, and J. P. Coon, "Optimal routing for multihop social-based 2D communications in the Internet of Things," *IEEE Internet of Things J.*, vol. 5, no. 3, pp. 1880–1889, Jun. 2018.
- [2] G. Chen, J. P. Coon, A. Mondal, B. Allen, and J. A. Chambers, "Performance analysis for multihop full-duplex IoT networks subject to poisson distributed interferers," *IEEE Internet of Things J.*, vol. 6, no. 2, pp. 3467–3479, Apr. 2019.
- [3] Z. Ding, Z. Yang, P. Fan, and H. V. Poor, "On the performance of non-orthogonal multiple access in 5G systems with randomly deployed users," *IEEE Signal Process. Lett.*, vol. 21, no. 12, pp. 1501–1505, Dec. 2014.
- [4] Z. Ding, Y. Liu, J. Choi, Q. Sun, M. Elkashlan, C. 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.
- [5] Z. Ding *et al.*, "A survey on non-orthogonal multiple access for 5G networks: Research challenges and future trends," *IEEE J. Sel. Areas in Commun.*, vol. 35, no. 10, pp. 2181–2195, Oct. 2017.
- [6] T. Lv, Y. Ma, J. Zeng, and P. T. Mathiopoulos, "Millimeter-wave NOMA transmission in cellular M2M communications for Internet of Things," *IEEE Internet of Things J.*, vol. 5, no. 3, pp. 1989–2000, Jun. 2018.
- [7] M. Liu, T. Song, and G. Gui, "Deep cognitive perspective: Resource allocation for NOMA-Based heterogeneous IoT with imperfect SIC," *IEEE Internet of Things J.*, vol. 6, no. 2, pp. 2885–2894, Apr. 2019.
- [8] L. P. Qian, A. Feng, Y. Huang, Y. Wu, B. Ji, and Z. Shi, "Optimal SIC ordering and computation resource allocation in MEC-Aware NOMA NB-IoT networks," *IEEE Internet of Things J.*, vol. 6, no. 2, pp. 2806–2816, Apr. 2019.
- [9] M. Alkhawatrah, Y. Gong, G. Chen, S. Lambbotharan, and J. A. Chambers, "Buffer-aided relay selection for cooperative NOMA in the Internet of things," *IEEE Internet of Things J.*, vol. 6, no. 3, pp. 5722–5731, Jun. 2019.
- [10] J. Kim and I. Lee, "Non-orthogonal multiple access in coordinated direct and relay transmission," *IEEE Commun. Lett.*, vol. 19, no. 11, pp. 2037–2040, Nov. 2015.
- [11] J. Men and J. Ge, "Non-orthogonal multiple access for multiple-antenna relaying networks," *IEEE Commun. Lett.*, vol. 19, no. 10, pp. 1686–1689, Oct. 2015.
- [12] Z. Yang, Z. Ding, Y. Wu, and P. Fan, "Novel relay selection strategies for cooperative NOMA," *IEEE Trans. Veh. Tech.*, vol. 66, no. 11, pp. 10114–10123, Nov. 2017.
- [13] P. Xu, Z. Yang, Z. Ding, and Z. Zhang, "Optimal relay selection schemes for cooperative NOMA," *IEEE Trans. Veh. Tech.*, vol. 67, no. 8, pp. 7851–7855, Aug. 2018.
- [14] Z. Ding, M. Peng, and H. V. Poor, "Cooperative non-orthogonal multiple access in 5G systems," *IEEE Commun. Lett.*, vol. 19, no. 8, pp. 1462–1465, Aug. 2015.
- [15] Z. Zhang, Z. Ma, M. Xiao, Z. Ding, and P. Fan, "Full-duplex device-to-device-aided cooperative non-orthogonal multiple access," *IEEE Trans. Veh. Tech.*, vol. 66, no. 5, pp. 4467–4471, May 2016.
- [16] K. Janghel and S. Prakriya, "Performance of adaptive OMA/cooperative-NOMA scheme with user selection," *IEEE Commun. Lett.*, vol. 22, no. 10, pp. 2092–2095, Oct. 2018.
- [17] Q. Y. Liao and C. Y. Leow, "Successive user relaying in cooperative NOMA system," *IEEE Wireless Commun. Lett.*, vol. 8, no. 3, pp. 921–924, Jun. 2019.
- [18] Y. Liu, Z. Ding, M. Elkashlan, and H. V. Poor, "Cooperative non-orthogonal multiple access with simultaneous wireless information and power transfer," *IEEE J. Sel. Areas in Commun.*, vol. 34, no. 4, pp. 938–953, Apr. 2016.
- [19] I. Krikidis, T. Charalambous, and J. S. Thompson, "Buffer-aided relay selection for cooperative diversity systems without delay constraints," *IEEE Trans. Wireless Commun.*, vol. 11, no. 5, pp. 1957–1967, May 2012.
- [20] T. Charalambous, N. Nomikos, I. Krikidis, D. Vouyioukas, and M. Johansson, "Modeling buffer-aided relay selection in networks with direct transmission capability," *IEEE Commun. Lett.*, vol. 19, no. 4, pp. 649–652, Apr. 2015.
- [21] Z. Tian, G. Chen, Y. Gong, Z. Chen, and J. A. Chambers, "Buffer-aided max-link relay selection in amplify-and-forward cooperative networks," *IEEE Trans. Veh. Tech.*, vol. 64, no. 2, pp. 553–565, Feb. 2015.
- [22] G. Chen, Z. Tian, Y. Gong, and J. A. Chambers, "Decode-and-forward buffer-aided relay selection in cognitive relay networks," *IEEE Trans. Veh. Tech.*, vol. 63, no. 9, pp. 4723–4728, Nov. 2014.
- [23] N. Nomikos, T. Charalambous, I. Krikidis, D. N. Skoutas, D. Vouyioukas, and M. Johansson, "A buffer-aided successive opportunistic relay selection scheme with power adaptation and inter-relay interference cancellation for cooperative diversity systems," *IEEE Trans. Commun.*, vol. 63, no. 5, pp. 1623–1634, May 2015.
- [24] P. Xu, Z. Ding, I. Krikidis, and X. Dai, "Achieving optimal diversity gain in buffer-aided relay networks with small buffer size," *IEEE Trans. Veh. Tech.*, vol. 65, no. 10, pp. 8788–8794, Oct. 2016.
- [25] S. Luo and K. C. Teh, "Buffer state based relay selection for buffer-aided cooperative relaying systems," *IEEE Trans. Wireless Commun.*, vol. 14, no. 10, pp. 5430–5439, Oct. 2015.
- [26] Z. Tian, Y. Gong, G. Chen, and J. A. Chambers, "Buffer-aided relay selection with reduced packet delay in cooperative networks," *IEEE Trans. Veh. Tech.*, vol. 66, no. 3, pp. 2567–2575, Mar. 2017.
- [27] D. Sui, F. Hu, W. Zhou, M. Shao, and M. Chen, "Relay selection for radio frequency energy-harvesting wireless body area network with buffer," *IEEE Internet of Things J.*, vol. 5, no. 2, pp. 1100–1107, Apr. 2018.
- [28] S. Luo and K. C. Teh, "Adaptive transmission for cooperative NOMA system with buffer-aided relaying," *IEEE Commun. Lett.*, vol. 21, no. 4, pp. 937–940, Apr. 2017.
- [29] Q. Zhang, Z. Liang, Q. Li, and J. Qin, "Buffer-aided non-orthogonal multiple access relaying systems in Rayleigh fading channels," *IEEE Trans. Commun.*, vol. 65, no. 1, pp. 95–106, Jan. 2017.
- [30] P. Xu, J. Quan, Z. Yang, G. Chen, and Z. Ding, "Performance analysis of buffer-aided hybrid NOMA/OMA in cooperative uplink system," *IEEE Access*, vol. 7, pp. 168759–168773, Nov. 2019.
- [31] X. Lan, Y. Zhang, Q. Chen, and L. Cai, "Energy efficient buffer-aided transmission scheme in wireless powered cooperative NOMA relay network," *IEEE Trans. Commun.*, vol. 68, no. 3, pp. 1432–1447, 2020.
- [32] J. Li, X. Lei, P. D. Diamantoulakis, F. Zhou, P. Sarigiannidis, and G. K. Karagiannidis, "Resource allocation in buffer-aided cooperative non-orthogonal multiple access systems," *IEEE Trans. Commun.*, doi:10.1109/TCOMM.2020.3023458.
- [33] N. Nomikos, T. Charalambous, D. Vouyioukas, R. Wichman, and G. K. Karagiannidis, "Integrating broadcasting and NOMA in full-duplex buffer-aided opportunistic relay networks," *IEEE Trans. Veh. Technol.*, vol. 69, no. 8, pp. 9157–9162, 2020.

- [34] N. Nomikos, T. Charalambous, D. Vouyioukas, G. K. Karagiannidis, and R. Wichman, "Hybrid NOMA/OMA with buffer-aided relay selection in cooperative networks," *IEEE J. Sel. Topics in Signal Process.*, vol. 13, no. 3, pp. 524–537, Jun. 2019.
- [35] N. Zlatanov, R. Schober, and L. Lampe, "Buffer-aided relaying in a three node network," in *IEEE Int. Symp. Info. Theory Proc.*, Cambridge, MA, USA, Jul. 2012, pp. 781–785.
- [36] H. Nasir, N. Javaid, and W. Raza, "Study of buffer-aided cooperative NOMA using incremental relaying in wireless networks," *Physical Commun.*, vol. 39, p. 101011, Apr. 2020.
- [37] Z. Ding, P. Fan, and H. V. Poor, "Impact of user pairing on 5G nonorthogonal multiple-access downlink transmissions," *IEEE Trans. Veh. Technol.*, vol. 65, no. 8, pp. 6010–6023, Aug. 2016.
- [38] Y. Sun, D. W. K. Ng, Z. Ding, and R. Schober, "Optimal joint power and subcarrier allocation for full-duplex multicarrier non-orthogonal multiple access systems," *IEEE Trans. Commun.*, vol. 65, no. 3, pp. 1077–1091, 2017.
- [39] J. R. Norris and J. R. Norris, *Markov chains*. Cambridge university press, 1997, no. 2.
- [40] K. Azarian, H. El Gamal, and P. Schniter, "On the achievable diversity-multiplexing tradeoff in half-duplex cooperative channels," *IEEE Trans. Inf. Theory*, vol. 51, no. 12, pp. 4152–4172, Dec. 2005.
- [41] J. D. Little and S. C. Graves, "Little's law," in *Building intuition*. Springer, 2008, pp. 81–100.



Peng Xu (Member, IEEE) received the B.Eng. and the Ph.D. degrees in electronic and information engineering from the University of Science and Technology of China, Anhui, China, in 2009 and 2014, respectively. From June 2014 to July 2016, he was working as a postdoctoral researchers with the Department of Electronic Engineering and Information Science, University of Science and Technology of China, Hefei, China. He is currently an Associated Professor with the School of Communication and Information Engineering, Chongqing University of

Posts and Telecommunications (CQUPT), Chongqing, China. His current research interests include cooperative communications, information-theoretic secrecy, NOMA techniques and reconfigurable intelligent surface. Dr. Peng Xu received IEEE WIRELESS COMMUNICATIONS LETTERS Exemplary Reviewer 2015 and Excellent Paper of Chongqing Association for Science and Technology 2018.



Yunwu Wang (Student Member, IEEE) received the B.Sc. degree in Department of Optical Information Science and Technology from the College of Post and Telecommunication of Wuhan Institute of Technology, Wuhan, China, in 2018. He is currently pursuing the M.Sc. degree in electronic and communication engineering at the School of Communication and Information Engineering, Chongqing University of Posts and Telecommunications (CQUPT), Chongqing, China. His research interests include cooperative communication, non-orthogonal

multiple access (NOMA), and buffer-aided relaying.



Gaojie Chen (Senior Member, IEEE) received the B.Eng. and B.Ec. Degrees in electrical information engineering and international economics and trade from Northwest University, China, in 2006, and the M.Sc. (Hons.) and Ph.D. degrees in electrical and electronic engineering from Loughborough University, Loughborough, U.K., in 2008 and 2012, respectively. From 2008 to 2009, he was a Software Engineer with DT Mobile, Beijing, China. From 2012 to 2013, he was a Research Associate with the School of Electronic, Electrical and Systems Engineering, Loughborough University. He was a Research Fellow with 5GIC, Faculty of Engineering and Physical Sciences, University of Surrey, U.K., from 2014 to 2015. He was also a Research Associate with the Department of Engineering Science, University of Oxford, U.K., from 2015 to 2018. He is currently a Lecturer with the School of Engineering, University of Leicester, U.K. His current research interests include information theory, wireless communications, cooperative communications, cognitive radio, the Internet of Things, secrecy communications, and random geometric networks. He received the Exemplary Reviewer Certificates of the IEEE WIRELESS COMMUNICATIONS LETTERS in 2018 and the IEEE TRANSACTIONS ON COMMUNICATIONS in 2019. He serves as an Associate Editor for the IEEE COMMUNICATIONS LETTERS, IEEE WIRELESS COMMUNICATIONS LETTERS, IEEE JOURNAL ON SELECTED AREAS IN COMMUNICATIONS - MACHINE LEARNING IN COMMUNICATIONS AND NETWORKS and *Electronics Letters* (IET).



Gaofeng Pan (Senior Member, IEEE) received his B.Sc. in Communication Engineering from Zhengzhou University, Zhengzhou, China, in 2005, and the Ph.D. degree in Communication and Information Systems from Southwest Jiaotong University, Chengdu, China, in 2011. Since August 2019, he has been a Visiting Researcher with the Communication Theory Lab, King Abdullah University of Science and Technology (KAUST), Thuwal, Saudi Arabia. He has also been a Professor with school of cyberspace science and technology, Beijing Institute of Technology, China, since Nov. 2020. His research interests include communications theory, signal processing, and protocol design.



Zhiguo Ding (Fellow, IEEE) received his B.Eng in Electrical Engineering from the Beijing University of Posts and Telecommunications in 2000, and the Ph.D degree in Electrical Engineering from Imperial College London in 2005. From Jul. 2005 to Apr. 2018, he was working in Queen's University Belfast, Imperial College, Newcastle University and Lancaster University. Since Apr. 2018, he has been with the University of Manchester as a Professor in Communications. From Oct. 2012 to Sept. 2021, he has also been an academic visitor in Princeton

University.

Dr Ding' research interests are B5G networks, machine learning, cooperative and energy harvesting networks and statistical signal processing. He is serving as an Area Editor for the IEEE OPEN JOURNAL OF THE COMMUNICATIONS SOCIETY, an Editor for it IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY, and IEEE OPEN JOURNAL OF THE SIGNAL PROCESSING SOCIETY, and was an Editor for IEEE TRANSACTIONS ON COMMUNICATIONS, IEEE WIRELESS COMMUNICATION LETTERS, and IEEE COMMUNICATION LETTERS. He received the best paper award in IET ICWMC-2009 and IEEE WCSP-2014, the EU Marie Curie Fellowship 2012-2014, the Top IEEE TVT Editor 2017, IEEE Heinrich Hertz Award 2018, IEEE Jack Neubauer Memorial Award 2018, IEEE Best Signal Processing Letter Award 2018, and Friedrich Wilhelm Bessel Research Award 2020. He is a Fellow of the IEEE and Web of Science Highly Cited Researcher in two categories 2020.