

Stable Throughput Regions of Opportunistic NOMA and Cooperative NOMA With Full-Duplex Relaying

Yong Zhou¹, Member, IEEE, Vincent W.S. Wong², Fellow, IEEE, and Robert Schober, Fellow, IEEE

Abstract—In this paper, we consider downlink non-orthogonal multiple access (NOMA) transmission with dynamic traffic arrival for spatially random users of different priorities. By exploiting limited channel state information, we propose an opportunistic NOMA scheme to enable NOMA for high- and low-priority users when high-priority users experience good channel conditions. Opportunistic NOMA improves the transmission opportunities of low-priority users while reducing the adverse effect of NOMA on high-priority users. Moreover, we propose a cooperative NOMA scheme with full-duplex relaying, where low-priority users act as full-duplex relays to assist the high-priority users. The high-priority user constructively combines the signal and its delayed version transmitted by the base station and a selected relay, respectively. The adopted relay selection scheme takes into account the users' spatial distribution, queue status, and channel conditions. By using tools from queuing theory and stochastic geometry, we derive the stable throughput regions of both proposed schemes. Furthermore, we derive the conditions under which the proposed NOMA schemes achieve larger stable throughput regions than orthogonal multiple access (OMA). At the expense of a higher implementation complexity and with appropriate parameter setting, cooperative NOMA with full-duplex relaying achieves a larger stable throughput region than opportunistic NOMA, which in turn outperforms OMA.

Index Terms—Non-orthogonal multiple access, stable throughput, dynamic traffic arrival, full-duplex relaying, spatially random users.

I. INTRODUCTION

TO MEET the rapidly increasing traffic demand caused by the proliferation of mobile devices and data intensive applications, non-orthogonal multiple access (NOMA) [2] has been proposed as a promising technique to enhance the spectral efficiency of the fifth generation (5G)

Manuscript received February 9, 2017; revised October 3, 2017; accepted May 2, 2018. This work was supported by the Natural Sciences and Engineering Research Council of Canada. This paper was presented in part at the IEEE International Conference on Communications, Paris, France, May 2017 [1]. The associate editor coordinating the review of this paper and approving it for publication was R. K. Ganti. (*Corresponding author: Vincent W.S. Wong.*)

Y. Zhou was with the Department of Electrical and Computer Engineering, The University of British Columbia, Vancouver, BC V6T 1Z4, Canada. He is now with the School of Information Science and Technology, ShanghaiTech University, Shanghai 201210, China (e-mail: zhouyong@shanghaitech.edu.cn).

V. W.S. Wong is with the Department of Electrical and Computer Engineering, The University of British Columbia, Vancouver, BC V6T 1Z4, Canada (e-mail: vincentw@ece.ubc.ca).

R. Schober is with the Institute for Digital Communications, Friedrich-Alexander University of Erlangen–Nuremberg, 91058 Erlangen, Germany (e-mail: robert.schober@fau.de).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TWC.2018.2837014

cellular network. With NOMA, multiple users can simultaneously be served by exploiting the power domain rather than the time and frequency domains as in orthogonal multiple access (OMA). By appropriately allocating the transmit power at the base station to multiple users with diverse channel conditions, NOMA can achieve a balance between network throughput and user fairness.

NOMA has recently received considerable research interest [3]–[9]. Specifically, the system-level performance of downlink NOMA transmission is evaluated in [3], which shows that transmit power allocation and user pairing are two important design aspects of NOMA. An optimal power allocation strategy is proposed in [4] to maximize the sum rate of multiple-input multiple-output (MIMO) NOMA networks. The authors in [5] formulate a joint transmit power and subcarrier allocation problem for maximization of the sum rate of multi-carrier NOMA networks and solve the problem using matching theory. The impact of user pairing on the performance of NOMA is investigated in [6], which shows that NOMA achieves a better performance when the paired NOMA users experience more distinct channel conditions. The authors in [7] derive the outage probability of MIMO-NOMA for both uplink and downlink transmission. In addition, the outage probability of a cooperative NOMA scheme is analyzed in [8], where a relay is selected to forward packets to paired NOMA users having different priorities and the low-priority user is served in an opportunistic manner. However, all of the aforementioned studies focus on resource allocation and performance analysis for NOMA with backlogged traffic.

Full-duplex communication can enhance the spectral efficiency by allowing the radios to simultaneously transmit and receive on the same frequency channel. The main challenge for realizing full-duplex communication is the self-interference due to signal leakage, which significantly degrades the performance gain achieved by full-duplexing [10]. Nevertheless, with the advancement of analog and digital self-interference cancelation techniques, full-duplex radios have been successfully implemented [11]. The rate region of full-duplex links in orthogonal frequency division multiplexing systems is analyzed in [12]. The authors in [13] develop a joint power and subcarrier allocation policy to maximize the weighted sum throughput of multi-carrier NOMA systems, where the full-duplex base station simultaneously serves multiple uplink and downlink users. Furthermore, full-duplex relaying has recently attracted significant interest [14], [15]. The authors in [14] compare the spectral efficiency of half- and full-duplex relaying strategies, and propose a joint opportunistic mode

selection and transmit power adaptation scheme to optimize spectral efficiency. However, the performance of full-duplex relaying in NOMA systems with dynamic traffic arrival and spatially random relays has not been studied yet.

Different from the aforementioned studies, we consider downlink NOMA transmission with dynamic traffic arrival and spatially random users of different priorities. For dynamic traffic arrival, the *stable throughput region* [16]–[19] is an important performance metric and defined as the set of achievable packet arrival rates given that all queues are stable. However, according to the NOMA principle, a low-priority user is allowed to share the frequency channel and transmit power with a high-priority user, which may reduce the reception reliability of the high-priority user and lead to queue instability. In NOMA, the low-priority user, which is allocated a lower transmit power, needs to decode the signal intended for the high-priority user first before decoding its own signal. Hence, the low-priority user can act as a relay and assist the transmission of the high-priority user. However, when half-duplex relaying is used, an additional time slot is required for packet forwarding, which reduces the spectral efficiency. Full-duplex relaying has the potential to mitigate this disadvantage. The performance gain achieved by full-duplex relaying can be further improved by relay selection, where the selection should take into account the residual self-interference, the queue status, and the spatial distribution of the potential relays. Considering dynamic traffic arrival together with NOMA leads to interacting queues, which complicates the performance analysis. In particular, the service process of a given queue depends on the status of the other queue, as the status of both queues determines whether NOMA can be enabled. Furthermore, channel state information (CSI) plays an important role in designing user pairing and transmit power allocation strategies. As full CSI is difficult to obtain in practice, the impact of limited CSI [20] on the performance of NOMA should be investigated.

To address the aforementioned issues, we first propose an *opportunistic NOMA* scheme exploiting limited CSI, where NOMA for high- and low-priority users is enabled only if the channel gain between the base station and the high-priority user does not fall below a certain threshold. NOMA for the low-priority users is also enabled by exploiting the differences of the low-priority users' distances to the base station. By appropriately setting the threshold to trigger NOMA, the opportunistic NOMA scheme improves the transmission opportunities of the low-priority users without degrading the performance of the high-priority users. Furthermore, we propose a *cooperative NOMA scheme with full-duplex relaying*, where the low-priority users act as full-duplex relays to help forward packets to the high-priority users. By exploiting cooperative diversity to enhance the probability of successful packet reception at the high-priority users, the number of packet retransmissions for the high-priority users is reduced, which in turn further improves the transmission opportunities of the low-priority users. The main contributions of this paper are summarized as follows:

- We develop a theoretical performance analysis framework for downlink NOMA transmission with dynamic traffic arrival

and spatially random users of different priorities. This analytical framework provides a better understanding of the benefits and limitations of NOMA.

- We decouple the interacting queues caused by dynamic traffic arrival and NOMA by allowing empty queues to contribute dummy packets. Tools from queueing theory and stochastic geometry are applied to characterize the stable throughput region of opportunistic NOMA.

- We derive the stable throughput region of cooperative NOMA with full-duplex relaying, taking into account the residual self-interference, spatially random low-priority users, and relay selection. Studying both opportunistic NOMA and cooperative NOMA with full-duplex relaying provides insights regarding the tradeoff between network performance and implementation complexity. We also derive the conditions under which the proposed NOMA schemes achieve larger stable throughput regions than OMA.

- Simulation results validate the analysis of the probabilities of successful packet reception. Numerical results show that, with appropriate parameter setting, both proposed NOMA schemes can outperform OMA, and cooperative NOMA with full-duplex relaying can achieve a larger stable throughput region than opportunistic NOMA at the expense of a higher implementation complexity. The impact of the relevant design and system parameters (e.g., the threshold to trigger NOMA and the power allocation coefficients) on the stable throughput regions of the proposed NOMA schemes is also evaluated.

The rest of this paper is organized as follows. We describe the system model in Section II. In Section III, we present the opportunistic NOMA scheme and derive its stable throughput region. We describe the cooperative NOMA scheme with full-duplex relaying and characterize its stable throughput region in Section IV. In Section V, we present the conditions under which the proposed NOMA schemes achieve larger stable throughput regions than OMA. Numerical results are provided in Section VI. Finally, Section VII concludes this paper.

II. SYSTEM MODEL

Consider a downlink transmission scenario consisting of one base station and multiple users, as shown in Fig. 1(a). Base station S is located at the center of the circular network coverage area with radius r . Over a single frequency channel, time is divided into slots of constant durations. Users are categorized into two groups with different priorities, i.e., K low-priority users in set \mathcal{U}^L and M high-priority users in set \mathcal{U}^H . The locations of low-priority users are assumed to follow a binomial point process (BPP) [21], [22]. Specifically, for each time slot, K low-priority users are independently and uniformly distributed within the network coverage area. On the other hand, the high-priority users are located r_H meters away¹

¹The proposed framework can be extended to the case where the high-priority users also have random distances to the base station by first conditioning on the distance and then taking the expectation over the high-priority user distance distribution. The resulting analytical expressions involve an additional integral compared to the results obtained for the fixed high-priority user distance considered in this paper. Fixed user distances were also assumed in other works in the literature, e.g., [23]–[25], as this approach simplifies the analytical expressions without compromising the insights that can be obtained, as demonstrated in [26].

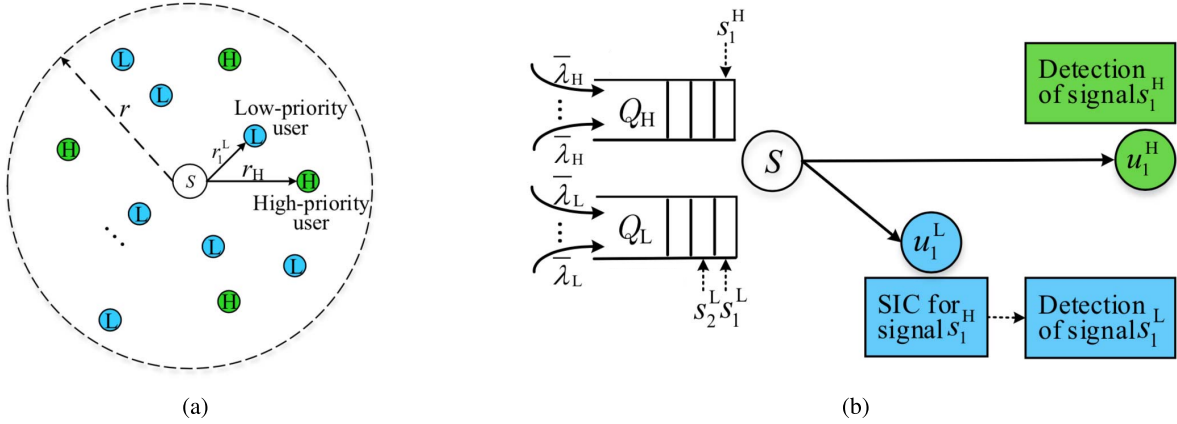


Fig. 1. (a) Illustration of the network topology for downlink transmission with spatially random users of different priorities, where base station S serves M high-priority users and K low-priority users. (b) Illustration of the queueing and signal reception models for downlink transmission with dynamic traffic arrival. Base station S transmits the first packet from queue Q_H and the first packet from queue Q_L to high-priority user u_1^H and low-priority user u_1^L using NOMA, respectively.

188 from base station S in a random direction. Base station S and
 189 all users have a single antenna.

190 Base station S is equipped with two queues of infinite
 191 size, denoted as Q_H and Q_L , which store the packets to be
 192 transmitted to the high- and low-priority users, respectively,
 193 as shown in Fig. 1(b). The packet arrival at base station S
 194 for each user follows an independent and stationary process. For
 195 ease of presentation, the average arrival rates of users having
 196 the same priority are assumed to be identical, but the analysis
 197 can be extended to a general scenario with diverse average
 198 arrival rates. The average arrival rates of queues Q_H and Q_L
 199 are given by $\bar{\lambda}_H = M\bar{\lambda}_H$ and $\bar{\lambda}_L = K\bar{\lambda}_L$ (packets per time
 200 slot), where $\bar{\lambda}_H$ and $\bar{\lambda}_L$ denote the average arrival rates for
 201 each high- and low-priority user, respectively. Packets for users
 202 having the same priority have the same size in bits and are
 203 served in a first-in first-out (FIFO) manner. Each packet is
 204 transmitted in one time slot.

205 The channel between any two transceivers suffers from
 206 path loss and Rayleigh fading. A packet can be successfully
 207 decoded only if the received signal-to-interference-plus-noise
 208 ratio (SINR) is not smaller than a required reception threshold.
 209 Upon successfully (or erroneously) receiving a packet from
 210 base station S , the corresponding receiver sends feedback
 211 that indicates the packet success or failure to base station S
 212 via an error- and delay-free control channel. After successful
 213 reception, the packet is removed from the queue at base
 214 station S . Otherwise, base station S retransmits the packet
 215 until it is successfully decoded. We denote $Q_H(t)$ and $Q_L(t)$
 216 as the queue lengths of Q_H and Q_L in time slot t , respectively.
 217 A queue is said to be stable if its queue length has a limiting
 218 distribution as time goes to infinity. For high-priority queue,
 219 we have $\lim_{t \rightarrow \infty} \mathbb{P}(Q_H(t) < l) = F(l)$ and $\lim_{l \rightarrow \infty} F(l) = 1$.
 220 If the arrival and service processes of a queue are jointly
 221 stationary and ergodic, by Loynes' theorem [27], the sufficient
 222 condition for the stability of queue Q_H is that $\lambda_H < \mu_H$,
 223 where μ_H (packets per time slot) is the average service
 224 rate of queue Q_H . The network is stable when both queues
 225 Q_H and Q_L are stable. In this work, the stable throughput

226 region is defined as the set of arrival rates of queues Q_H and
 227 Q_L that lead to a stable network for fixed power allocation
 228 coefficients and threshold to trigger NOMA. The full stable
 229 throughput region refers to the union of the stable throughput
 230 regions over all possible values of the power allocation coef-
 231 ficients and threshold to trigger NOMA.

232 In order to reduce the implementation complexity,
 233 we consider the case when two users are paired for NOMA
 234 transmission. Such a two-user NOMA scheme is included in
 235 the 3rd Generation Partnership Project (3GPP) standard [28]
 236 and considered in [4] and [6]–[8]. We denote the intended
 237 receivers of the first packet from queue Q_H and the first
 238 packet from queue Q_L by u_1^H and u_1^L , respectively. When
 239 NOMA is performed to serve users u_1^H and u_1^L in time slot t ,
 240 the superimposed signal transmitted by base station S is
 241 $\alpha_H \sqrt{P_S} s_1^H(t) + \alpha_L \sqrt{P_S} s_1^L(t)$, where P_S denotes the transmit
 242 power of base station S , α_H and α_L denote the transmit
 243 power allocation coefficients for the high- and low-priority
 244 users, respectively, and $s_1^H(t)$ and $s_1^L(t)$ denote the signals
 245 intended for users u_1^H and u_1^L in time slot t , respectively,
 246 with $\mathbb{E}(|s_1^H(t)|^2) = \mathbb{E}(|s_1^L(t)|^2) = 1$. Here, $\mathbb{E}(\cdot)$ denotes
 247 statistical expectation. The paired NOMA users are ordered
 248 according to their priorities for being served [8]. As user u_1^H
 249 has a higher priority, we have $\alpha_H > \alpha_L$ and $\alpha_H^2 + \alpha_L^2 = 1$.
 250 Before transmission begins, the base station informs user
 251 u_1^L that it is expected to perform successive interference
 252 cancelation (SIC) by sending a corresponding control informa-
 253 tion, which includes information about the allocated transmit
 254 power and is attached to the user's scheduling information,
 255 as suggested in [28, pp. 15].

256 The superimposed signal received at user u_1^a , $a \in \{H, L\}$,
 257 in time slot t is given by

$$258 y_1^a(t) = (\alpha_H s_1^H(t) + \alpha_L s_1^L(t)) \sqrt{P_S} h_1^a(t) \sqrt{\ell(x_1^a)} + n_1^a(t), \quad 259 \quad (1)$$

260 where $h_1^a(t)$ denotes the Rayleigh fading channel gain between
 261 base station S and user u_1^a in time slot t , $n_1^a(t)$ denotes the

additive white Gaussian noise (AWGN) at user u_1^a with zero mean and variance σ^2 in time slot t , x_1^a denotes the location of user u_1^a , $\ell(x_1^a) = (1 + (r_1^a)^\beta)^{-1}$ and r_1^a denote the non-singular path loss and the distance between base station S and user u_1^a , respectively, and β denotes the path loss exponent. Hence, $|h_1^a(t)|^2$ is an exponential random variable with unit mean.

After receiving the signal from base station S , high-priority user u_1^H treats the signal intended for low-priority user u_1^L as interference and decodes its own signal based on SINR

$$\Gamma_{H1|L1}(t, \alpha_H) = \frac{\alpha_H^2 P_S |h_1^H(t)|^2 \ell(x_1^H)}{\alpha_L^2 P_S |h_1^L(t)|^2 \ell(x_1^L) + \sigma^2}, \quad (2)$$

where $\Gamma_{H1|L1}(t, \alpha_H)$ denotes the SINR of signal $s_1^H(t)$ observed at high-priority user u_1^H when paired with low-priority user u_1^L in time slot t .

Low-priority user u_1^L first decodes the signal intended for high-priority user u_1^H with SINR

$$\Gamma_{H1 \rightarrow L1}(t, \alpha_H) = \frac{\alpha_H^2 P_S |h_1^H(t)|^2 \ell(x_1^H)}{\alpha_L^2 P_S |h_1^L(t)|^2 \ell(x_1^L) + \sigma^2}, \quad (3)$$

where $\Gamma_{H1 \rightarrow L1}(t, \alpha_H)$ denotes the SINR of signal $s_1^H(t)$ observed at user u_1^L in time slot t .

If low-priority user u_1^L successfully decodes signal $s_1^H(t)$, i.e., $\Gamma_{H1 \rightarrow L1}(t, \alpha_H) \geq \Gamma_{th}^H$, where Γ_{th}^H denotes the threshold required to successfully decode the packets intended for the high-priority users, then low-priority user u_1^L removes signal $s_1^H(t)$ from received signal $y_1^L(t)$ by applying SIC, and decodes its own signal with signal-to-noise ratio (SNR)

$$\Gamma_{L1}(t, \alpha_L) = \frac{\alpha_L^2 P_S |h_1^L(t)|^2 \ell(x_1^L)}{\sigma^2}, \quad (4)$$

where $\Gamma_{L1}(t, \alpha_L)$ denotes the SNR of signal $s_1^L(t)$ observed at user u_1^L in time slot t .

When NOMA is enabled, users u_1^H and u_1^L can successfully decode their own signals if events $\{\Gamma_{H1|L1}(t, \alpha_H) \geq \Gamma_{th}^H\}$ and $\{\Gamma_{H1 \rightarrow L1}(t, \alpha_H) \geq \Gamma_{th}^H, \Gamma_{L1}(t, \alpha_L) \geq \Gamma_{th}^L\}$ occur, respectively, where Γ_{th}^L denotes the threshold required to successfully decode the packets intended for the low-priority users. Base station S simultaneously serves users u_1^H and u_1^L , at the cost of reducing the probability of successful packet reception at high-priority user u_1^H . Specifically, by sharing the frequency channel and transmit power, the received SINR at high-priority user u_1^H decreases, i.e., $\Gamma_{H1|L1}(t, \alpha_H) < \Gamma_{H1}(t, 1) = P_S |h_1^H(t)|^2 \ell(x_1^H) / \sigma^2$. Hence, to guarantee the stability of queue Q_H , NOMA cannot always be enabled, especially when the average arrival rate λ_H is large.

To facilitate our analysis, for the remainder of this paper, we make the following assumptions. The protocol overhead due to feedback from the users to the base station is much smaller than the packet size and is neglected. The fading coefficients are assumed to remain invariant during one time slot and vary independently over different time slots and across different links, as in [23]–[26] and [29]. At the end of each time slot $t \in \mathbb{Z}^+$, the locations of the low-priority users change according to a high mobility random walk model within

the network coverage area as in [23]–[25] and [29]. Hence, the displacement theorem [30] can be applied and the user locations are independent across time slots, which enables the derivation of tractable performance results, providing useful insights on the network performance.

III. OPPORTUNISTIC NOMA

In this section, we propose an opportunistic NOMA scheme to improve the transmission opportunities of low-priority users while reducing the adverse effect of NOMA on high-priority users, and characterize the stable throughput region. We assume that only limited instantaneous CSI is available at base station S . First, when queue Q_H is non-empty in time slot t , one bit of information is sent back from high-priority user u_1^H to base station S . In particular, high-priority user u_1^H sends feedback 1 to base station S if the instantaneous channel gain, $|h_1^H(t)|^2 \ell(x_1^H)$, is not less than a threshold, θ , and sends feedback 0 to base station S otherwise. Second, when queue Q_H is empty in time slot t , users u_1^L and u_2^L send back their distances to base station S , where u_2^L denotes the intended receiver of the second packet from queue Q_L when available. Based on the limited CSI, NOMA is enabled by base station S in an opportunistic manner.

We denote the opportunistic NOMA system as Φ^{ON} , where base station S transmits the first packet from queue Q_H whenever it is non-empty due to its high priority to be served. The packet transmissions depend on the status of queues Q_H and Q_L , and are discussed in the following.

Case 1: If $Q_H(t) > 0$ and $Q_L(t) > 0$, then base station S transmits the first packet from queue Q_H and the first packet from queue Q_L to users u_1^H and u_1^L , respectively, using NOMA with fixed power allocation coefficients (α_H^2, α_L^2) when $|h_1^H(t)|^2 \ell(x_1^H) \geq \theta$, and transmits the first packet from queue Q_H to user u_1^H using OMA with power P_S when $|h_1^H(t)|^2 \ell(x_1^H) < \theta$.

Case 2: If $Q_H(t) > 0$ and $Q_L(t) = 0$, then base station S transmits the first packet from queue Q_H to user u_1^H using OMA² with power P_S .

Case 3: If $Q_H(t) = 0$ and $Q_L(t) > 0$, then base station S transmits the first and second packets from queue Q_L to users u_1^L and u_2^L , respectively, using NOMA when the first two packets are intended for different users (i.e., $u_1^L \neq u_2^L$), and transmits the first packet from queue Q_L to user u_1^L using OMA with power P_S when the first two packets are intended for the same user (i.e., $u_1^L = u_2^L$) or $Q_L(t) = 1$.

The average service rate of queue Q_H depends on the status of queue Q_L . When queue Q_L is empty, base station S transmits the first packet from queue Q_H to user u_1^H using OMA. When queue Q_L is non-empty, base station S transmits the first packet from queue Q_H to user u_1^H using NOMA with probability $\mathbb{P}\left(|h_1^H(t)|^2 \ell(x_1^H) \geq \theta\right) = \exp\left(-\theta\left(1 + r_1^\beta\right)\right)$. Similarly, the average service rate of queue Q_L also depends

²Note that NOMA for different high-priority users is not enabled in this paper, as the probability that different high-priority users experience very different channel conditions is low. However, different user channel conditions are crucial for achieving a gain with NOMA [6]. A similar setting is also considered in [7].

on the status of queue Q_H . Hence, queues Q_H and Q_L interact with each other and their average service rates cannot be directly calculated. In this context, stochastic dominance [31] is a useful tool and can be used to decouple the interacting queues and to characterize the stable throughput region. By using stochastic dominance, we construct two *dominant systems* Φ_1^{ON} and Φ_2^{ON} based on the original opportunistic NOMA system Φ^{ON} . In the following, we derive the stable throughput regions of dominant systems Φ_1^{ON} and Φ_2^{ON} , and then show that the stable throughput region of the original opportunistic NOMA system Φ^{ON} is equal to the union of the stable throughput regions of dominant systems Φ_1^{ON} and Φ_2^{ON} .

A. Stable Throughput Region of Dominant System Φ_1^{ON}

Dominant system Φ_1^{ON} : If queue Q_L is empty, then queue Q_L contributes a dummy packet when high-priority user u_1^H sends feedback 1 to base station S , while queue Q_H acts in the same manner as in the original opportunistic NOMA system Φ^{ON} .

In dominant system Φ_1^{ON} , the service process of queue Q_H depends on the condition of the channel between base station S and user u_1^H . Base station S transmits the first packet from queue Q_H to user u_1^H using OMA and NOMA when $|h_1^H(t)|^2 \ell(x_1^H) < \theta$ and $|h_1^H(t)|^2 \ell(x_1^H) \geq \theta$, respectively. Note that the average probability of successful packet reception at each high-priority user is the same. Hence, the average service rate of queue Q_H , denoted as μ_H^{ON1} , is given by

$$\begin{aligned}
 \mu_H^{\text{ON1}} &= \mathbb{P}\left(\Gamma_{H1}(t, 1) \geq \Gamma_{\text{th}}^H, |h_1^H(t)|^2 \ell(x_1^H) < \theta\right) \\
 &+ \mathbb{P}\left(\Gamma_{H1|L1}(t, \alpha_H) \geq \Gamma_{\text{th}}^H, |h_1^H(t)|^2 \ell(x_1^H) \geq \theta\right), \quad (5)
 \end{aligned}$$

where the first and second terms of the right-hand side of (5) represent the probabilities of successful packet reception at high-priority user u_1^H when OMA and NOMA are enabled, denoted as $q_{H1}^{\text{OMA}}(\theta)$ and $q_{H1|L1}^{\text{ON}}(\alpha_H, \theta)$, respectively. The following lemma provides the stability condition for queue Q_H in dominant system Φ_1^{ON} .

Lemma 1: In dominant system Φ_1^{ON} , queue Q_H is stable if

$$\begin{aligned}
 \lambda_H < \mu_H^{\text{ON1}} &= \exp\left(-\rho_H(1+r_H^\beta)\right) - \exp\left(-\theta(1+r_H^\beta)\right) \\
 &+ \exp\left(-\max\left\{\frac{\rho_H}{\alpha_H^2 - \Gamma_{\text{th}}^H \alpha_L^2}, \theta\right\}(1+r_H^\beta)\right), \quad (6)
 \end{aligned}$$

where $\rho_H = \Gamma_{\text{th}}^H \sigma^2 / P_S$, $\theta > \rho_H$, and $\alpha_H^2 > \Gamma_{\text{th}}^H \alpha_L^2$.

Proof: Please refer to Appendix A. \blacksquare

The service process of queue Q_L in dominant system Φ_1^{ON} depends on the status of queue Q_H . If queue Q_H is non-empty, then base station S transmits the first packet from queue Q_L to user u_1^L using NOMA when $|h_1^H(t)|^2 \ell(x_1^H) \geq \theta$. If queue Q_H is empty, then base station S transmits the first and second packets from queue Q_L to users u_1^L and u_2^L using NOMA when those two packets are intended for different users (which occurs with probability $1 - \frac{1}{K}$), and transmits the first packet from queue Q_L to user u_1^L using OMA when the first two packets are intended for the same user (which occurs with probability $\frac{1}{K}$). As all low-priority users follow the same location distribution, the average probability of successful

packet reception at each low-priority user is the same. The average service rate of queue Q_L , denoted as μ_L^{ON1} , can be expressed as

$$\begin{aligned}
 \mu_L^{\text{ON1}} &= \mathbb{P}(Q_H(t) > 0) \mathbb{P}\left(|h_1^H(t)|^2 \ell(x_1^H) \geq \theta\right) q_{L1|H1}^{\text{ON}}(\alpha_L) \\
 &+ \mathbb{P}(Q_H(t) = 0) \left(\left(1 - \frac{1}{K}\right) q_{L1L2}^{\text{ON}} + \frac{1}{K} q_{L1}^{\text{OMA}} \right), \quad (7)
 \end{aligned}$$

where $\mathbb{P}(Q_H(t) > 0) = \lambda_H / \mu_H^{\text{ON1}}$, $q_{L1|H1}^{\text{ON}}(\alpha_L)$ is the probability of successful packet reception at user u_1^L with power allocation coefficient α_L when paired with user u_1^H , q_{L1L2}^{ON} is the summation of the probabilities of successful packet reception at users u_1^L and u_2^L using NOMA, and q_{L1}^{OMA} denotes the probability of successful packet reception at user u_1^L using OMA. For two paired low-priority users (i.e., u_1^L and u_2^L), the transmit power allocation coefficients for the users closer to and farther from the base station are denoted as α_n and α_f , respectively, with $\alpha_n^2 + \alpha_f^2 = 1$. The following lemma presents the stability condition for queue Q_L in dominant system Φ_1^{ON} .

Lemma 2: In dominant system Φ_1^{ON} , queue Q_L is stable if

$$\begin{aligned}
 \lambda_L < \mu_L^{\text{ON1}} &= \frac{\lambda_H}{\mu_H^{\text{ON1}}} \exp\left(-\theta(1+r_H^\beta)\right) q_{L1|H1}^{\text{ON}}(\alpha_L) \\
 &+ \left(1 - \frac{\lambda_H}{\mu_H^{\text{ON1}}}\right) \eta, \quad (8)
 \end{aligned}$$

where μ_H^{ON1} is given in (6),

$$q_{L1|H1}^{\text{ON}}(\alpha_L) = \frac{2}{r^2 \beta} N_1^{-2/\beta} \exp(-N_1) \gamma\left(\frac{2}{\beta}, N_1 r^\beta\right), \quad (9)$$

$$\begin{aligned}
 \eta &= \left(1 - \frac{1}{K}\right) \left(q_{L1|Ln}^{\text{ON}}(\alpha_f) + q_{L1|Lf}^{\text{ON}}(\alpha_n)\right) \\
 &+ \frac{1}{K} q_{L1}^{\text{OMA}}, \quad (10)
 \end{aligned}$$

$$q_{L1|Ln}^{\text{ON}}(\alpha_f) = \frac{4}{r^4 \beta} N_2^{-4/\beta} \exp(-N_2) \gamma\left(\frac{4}{\beta}, N_2 r^\beta\right), \quad (11)$$

$$\begin{aligned}
 q_{L1|Lf}^{\text{ON}}(\alpha_n) &= \frac{4}{r^2 \beta} N_3^{-2/\beta} \exp(-N_3) \gamma\left(\frac{2}{\beta}, N_3 r^\beta\right) \\
 &- \frac{4}{r^4 \beta} N_3^{-4/\beta} \gamma\left(\frac{4}{\beta}, N_3 r^\beta\right), \quad (12)
 \end{aligned}$$

$$q_{L1}^{\text{OMA}} = \frac{2}{r^2 \beta} \rho_L^{-2/\beta} \exp(-\rho_L) \gamma\left(\frac{2}{\beta}, \rho_L r^\beta\right), \quad (13)$$

$N_1 = \max\left\{\frac{\rho_H}{\alpha_H^2 - \Gamma_{\text{th}}^H \alpha_L^2}, \frac{\rho_L}{\alpha_f^2}\right\}$, $N_2 = \frac{\rho_L}{\alpha_f^2 - \Gamma_{\text{th}}^L \alpha_n^2}$, $N_3 = \max\left\{\frac{\rho_L}{\alpha_f^2 - \Gamma_{\text{th}}^L \alpha_n^2}, \frac{\rho_L}{\alpha_n^2}\right\}$, $\rho_L = \Gamma_{\text{th}}^L \sigma^2 / P_S$, $\alpha_H^2 > \Gamma_{\text{th}}^H \alpha_L^2$, $\alpha_f^2 > \Gamma_{\text{th}}^L \alpha_n^2$, and $\gamma(w, v) = \int_0^v e^{-z} z^{w-1} dz$ is the lower incomplete Gamma function [32].

Proof: Please refer to Appendix B. \blacksquare

Dominant system Φ_1^{ON} is stable if both queues Q_H and Q_L are stable, i.e., both $\lambda_H < \mu_H^{\text{ON1}}$ and $\lambda_L < \mu_L^{\text{ON1}}$ hold. As a result, based on Lemmas 1 and 2, the stable throughput region

of dominant system Φ_1^{ON} , denoted as $\mathcal{R}_1^{\text{ON}}$, is given by

$$\mathcal{R}_1^{\text{ON}} = \left\{ (\lambda_H, \lambda_L) : \frac{\lambda_L}{\eta} + \frac{(\eta - \exp(-\theta(1+r_H^\beta))) q_{L1|H1}^{\text{ON}}(\alpha_L)}{\eta \mu_H^{\text{ON1}}} \lambda_H < 1, \right. \\ \left. \text{for } 0 \leq \lambda_H < \mu_H^{\text{ON1}} \right\}. \quad (14)$$

B. Stable Throughput Region of Dominant System Φ_2^{ON}

Dominant system Φ_2^{ON} : If queue Q_H is empty, then queue Q_H contributes a dummy packet, while queue Q_L acts in the same manner as in the original opportunistic NOMA system Φ^{ON} .

In dominant system Φ_2^{ON} , base station S transmits the first packet from queue Q_L to user u_1^L using NOMA when $|h_1^H(t)|^2 \ell(x_1^H) \geq \theta$. The average service rate of queue Q_L can be expressed as $\mu_L^{\text{ON2}} = \exp(-\theta(1+r_H^\beta)) q_{L1|H1}^{\text{ON}}(\alpha_L)$. Queue Q_L in dominant system Φ_2^{ON} is stable if $\lambda_L < \mu_L^{\text{ON2}}$. The service process of queue Q_H in dominant system Φ_2^{ON} depends on queue Q_L . If queue Q_L is empty, then base station S transmits the first packet from queue Q_H to user u_1^H using OMA. If queue Q_L is non-empty, then base station S transmits the first packet from queue Q_H to user u_1^H using NOMA and OMA when $|h_1^H(t)|^2 \ell(x_1^H) \geq \theta$ and $|h_1^H(t)|^2 \ell(x_1^H) < \theta$, respectively. The average service rate of queue Q_H in dominant system Φ_2^{ON} is given by

$$\mu_H^{\text{ON2}} = \mathbb{P}(Q_L(t) = 0) \mu_H^{\text{OMA}} + \mathbb{P}(Q_L(t) > 0) (q_{H1}^{\text{OMA}}(\theta) + q_{H1|L1}^{\text{ON}}(\alpha_H, \theta)), \quad (15)$$

where $\mu_H^{\text{OMA}} = \exp(-\rho_H(1+r_H^\beta))$, $\mathbb{P}(Q_L(t) = 0) = 1 - \lambda_L/\mu_L^{\text{ON2}}$, and $q_{H1}^{\text{OMA}}(\theta)$ and $q_{H1|L1}^{\text{ON}}(\alpha_H, \theta)$ are given in (38) and (39), respectively. Queue Q_H in dominant system Φ_2^{ON} is stable if $\lambda_H < \mu_H^{\text{ON2}}$.

The stable throughput region of dominant system Φ_2^{ON} , denoted as $\mathcal{R}_2^{\text{ON}}$, is given by

$$\mathcal{R}_2^{\text{ON}} = \left\{ (\lambda_H, \lambda_L) : \frac{\lambda_H}{\mu_H^{\text{OMA}}} + \frac{(\mu_H^{\text{OMA}} - q_{H1|L1}^{\text{ON}}(\alpha_H, \theta) - q_{H1}^{\text{OMA}}(\theta)) \lambda_L}{\exp(-\theta(1+r_H^\beta)) \mu_H^{\text{OMA}} q_{L1|H1}^{\text{ON}}(\alpha_L)} < 1, \right. \\ \left. \text{for } 0 \leq \lambda_L < \exp(-\theta(1+r_H^\beta)) q_{L1|H1}^{\text{ON}}(\alpha_L) \right\}. \quad (16)$$

The following theorem presents the stable throughput region of the original opportunistic NOMA system Φ^{ON} .

Theorem 1: The stable throughput region of the original dominant NOMA system Φ^{ON} for fixed power allocation coefficients and threshold θ , denoted as \mathcal{R}^{ON} , is equal to the union of the stable throughput regions of dominant systems Φ_1^{ON} and Φ_2^{ON} , i.e., $\mathcal{R}^{\text{ON}} = \mathcal{R}_1^{\text{ON}} \cup \mathcal{R}_2^{\text{ON}}$.

Proof: Please refer to Appendix C. ■

Due to the complexity of analytically deriving the full stable throughput region, we resort to numerical analysis to obtain the full stable throughput region in Section VI, as in [17] and [18].

IV. COOPERATIVE NOMA WITH FULL-DUPLEX RELAYING

The proposed opportunistic NOMA scheme enhances the stable throughput region by providing more transmission opportunities to the low-priority users without improving the performance of the high-priority users. In this section, we propose a cooperative NOMA scheme with full-duplex relaying to improve the reception reliability of the high-priority users with the help of the low-priority users. By exploiting the cooperative diversity gain, the transmission opportunities of the low-priority users can be further increased. The cooperative NOMA system with full-duplex relaying, denoted as Φ^{FCN} , is described as follows.

Case 1: If $Q_H(t) > 0$ and $Q_L(t) > 0$, then base station S transmits the first packet from queue Q_H and the first packet from queue Q_L to users u_1^H and u_1^L , respectively, using cooperative NOMA with fixed power allocation coefficients (α_H^2, α_L^2) . Before transmission begins, base station S informs low-priority user u_1^L to act as a full-duplex relay. In accordance with the NOMA decoding strategy, low-priority user u_1^L decodes signal $s_1^H(t)$ intended for high-priority user u_1^H before performing SIC. By utilizing suitable channel coding (e.g., convolutional coding), low-priority user u_1^L can decode signal $s_1^H(t)$ after a delay of δ symbol durations. Hence, after δ symbol durations, low-priority user u_1^L , which is assumed to be a full-duplex node, simultaneously receives the superimposed signal from the base station and forwards the delayed version of signal $s_1^H(t)$ to high-priority user u_1^H [14], [33]. Full-duplex relaying prototypes have been reported in the literature, e.g., [34]. High-priority user u_1^H constructively combines and decodes the signal transmitted by base station S and its delayed version forwarded by user u_1^L .³ At the end of time slot t , low-priority user u_1^L performs SIC to remove the contribution of signal $s_1^H(t)$ from its received signal, and then decodes its own signal $s_1^L(t)$. The delay δ can be made much smaller than the packet size, and hence, it is neglected for the analysis in this paper. On the other hand, if user u_1^L cannot successfully decode signal $s_1^H(t)$, then user u_1^L remains silent in time slot t and decodes the signals without suffering from the self-interference caused by full-duplex relaying.

Case 2: If $Q_H(t) > 0$ and $Q_L(t) = 0$, then base station S transmits the first packet from queue Q_H to high-priority user u_1^H using cooperative OMA. Among all low-priority users, the low-priority user that can decode signal $s_1^H(t)$ from base station S and has the best channel condition with respect to high-priority user u_1^H is selected as the best relay. The best relay forwards the delayed version of the signal to user u_1^H in the same time slot. Various efficient relay selection schemes have been proposed in the literature. High-priority user u_1^H

³The constructive combination of the signals from the direct and full-duplex forwarding links has recently been implemented in [34] based on a constructive filter and a Viterbi-style decoder.

constructively combines and decodes the signal received from base station S and its delayed version received from the best relay. If no low-priority user can successfully decode signal $s_1^H(t)$, then user u_1^H decodes signal $s_1^H(t)$ only based on the signal transmitted by base station S .

Case 3: If $Q_H(t) = 0$ and $Q_L(t) > 0$, then base station S transmits the first and second packets from queue Q_L to users u_1^L and u_2^L , respectively, using NOMA when the first two packets are intended for different users, and transmits the first packet from queue Q_L to user u_1^L using OMA with power P_S when the first two packets are intended for the same user or $Q_L(t) = 1$.

Queues Q_H and Q_L in the cooperative NOMA system with full-duplex relaying Φ^{FCN} interact with each other, as the average service rate of queue Q_H (Q_L) depends on the status of queue Q_L (Q_H). When queue Q_L is non-empty, base station S transmits the first packet from queue Q_H using cooperative NOMA. When queue Q_L is empty, base station S transmits the first packet from queue Q_H using cooperative OMA. The probabilities of successful packet reception at user u_1^H under these two conditions are different. Thus, their average service rates cannot be directly calculated. To decouple the interacting queues and facilitate the derivation of the stable throughput region, we construct two dominant systems, denoted as Φ_1^{FCN} and Φ_2^{FCN} , by using the concept of stochastic dominance, as discussed in the following.

A. Stable Throughput Region of Dominant System Φ_1^{FCN}

Dominant system Φ_1^{FCN} : If queue Q_L is empty, then queue Q_L contributes a dummy packet, while queue Q_H acts in the same manner as in the cooperative NOMA system with full-duplex relaying Φ^{FCN} . In dominant system Φ_1^{FCN} , a randomly selected low-priority user u_1^L acts as a full-duplex relay in time slot t when the following condition is satisfied:

$$\begin{aligned} \Gamma_{H1 \rightarrow L1}^{\text{FCN}}(t, \alpha_H) &= \frac{\alpha_H^2 P_S |h_1^L(t)|^2 \ell(x_1^L)}{\alpha_L^2 P_S |h_1^L(t)|^2 \ell(x_1^L) + \zeta P_L + \sigma^2} \\ &\geq \Gamma_{\text{th}}^H, \end{aligned} \quad (17)$$

where $\Gamma_{H1 \rightarrow L1}^{\text{FCN}}(t, \alpha_H)$ denotes the SINR of signal $s_1^H(t)$ observed at user u_1^L in time slot t when cooperative NOMA is enabled, ζ denotes the residual self-interference-to-power ratio due to imperfect self-interference cancellation, and P_L is the transmit power of the low-priority users.

The service process of queue Q_H depends on the value of $\Gamma_{H1 \rightarrow L1}^{\text{FCN}}(t, \alpha_H)$. Base station S transmits the first packet from queues Q_H to user u_1^H using NOMA and cooperative NOMA when $\Gamma_{H1 \rightarrow L1}^{\text{FCN}}(t, \alpha_H) < \Gamma_{\text{th}}^H$ and $\Gamma_{H1 \rightarrow L1}^{\text{FCN}}(t, \alpha_H) \geq \Gamma_{\text{th}}^H$, respectively. Hence, the average service rate of queue Q_H in

dominant system Φ_1^{FCN} , denoted as μ_H^{FCN1} , is given by

$$\begin{aligned} \mu_H^{\text{FCN1}} &= \mathbb{P}(\Gamma_{H1|L1}^H(t, \alpha_H) \geq \Gamma_{\text{th}}^H, \Gamma_{H1 \rightarrow L1}^{\text{FCN}}(t, \alpha_H) < \Gamma_{\text{th}}^H) \\ &\quad + \mathbb{P}(\Gamma_{H1|L1}^{\text{FCN}}(t, \alpha_H) \geq \Gamma_{\text{th}}^H, \Gamma_{H1 \rightarrow L1}^{\text{FCN}}(t, \alpha_H) \geq \Gamma_{\text{th}}^H), \end{aligned} \quad (18)$$

where $\Gamma_{H1|L1}^H(t, \alpha_H)$ is given in (2). The SINR of signal $s_1^H(t)$ observed at user u_1^H in time slot t when cooperative NOMA is enabled, denoted as $\Gamma_{H1|L1}^{\text{FCN}}(t, \alpha_H)$, can be expressed as

$$\begin{aligned} \Gamma_{H1|L1}^{\text{FCN}}(t, \alpha_H) &= \frac{\alpha_H^2 P_S |h_1^H(t)|^2 \ell(x_1^H) + P_L |g_{1,1}^{\text{HL}}(t)|^2 \ell(x_1^H - x_1^L)}{\alpha_L^2 P_S |h_1^H(t)|^2 \ell(x_1^H) + \sigma^2}, \\ &\text{if } \Gamma_{H1 \rightarrow L1}^{\text{FCN}}(t, \alpha_H) \geq \Gamma_{\text{th}}^H, \end{aligned} \quad (19)$$

where $g_{1,1}^{\text{HL}}(t)$ and $\ell(x_1^H - x_1^L)$ denote the Rayleigh fading channel gain and non-singular path loss between users u_1^H and u_1^L in time slot t , respectively. The following lemma provides the stability condition for queue Q_H in dominant system Φ_1^{FCN} .

Lemma 3: In dominant system Φ_1^{FCN} , queue Q_H is stable if

$$\begin{aligned} \lambda_H < \mu_H^{\text{FCN1}} &= \exp\left(-\frac{\rho_H(1+r_H^\beta)}{\alpha_H^2 - \Gamma_{\text{th}}^H \alpha_L^2}\right) \\ &\quad \times \left(1 - \frac{2}{r^2 \beta} N_4^{-2/\beta} \exp(-N_4) \gamma\left(\frac{2}{\beta}, N_4 r^\beta\right)\right) \\ &\quad + C(\alpha_H) + \frac{2N_5}{r^2 \beta} N_4^{-2/\beta} \exp(-N_4) \gamma\left(\frac{2}{\beta}, N_4 r^\beta\right), \end{aligned} \quad (20)$$

where $N_4 = \frac{(\zeta P_L + \sigma^2) \Gamma_{\text{th}}^H}{(\alpha_H^2 - \Gamma_{\text{th}}^H \alpha_L^2) P_S}$, $\ell(x_1^H - x_1^L) = \left(1 + (r_H^2 + (r_1^L)^2 - 2 r_H r_1^L \cos \tau_1^L)^{\beta/2}\right)^{-1}$, $N_5 = \exp\left(-\frac{\Gamma_{\text{th}}^H \sigma^2}{Z}\right)$, $Z = (\alpha_H^2 - \Gamma_{\text{th}}^H \alpha_L^2) P_S \ell(x_1^H)$, $\ell(x_1^L) = \left(1 + (r_1^L)^\beta\right)^{-1}$, $\alpha_H^2 > \Gamma_{\text{th}}^H \alpha_L^2$, and $C(\alpha_H)$ is given in (21), as shown at the bottom of this page.

Proof: Please refer to Appendix D. ■

The service process of queue Q_L can also be divided into two cases: a) if queue Q_H is non-empty, then base station S transmits the first packet from queue Q_L to user u_1^L using cooperative NOMA; b) if queue Q_H is empty, then base station S transmits the first two packets from queue Q_L to users u_1^L and u_2^L using NOMA when $u_1^L \neq u_2^L$ and the first packet from queue Q_L to user u_1^L using OMA when $u_1^L = u_2^L$. The average service rate of queue Q_L , denoted as μ_L^{FCN1} , is

$$\begin{aligned} \mu_L^{\text{FCN1}} &= \mathbb{P}(Q_H(t) > 0) q_{L1|H1}^{\text{FCN}}(\alpha_L) + \mathbb{P}(Q_H(t) = 0) \\ &\quad \times \left(\left(1 - \frac{1}{K}\right) q_{L1L2}^{\text{ON}} + \frac{1}{K} q_{L1}^{\text{OMA}} \right), \end{aligned} \quad (22)$$

$$C(\alpha_H) = \frac{1}{\pi r^2} \int_0^r \int_0^{2\pi} \frac{\exp(-N_4/\ell(x_1^L))}{1 - \frac{Z}{P_L \ell(x_1^H - x_1^L)}} \left(\exp\left(-\frac{\Gamma_{\text{th}}^H \sigma^2}{P_L \ell(x_1^H - x_1^L)}\right) - N_5 \right) r_1^L dr_1^L d\tau_1^L \quad (21)$$

$$q_{L1|H1}^{\text{FCN}}(\alpha_L) = \mathbb{P}(\Gamma_{H1 \rightarrow L1}^{\text{FCN}}(t, \alpha_H) \geq \Gamma_{\text{th}}^H, \Gamma_{L1}^{\text{FCN}}(t, \alpha_L) \geq \Gamma_{\text{th}}^L) \\ + \mathbb{P}(\Gamma_{H1 \rightarrow L1}^{\text{FCN}}(t, \alpha_H) < \Gamma_{\text{th}}^H, \Gamma_{H1 \rightarrow L1}(t, \alpha_H) \geq \Gamma_{\text{th}}^H, \Gamma_{L1}^{\text{FCN}}(t, \alpha_L) \geq \Gamma_{\text{th}}^L) \quad (24)$$

where $\mathbb{P}(Q_H(t) > 0) = \lambda_H / \mu_H^{\text{FCN1}}$, $q_{L1|H1}^{\text{FCN}}(\alpha_L)$ is the probability of successful packet reception at user u_1^L when cooperative NOMA is enabled, and q_{L1L2}^{ON} and q_{L1}^{OMA} are given in (47) and (48), respectively.

Depending on whether or not user u_1^L forwards signal $s_1^H(t)$ to user u_1^H , the received SINR of signal $s_1^L(t)$ observed at user u_1^L in time slot t can be expressed as

$$\Gamma_{L1}^{\text{FCN}}(t, \alpha_L) = \begin{cases} \frac{\alpha_L^2 P_S |h_1^L(t)|^2 \ell(x_1^L)}{\zeta P_L + \sigma^2}, & \text{if } \Gamma_{H1 \rightarrow L1}^{\text{FCN}}(t, \alpha_H) \geq \Gamma_{\text{th}}^H, \\ \frac{\alpha_L^2 P_S |h_1^L(t)|^2 \ell(x_1^L)}{\sigma^2}, & \text{if } \Gamma_{H1 \rightarrow L1}^{\text{FCN}}(t, \alpha_H) < \Gamma_{\text{th}}^H. \end{cases} \quad (23)$$

As a result, we obtain (24), as shown at the top of this page.

The following lemma provides the stability condition for queue Q_L in dominant system Φ_1^{FCN} .

Lemma 4: In dominant system Φ_1^{FCN} , queue Q_L is stable if

$$\lambda_L < \mu_L^{\text{FCN1}} = \frac{\lambda_H}{\mu_H^{\text{FCN1}}} \left(\frac{2}{r^2 \beta} N_6^{-2/\beta} \exp(-N_6) \gamma\left(\frac{2}{\beta}, N_6 r^\beta\right) \right. \\ \left. + \frac{2}{r^2 \beta} N_1^{-2/\beta} \exp(-N_1) \gamma\left(\frac{2}{\beta}, N_1 r^\beta\right) \right. \\ \left. - \frac{2}{r^2 \beta} N_4^{-2/\beta} \exp(-N_4) \gamma\left(\frac{2}{\beta}, N_4 r^\beta\right) \right) \\ + \left(1 - \frac{\lambda_H}{\mu_H^{\text{FCN1}}}\right) \eta, \quad (25)$$

where $N_6 = \max\left\{\frac{(\zeta P_L + \sigma^2) \Gamma_{\text{th}}^H}{(\alpha_H^2 - \Gamma_{\text{th}}^H \alpha_L^2) P_S}, \frac{(\zeta P_L + \sigma^2) \Gamma_{\text{th}}^L}{\alpha_L^2 P_S}\right\}$, $\alpha_H^2 > \Gamma_{\text{th}}^H \alpha_L^2$, and η is given in (10).

Proof: Please refer to Appendix E. ■

Based on the average service rates of queues Q_H and Q_L , the stable throughput region of dominant system Φ_1^{FCN} , denoted as $\mathcal{R}_1^{\text{FCN}}$, can be expressed as

$$\mathcal{R}_1^{\text{FCN}} = \left\{ (\lambda_H, \lambda_L) : \frac{(\eta - q_{L1|H1}^{\text{FCN}}(\alpha_L)) \lambda_H}{\eta (q_{H1}^{\text{N}}(\alpha_H) + q_{H1}^{\text{FCN}}(\alpha_H))} + \frac{\lambda_L}{\eta} < 1, \right. \\ \left. \text{for } 0 \leq \lambda_H < q_{H1}^{\text{N}}(\alpha_L) + q_{H1}^{\text{FCN}}(\alpha_H) \right\}. \quad (26)$$

According to (26), stable throughput region $\mathcal{R}_1^{\text{FCN}}$ depends on the self-interference cancellation coefficient.

B. Stable Throughput Region of Dominant System Φ_2^{FCN}

Dominant system Φ_2^{FCN} : If queue Q_H is empty, then queue Q_H contributes a dummy packet, while queue Q_L acts in the same manner as in the cooperative NOMA system with full-duplex relaying Φ^{FCN} . In dominant system Φ_2^{FCN} , base station S transmits the first packet from queue Q_L to

user u_1^L using cooperative NOMA. The average service rate of queue Q_L , denoted as μ_L^{FCN2} , can be expressed as $\mu_L^{\text{FCN2}} = q_{L1|H1}^{\text{FCN}}(\alpha_L)$. Queue Q_L in dominant system Φ_2^{FCN} is stable if $\lambda_L < \mu_L^{\text{FCN2}}$. The service process of queue Q_H depends on queue Q_L . If queue Q_L is empty, then base station S transmits the first packet from queue Q_H to user u_1^H using cooperative OMA. If queue Q_L is non-empty, then base station S transmits the first packet from queue Q_H to user u_1^H using cooperative NOMA when $\Gamma_{H1 \rightarrow L1}^{\text{FCN}}(t, \alpha_H) \geq \Gamma_{\text{th}}^H$ and using NOMA when $\Gamma_{H1 \rightarrow L1}^{\text{FCN}}(t, \alpha_H) < \Gamma_{\text{th}}^H$. Thus, the average service rate of queue Q_H in dominant system Φ_2^{FCN} , denoted as μ_H^{FCN2} , is given by

$$\mu_H^{\text{FCN2}} = \mathbb{P}(Q_L(t) = 0) q_{H1}^{\text{FC}} \\ + \mathbb{P}(Q_L(t) > 0) (q_{H1}^{\text{N}}(\alpha_H) + q_{H1}^{\text{FCN}}(\alpha_H)), \quad (27)$$

where $\mathbb{P}(Q_L(t) = 0) = 1 - \lambda_L / \mu_L^{\text{FCN2}}$, q_{H1}^{FC} denotes the probability of successful packet reception at user u_1^H when cooperative OMA is enabled, and $q_{H1}^{\text{N}}(\alpha_H)$ and $q_{H1}^{\text{FCN}}(\alpha_H)$ are given in (49) and (52), respectively. When cooperative OMA is enabled, the low-priority users that can successfully decode signal $s_1^H(t)$ are referred to as *qualified relays*, which form the decoding set in time slot t , denoted as $\Omega(t)$ and given by

$$\Omega(t) = \{u_k^{\text{R}} \in \mathcal{U}^{\text{L}} : \Gamma_{H1 \rightarrow \text{R}k}^{\text{FC}}(t) \geq \Gamma_{\text{th}}^H\}, \quad (28)$$

where $u_k^{\text{R}} \in \mathcal{U}^{\text{L}}$ denotes the k -th full-duplex relay and $\Gamma_{H1 \rightarrow \text{R}k}^{\text{FC}}(t) = \frac{P_S |h_k^{\text{R}}(t)|^2 \ell(x_k^{\text{R}})}{\zeta P_L + \sigma^2}$.

We assume that, via coordination signaling between base station S and user u_1^H before the packet transmission, each qualified relay knows the instantaneous channel gain between itself and user u_1^H . If decoding set $\Omega(t)$ is empty, no low-priority user can help forward signal $s_1^H(t)$ to user u_1^H . On the other hand, if decoding set $\Omega(t)$ is non-empty, the qualified relay that has the best channel condition with respect to high-priority user u_1^H is selected as the best relay, i.e.,

$$u_b^{\text{R}} = \arg \max_{u_k^{\text{R}} \in \Omega(t)} \left\{ P_L |g_{1,k}^{\text{HR}}(t)|^2 \ell(x_1^{\text{H}} - x_k^{\text{R}}) \right\}. \quad (29)$$

User u_1^H can successfully decode signal $s_1^H(t)$ in time slot t if the received SNR is not less than the reception threshold, i.e.,

$$\Gamma_{H1 \text{R}b}^{\text{FC}}(t) = \frac{P_S |h_1^{\text{H}}(t)|^2 \ell(x_1^{\text{H}}) + P_L |g_{1,b}^{\text{HR}}(t)|^2 \ell(x_b^{\text{R}} - x_1^{\text{H}})}{\sigma^2} \\ \geq \Gamma_{\text{th}}^H, \quad (30)$$

where $\Gamma_{H1 \text{R}b}^{\text{FC}}(t)$ denotes the SNR of signal $s_1^H(t)$ observed at user u_1^H in time slot t when user u_b^{R} acts as the full-duplex relay.

The probability of successful packet transmission is the complement of the outage probability. In this context, an outage occurs when high-priority user u_1^H fails to decode

699 the packet after constructively combining the signals trans-
 700 mitted by the base station and the best relay u_b^R . By selecting
 701 the best relay, this outage event is equivalent to the event that
 702 all qualified relays are in outage, which means that no low-
 703 priority user satisfies the following condition:

$$704 \Gamma_{H1 \rightarrow Rk}^{FC} \geq \Gamma_{th}^H \quad \text{and} \quad \Gamma_{H1Rk}^{FC} \geq \Gamma_{th}^H, \quad \forall u_k^R \in \mathcal{U}^L. \quad (31)$$

705 The following lemma presents the stability condition for
 706 queue Q_H in dominant system Φ_2^{FCN} .

707 *Lemma 5:* In dominant system Φ_2^{FCN} , queue Q_H is stable if

$$708 \lambda_H < \mu_H^{FCN2} = \left(1 - \frac{\lambda_L}{\mu_L^{FCN2}}\right) q_{H1}^{FC} \\
 709 + \frac{\lambda_L}{\mu_L^{FCN2}} (q_{H1}^N(\alpha_H) + q_{H1}^{FCN}(\alpha_H)), \quad (32)$$

710 where $\mu_L^{FCN2} = q_{L1|H1}^{FCN}(\alpha_L)$, $q_{H1}^N(\alpha_H)$ and $q_{H1}^{FCN}(\alpha_H)$ are
 711 given in (49) and (52), respectively, and

$$712 q_{H1}^{FC} = \exp\left(-\rho_H(1+r_H^\beta)\right) + \sum_{j=1}^K \binom{K}{j} (-1)^{j+1} (C(1))^j \\
 713 \times \left(1 - \exp\left(-\rho_H(1+r_H^\beta)\right)\right)^{1-j}. \quad (33)$$

714 *Proof:* Please refer to Appendix F. \blacksquare

715 After deriving the average service rates of queues
 716 Q_H and Q_L , the stable throughput region of dominant
 717 system Φ_2^{FCN} , denoted as \mathcal{R}_2^{FCN} , can be expressed as

$$718 \mathcal{R}_2^{FCN} = \left\{ (\lambda_H, \lambda_L) : \frac{\lambda_H}{q_{H1}^{FC}} \right. \\
 719 + \frac{(q_{H1}^{FC} - q_{H1}^N(\alpha_H) - q_{H1}^{FCN}(\alpha_H)) \lambda_L}{q_{H1}^{FC} q_{L1|H1}^{FCN}(\alpha_L)} < 1, \\
 720 \left. \text{for } 0 \leq \lambda_L < q_{L1|H1}^{FCN}(\alpha_L) \right\}. \quad (34)$$

721 According to (34), stable throughput region \mathcal{R}_2^{FCN} depends
 722 on the number of low-priority users and the self-interference
 723 cancelation coefficient.

724 Based on the above derivations, the following theorem
 725 presents the stable throughput region of the cooperative
 726 NOMA system with full-duplex relaying Φ^{FCN} .

727 *Theorem 2:* The stable throughput region of the cooperative
 728 NOMA system with full-duplex relaying Φ^{FCN} for fixed
 729 power allocation coefficients, denoted as \mathcal{R}^{FCN} , is the union
 730 of the stable throughput regions of dominant systems Φ_1^{FCN}
 731 and Φ_2^{FCN} , i.e., $\mathcal{R}^{FCN} = \mathcal{R}_1^{FCN} \cup \mathcal{R}_2^{FCN}$.

732 *Proof:* The proof is similar to that of Theorem 1, and
 733 hence, it is omitted here. \blacksquare

734 Similarly, we resort to numerical analysis to obtain the full
 735 stable throughput region in Section VI.

736 V. COMPARISON OF NOMA AND OMA

737 In this section, we derive the stable throughput region of a
 738 baseline OMA scheme and the conditions under which the
 739 proposed NOMA schemes achieve larger stable throughput
 740 regions than the baseline OMA scheme.

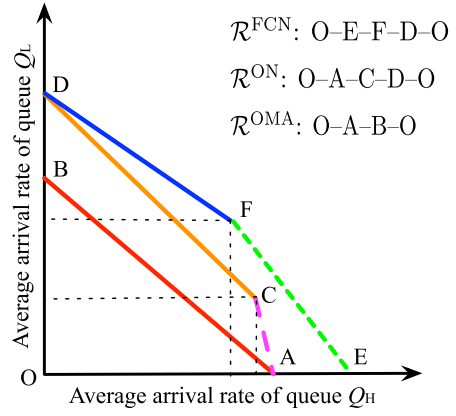


Fig. 2. Stable throughput regions of OMA system Φ^{OMA} , opportunistic NOMA system Φ^{ON} , and cooperative NOMA system with full-duplex relaying Φ^{FCN} .

741 A. Baseline Orthogonal Multiple Access Scheme

742 We consider a time division multiple access (TDMA) based
 743 OMA system, denoted as Φ^{OMA} , as a baseline, where base
 744 station S transmits one packet in one time slot. As queues
 745 Q_H and Q_L do not interact with each other when OMA is
 746 utilized, the stability conditions of these two queues can be
 747 separately analyzed. Base station S transmits the first packet
 748 from queue Q_H to high-priority user u_1^H whenever queue Q_H
 749 is not empty, regardless of the status of queue Q_L . The average
 750 service rate of queue Q_H in OMA system Φ^{OMA} is $\mu_H^{OMA} =$
 751 $\exp\left(-\rho_H(1+r_H^\beta)\right)$. When queue Q_H is empty, base station
 752 S transmits the first packet from queue Q_L to user u_1^L . The
 753 average service rate of queue Q_L is given by $\mu_L^{OMA} =$
 754 $\mathbb{P}(Q_H = 0) \mathbb{P}(\Gamma_{L1}(t, 1) \geq \Gamma_{th}) = (1 - \lambda_H/\mu_H^{OMA}) q_{L1}^{OMA}$,
 755 where q_{L1}^{OMA} is given in (48). The stable throughput region
 756 of OMA system Φ^{OMA} is given by

$$757 \mathcal{R}^{OMA} = \left\{ (\lambda_H, \lambda_L) : \frac{\lambda_H}{\exp\left(-\rho_H(1+r_H^\beta)\right)} + \frac{\lambda_L}{q_{L1}^{OMA}} < 1, \right. \\
 758 \left. \text{for } 0 \leq \lambda_H < \exp\left(-\rho_H(1+r_H^\beta)\right) \right\}. \quad (35)$$

759 For the queueing model under consideration, the perfor-
 760 mance comparison between the OMA scheme and the
 761 proposed NOMA scheme is fair in the sense that each user
 762 is served based on its priority and the order of its packets in
 763 the queue, but not based on its CSI.

764 B. Comparison Between NOMA and OMA

765 In the following, we present the conditions under which the
 766 proposed NOMA schemes achieve larger stable throughput
 767 regions than OMA. Based on the stable throughput regions
 768 given in (14), (16), (26), (34), and (35), Fig. 2 plots
 769 the stable throughput regions of OMA system Φ^{OMA}
 770 (i.e., \mathcal{R}^{OMA} : O-A-B-O), opportunistic NOMA system Φ^{ON}
 771 (i.e., \mathcal{R}^{ON} : O-A-C-D-O), and the cooperative NOMA system
 772 with full-duplex relaying Φ^{FCN} (i.e., \mathcal{R}^{FCN} : O-E-F-D-O).
 773 The coordinates of the corner points in Fig. 2 are

774 $O = (0, 0)$, $A = (\mu_H^{\text{OMA}}, 0)$, $B = (0, q_{L1}^{\text{OMA}})$, $C =$
 775 $(\mu_H^{\text{ON1}}, \xi q_{L1|H1}^{\text{ON}}(\alpha_L))$, $D = (0, \eta)$, $E = (q_{H1}^{\text{FC}}, 0)$, and $F =$
 776 $(q_{H1}^{\text{N}}(\alpha_H) + q_{H1}^{\text{FCN}}(\alpha_H), q_{L1|H1}^{\text{FCN}}(\alpha_L))$, where $\mu_H^{\text{OMA}} =$
 777 $\exp(-\rho_H(1+r_H^\beta))$ and $\xi = \exp(-\theta(1+r_H^\beta))$.

778 *Proposition 1:* The cooperative NOMA scheme with full-
 779 duplex relaying achieves a larger stable throughput region than
 780 OMA, i.e., $\mathcal{R}^{\text{OMA}} \subset \mathcal{R}^{\text{FCN}}$, when the following conditions
 781 hold:

$$782 \quad q_{L1|H1}^{\text{FCN}}(\alpha_L) > q_{L1}^{\text{OMA}} \left(1 - \frac{q_{H1}^{\text{N}}(\alpha_H) + q_{H1}^{\text{FCN}}(\alpha_H)}{\mu_H^{\text{OMA}}}\right), \quad (36)$$

$$783 \quad q_{L1L2}^{\text{ON}} > q_{L1}^{\text{OMA}}. \quad (37)$$

784 *Proof:* Please refer to Appendix G. ■

785 As $q_{L1|H1}^{\text{FCN}}(\alpha_L)$, $q_{H1}^{\text{N}}(\alpha_H)$, and $q_{H1}^{\text{FCN}}(\alpha_H)$ are functions of
 786 α_L , it is very difficult to derive a closed-form condition
 787 in terms of α_L from (36). Nonetheless, we can obtain all
 788 possible values of α_L that lead to $\mathcal{R}^{\text{OMA}} \subset \mathcal{R}^{\text{FCN}}$ by
 789 evaluating (36) numerically, as (37) does not depend on α_L .
 790 Moreover, the stable throughput region of the cooperative
 791 NOMA scheme with full-duplex relaying can be maximized
 792 by fixing λ_H and then maximizing the corresponding average
 793 service rate of queue Q_L , i.e., μ_L^{FCN1} or μ_L^{FCN2} , by optimizing
 794 the value of α_L .

795 *Proposition 2:* The opportunistic NOMA scheme achieves
 796 a larger stable throughput region than OMA, i.e., $\mathcal{R}^{\text{OMA}} \subset$
 797 \mathcal{R}^{ON} , when $q_{L1L2}^{\text{ON}} > q_{L1}^{\text{OMA}}$ and $\xi q_{L1|H1}^{\text{ON}}(\alpha_L) >$
 798 $q_{L1}^{\text{OMA}} \left(1 - \frac{\mu_H^{\text{ON1}}}{\mu_H^{\text{OMA}}}\right)$ hold.

799 *Proof:* The proof is similar to that of Proposition 1, and
 800 hence, it is omitted here. ■

801 VI. NUMERICAL RESULTS

802 In this section, we evaluate the stable throughput regions
 803 of opportunistic NOMA and cooperative NOMA with full-
 804 duplex relaying and compare them with the stable throughput
 805 region of baseline OMA. The radius of the circular network
 806 coverage area is $r = 1.3$ km, where $M = 4$ high-priority
 807 users are located $r_H = 1.2$ km away from base station S .
 808 The transmit powers (i.e., P_S and P_L) and noise power σ^2
 809 are set to be 1 W and -100 dBm, respectively. We consider
 810 Rayleigh fading channels and the path loss exponent β is set
 811 to be 4. The power allocation coefficients of the far and near
 812 users when NOMA is enabled to serve the first two packets
 813 from queue Q_L , (α_f^2, α_n^2) , are set to be $(0.8, 0.2)$.

814 Fig. 3 shows the impact of the number of low-priority
 815 users K and self-interference cancellation coefficient ζ on the
 816 probabilities of successful packet reception at the high-priority
 817 users when cooperative NOMA and OMA are employed
 818 (i.e., $q_{H1}^{\text{N}}(\alpha_H) + q_{H1}^{\text{FCN}}(\alpha_H)$ and q_{H1}^{FC}). The simulation (Sim)
 819 results match the analytical (Ana) results well, which validates
 820 the performance analysis. We observe that q_{H1}^{FC} increases
 821 with K , as the probability of selecting a full-duplex relay
 822 with good channel condition with respect to user u_1^H becomes
 823 higher because of the spatial diversity gain. On the other hand,
 824 $q_{H1}^{\text{N}}(\alpha_H) + q_{H1}^{\text{FCN}}(\alpha_H)$ does not change with K , as the intended
 825 receiver of the first packet from queue Q_L is selected to act
 826 as a full-duplex relay when it can successfully decode signal

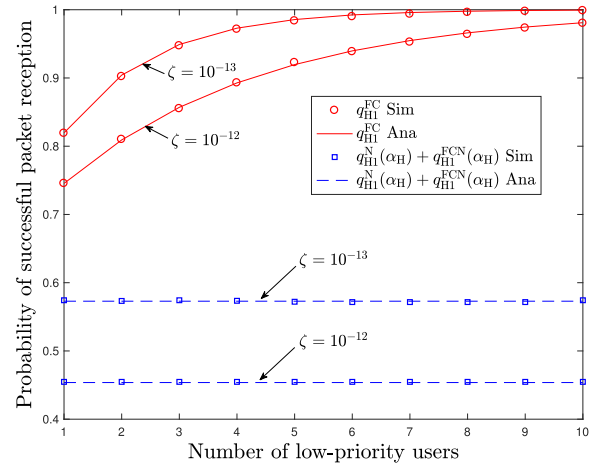


Fig. 3. Probabilities of successful packet reception at user u_1^H versus the number of low-priority users K for different self-interference cancellation coefficients, ζ , when $(\alpha_H^2, \alpha_L^2) = (0.8, 0.2)$ and $\Gamma_{\text{th}}^H = 2$.

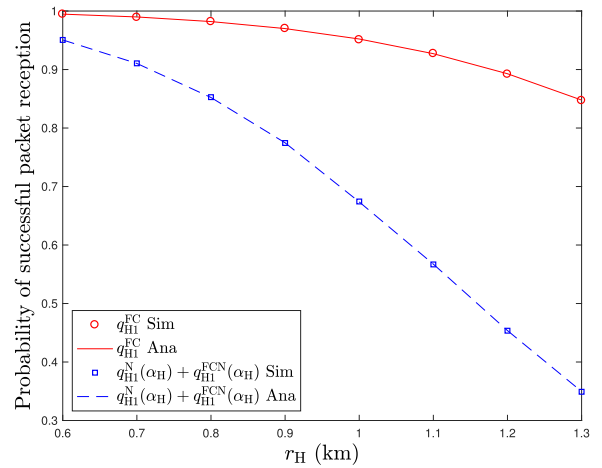


Fig. 4. Probabilities of successful packet reception at user u_1^H versus its distance with respect to base station, r_H , when $(\alpha_H^2, \alpha_L^2) = (0.8, 0.2)$, $\Gamma_{\text{th}}^H = 2$, $\zeta = 10^{-12}$, and $K = 4$.

827 $s_1^H(t)$ received from base station S , regardless of its channel
 828 condition with respect to user u_1^H . With better self-interference
 829 cancellation (i.e., a smaller value of ζ), the probability of
 830 successful packet reception at user u_1^H increases for both
 831 cooperative NOMA and OMA, as the SINR of signal $s_1^H(t)$
 832 at the low-priority users becomes larger and in turn the
 833 probability of selecting a reliable full-duplex relay increases.

834 Fig. 4 illustrates the impact of the distance between the base
 835 station and the high-priority users, r_H , on the probabilities of
 836 successful packet reception at the high-priority users (i.e., q_{H1}^{FC}
 837 and $q_{H1}^{\text{N}}(\alpha_H) + q_{H1}^{\text{FCN}}(\alpha_H)$). As r_H increases, both probabilities
 838 decrease because of the larger path loss. As the relay with
 839 the best channel condition with respect to user u_1^H is selected
 840 for cooperative NOMA, the gap between q_{H1}^{FC} and $q_{H1}^{\text{N}}(\alpha_H) +$
 841 $q_{H1}^{\text{FCN}}(\alpha_H)$ becomes larger as r_H increases. In addition, q_{H1}^{FC} is
 842 always larger than $q_{H1}^{\text{N}}(\alpha_H) + q_{H1}^{\text{FCN}}(\alpha_H)$ because of the higher
 843 base station transmit power for the high-priority user as well
 844 as the higher spatial diversity gain due to relay selection.

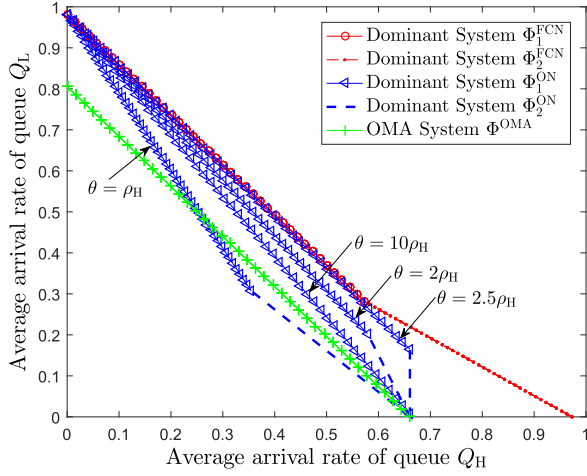


Fig. 5. Stable throughput region for different values of threshold θ when $(\alpha_H^2, \alpha_L^2) = (0.8, 0.2)$, $\Gamma_{th}^H = 2$, $\Gamma_{th}^L = 2.5$, $K = 4$, and $\zeta = 10^{-13}$.

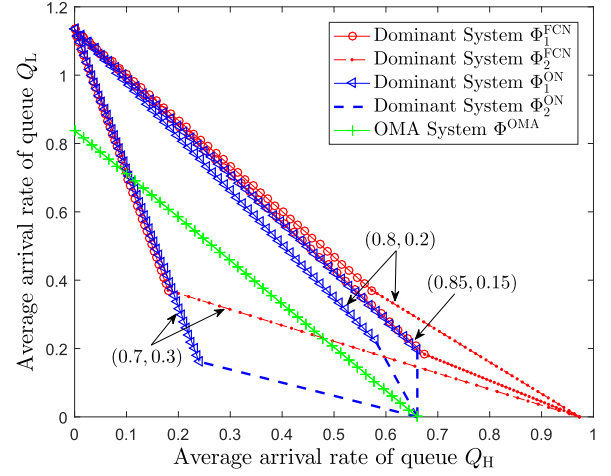


Fig. 6. Stable throughput region for different values of power allocation coefficients (α_H^2, α_L^2) when $\theta = 2\rho_H$, $\Gamma_{th}^H = \Gamma_{th}^L = 2$, $K = 4$, and $\zeta = 10^{-13}$.

845 Fig. 5 plots the stable throughput region for various values
 846 of threshold θ . For opportunistic NOMA system Φ^{ON} and
 847 OMA system Φ^{OMA} , the maximum achievable λ_H is the same,
 848 as the stable throughput region of opportunistic NOMA system
 849 Φ^{ON} is equal to the union of the stable throughput regions
 850 of dominant systems Φ_1^{ON} and Φ_2^{ON} . The stable throughput
 851 region of opportunistic NOMA system Φ^{ON} depends on θ .
 852 As θ increases, the probability of enabling NOMA transmis-
 853 sion decreases, which reduces the transmission opportunities
 854 of the low-priority users. For larger θ , the packet retrans-
 855 mission probability of high-priority users decreases, which in turn
 856 provides more transmission opportunities to the low-priority
 857 users. With an appropriate choice of θ to balance these two
 858 aspects, e.g., $\theta = 2.5\rho_H$ in Fig. 5, opportunistic NOMA
 859 can achieve a much larger stable throughput region than OMA.
 860 By enabling the low-priority users to act as full-duplex relays
 861 and assist the high-priority users, the maximum achievable λ_H
 862 of the cooperative NOMA system with full-duplex relaying
 863 Φ^{FCN} is even larger than that of opportunistic NOMA system
 864 Φ^{ON} due to the cooperative diversity gain. The enhanced
 865 packet reception reliability for the high-priority users can be
 866 exploited to provide more transmission opportunities to the
 867 low-priority users. Thereby, the stable throughput region is
 868 further enlarged.

869 Fig. 6 shows the stable throughput region for various values
 870 of transmit power allocation coefficients (α_H^2, α_L^2) . When
 871 $(\alpha_H^2, \alpha_L^2) = (0.7, 0.3)$, condition (36) does not hold, and
 872 hence, the stable throughput region of the cooperative NOMA
 873 system with full-duplex relaying Φ^{FCN} is not larger than that
 874 of OMA system Φ^{OMA} . This is because the value of α_H^2 is not
 875 large enough for the low-priority users to successfully decode
 876 the signals intended for the high-priority users, which is the
 877 prerequisite for performing SIC. By increasing α_H^2 to 0.8,
 878 the conditions given in Propositions 1 and 2 hold, and hence
 879 the stable throughput regions of both opportunistic NOMA and
 880 cooperative NOMA with full-duplex relaying become larger
 881 than that of OMA. However, by further increasing α_H^2 from
 882 0.8 to 0.85, the stable throughput region of the cooperative
 883 NOMA system with full-duplex relaying Φ^{FCN} decreases.

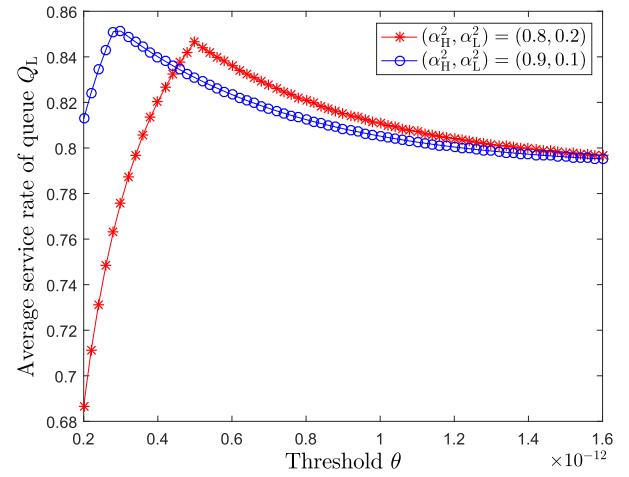


Fig. 7. Average service rate of queue Q_L in opportunistic NOMA system versus threshold θ and power allocation coefficients (α_H^2, α_L^2) when $\lambda_H = 0.2$, $\Gamma_{th}^H = \Gamma_{th}^L = 2$, and $K = 4$.

884 This is because the increased transmission opportunities of
 885 the low-priority users cannot compensate for the reduction of
 886 successful packet reception at the low-priority users due to the
 887 lower transmit power.

888 Fig. 7 shows the impact of threshold θ and power alloca-
 889 tion coefficients on the average service rate of queue Q_L .
 890 If $(\alpha_H^2, \alpha_L^2) = (0.8, 0.2)$, the average service rate of queue
 891 Q_L increases with θ when $\theta < 0.5 \times 10^{-12}$. By enabling
 892 NOMA when the channel gain between base station S and
 893 the high-priority users is larger, fewer packet retransmissions
 894 are required for high-priority users, which in turn improves
 895 the transmission opportunities of low-priority users. The
 896 average service rate of queue Q_L decreases with θ when
 897 $\theta > 0.5 \times 10^{-12}$ and converges to 0.795, as the probability
 898 that NOMA is enabled decreases. By increasing α_H^2 to 0.9,
 899 the optimal threshold θ that maximizes the average service
 900 rate of queue Q_L becomes smaller, as allocating more transmit
 901 power to the high-priority users allows NOMA to be used for
 902 smaller channel gain.

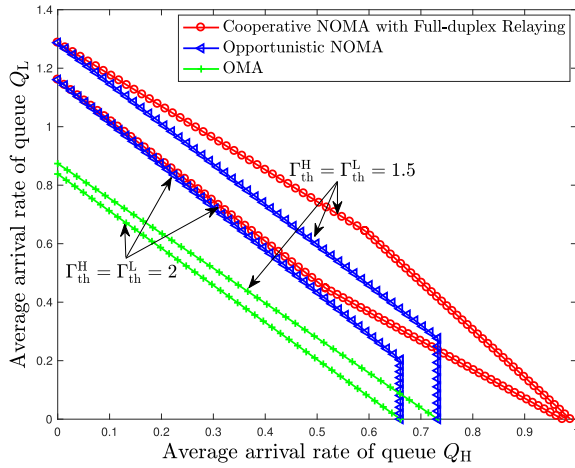


Fig. 8. Full stable throughput region for different values of reception thresholds Γ_{th}^H and Γ_{th}^L when $K = 4$ and $\zeta = 10^{-13}$.

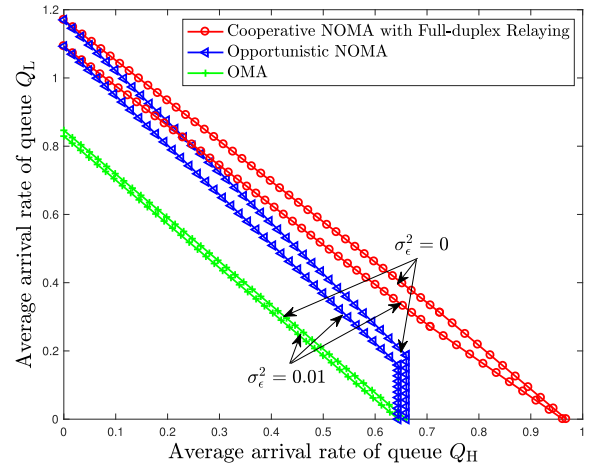


Fig. 10. Impact of imperfect CSI on full stable throughput region when $\Gamma_{th}^H = \Gamma_{th}^L = 2$, $K = 4$, and $\zeta = 10^{-13}$.

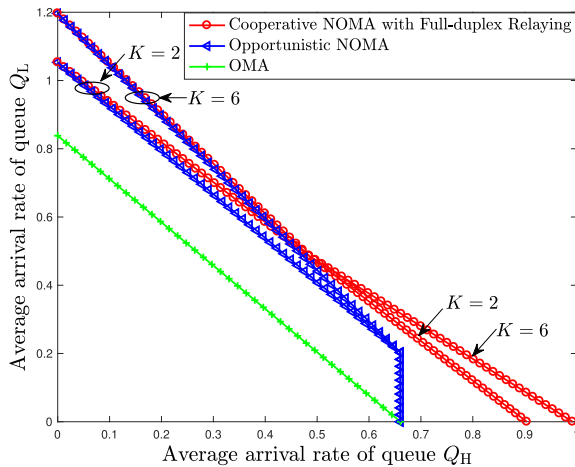


Fig. 9. Full stable throughput region for different numbers of low-priority users K when $\Gamma_{th}^H = \Gamma_{th}^L = 2$ and $\zeta = 10^{-13}$.

In Fig. 8, we study the impact of reception thresholds Γ_{th}^H and Γ_{th}^L on the full stable throughput regions of both proposed NOMA schemes. For a smaller reception threshold, the maximum achievable λ_H and λ_L in all schemes increase, as the probability of successful packet reception at each user increases. With a smaller reception threshold, the probability of queue Q_H being empty is higher, which leads to more time slots being available for the base station to serve queue Q_L using NOMA. Hence, the performance gap between the opportunistic NOMA system (cooperative NOMA system with full-duplex relaying) and the OMA system becomes larger when the reception thresholds are smaller. The average service rate of queue Q_L can exceed 1 as NOMA can simultaneously serve two packets from queue Q_L when queue Q_H is empty.

Fig. 9 shows the impact of the number of low-priority users K on the full stable throughput region. For larger K , the maximum achievable λ_L of both NOMA schemes increases, as the probability that NOMA can serve the packets from queue Q_L increases. The maximum achievable λ_H of the opportunistic NOMA system does not depend on K .

The maximum achievable λ_H of the cooperative NOMA system with full-duplex relaying increases with K , as more low-priority users are available and the probability of selecting a reliable relay becomes higher. Thus, the full stable throughput regions of both NOMA systems increase with K .

Fig. 10 shows the impact of imperfect CSI estimation on the full stable throughput region. Adopting the model for imperfect CSI in [35], we have $h_i^a(t) = \hat{h}_i^a(t) + \epsilon_i^a(t)$, $i \in \{1, 2\}$, $a \in \{H, L\}$, where $\hat{h}_i^a(t)$ denotes the estimate of $h_i^a(t)$ at user u_i^a and $\epsilon_i^a(t)$ is the complex Gaussian channel estimation error at user u_i^a with zero mean and variance σ_ϵ^2 in time slot t . The value of variance σ_ϵ^2 reflects the accuracy of channel estimation. We first obtain the average service rates of queues Q_H and Q_L for all considered dominant systems via simulations, which are then used to plot the stable throughput regions in Fig. 10. Results show that, by increasing the value of σ_ϵ^2 from 0 to 0.01, the full stable throughput regions of all considered schemes become smaller. This is because channel estimation errors lead to additional interference as in the SINR expression, which reduces the probability of successful packet reception at each user. In addition, the impact of imperfect CSI on the performance of NOMA is greater compared to OMA, as the SIC at the user being allocated a lower transmit power in NOMA is negatively affected by imperfect CSI.

VII. CONCLUSION

In this paper, we studied the performance of downlink NOMA transmission with dynamic traffic arrival and spatially random users of different priorities. To reduce the adverse effect of NOMA on high-priority users, we proposed an opportunistic NOMA scheme requiring only limited CSI at the base station. Moreover, we proposed a cooperative NOMA scheme with full-duplex relaying, where the low-priority users assist the high-priority users to enhance the network performance. By utilizing tools from queueing theory and stochastic geometry, we characterized the stable throughput regions of both proposed NOMA schemes by constructing dominant systems to decouple the interacting queues. Simulation results validated the performance analysis. With appropriate parameter setting,

the proposed NOMA schemes can significantly improve the transmission opportunities and enhance the stable throughput region compared to OMA.

There are several interesting topics for future work. First, the performance analysis of cooperative NOMA with full-duplex relaying can be extended to multi-cell networks, where the interference from the base stations and the full-duplex relays in other cells has to be taken into account. Second, the proposed performance analysis framework can be extended to the case where each user exploits retransmission diversity [36]. Third, the proposed framework can be extended to the case where more than two users are paired for NOMA transmission.

APPENDIX A PROOF OF LEMMA 1

When OMA is enabled, the probability of successful packet reception at user u_1^H is given by

$$\begin{aligned} q_{H1}^{\text{OMA}}(\theta) &= \mathbb{P}\left(\frac{\rho_H}{\ell(x_1^H)} \leq |h_1^H(t)|^2 < \frac{\theta}{\ell(x_1^H)}\right) \\ &\stackrel{(a)}{=} \exp\left(-\rho_H(1+r_H^\beta)\right) - \exp\left(-\theta(1+r_H^\beta)\right), \end{aligned} \quad (38)$$

if $\theta > \rho_H$, otherwise, $q_{H1}^{\text{OMA}}(\theta) = 0$, where (a) follows from the exponential distribution of $|h_1^H(t)|^2$.

On the other hand, when NOMA is enabled, the probability of successful packet reception at high-priority user u_1^H can be expressed as

$$\begin{aligned} q_{H1|L1}^{\text{ON}}(\alpha_H, \theta) &= \mathbb{P}\left(|h_1^H(t)|^2 \geq \max\left\{\frac{\rho_H}{\alpha_H^2 - \Gamma_{\text{th}}^H \alpha_L^2}, \theta\right\} \frac{1}{\ell(x_1^H)}\right) \\ &= \exp\left(-\max\left\{\frac{\rho_H}{\alpha_H^2 - \Gamma_{\text{th}}^H \alpha_L^2}, \theta\right\} (1+r_H^\beta)\right), \end{aligned} \quad (39)$$

if $\alpha_H^2 > \Gamma_{\text{th}}^H \alpha_L^2$, otherwise, $q_{H1|L1}^{\text{ON}}(\alpha_H, \theta) = 0$.

By substituting (38) and (39) into (5), the average service rate of queue Q_H is given by

$$\mu_H^{\text{ON1}} = \begin{cases} q_{H1|L1}^{\text{ON}}(\alpha_H, \theta), & \text{if } \theta \leq \rho_H, \\ q_{H1}^{\text{OMA}}(\theta) + q_{H1|L1}^{\text{ON}}(\alpha_H, \theta), & \text{if } \theta > \rho_H. \end{cases} \quad (40)$$

According to (40), we set $\theta > \rho_H$ to achieve a higher value of μ_H^{ON1} . By Loynes' theorem [27], queue Q_H is stable if (6) holds.

APPENDIX B PROOF OF LEMMA 2

The probability of successful packet reception at low-priority user u_1^L when paired with high-priority user u_1^H to perform NOMA is given by

$$\begin{aligned} q_{L1|H1}^{\text{ON}}(\alpha_L) &= \mathbb{P}\left(\Gamma_{H1 \rightarrow L1}(t, \alpha_H) \geq \Gamma_{\text{th}}^H, \Gamma_{L1}(t, \alpha_L) \geq \Gamma_{\text{th}}^L\right) \\ &= \mathbb{P}\left(|h_1^L(t)|^2 \geq \frac{\rho_H}{(\alpha_H^2 - \Gamma_{\text{th}}^H \alpha_L^2) \ell(x_1^L)}, |h_1^L(t)|^2 \geq \frac{\rho_L}{\alpha_L^2 \ell(x_1^L)}\right) \\ &= \mathbb{E}_{x_1^L} \left[\exp\left(-N_1/\ell(x_1^L)\right) \right], \end{aligned} \quad (41)$$

where $\mathbb{E}_{x_1^L}[\cdot]$ denotes the expectation over location coordinate x_1^L of low-priority user u_1^L . The probability density function (PDF) of the location of low-priority user u_1^L is given by $f(x_1^L) = 1/(\pi r^2)$. Hence, we have

$$\begin{aligned} q_{L1|H1}^{\text{ON}}(\alpha_L) &= \frac{2}{r^2} \int_0^r \exp\left(-N_1(1+(r_1^L)^\beta)\right) r_1^L dr_1^L \\ &= \frac{2}{r^2 \beta} N_1^{-2/\beta} \exp(-N_1) \gamma\left(\frac{2}{\beta}, N_1 r^\beta\right). \end{aligned} \quad (42)$$

When queue Q_H is empty and the first two packets from queue Q_L are intended for different users, base station S transmits the first and second packets from queue Q_L using NOMA based on the distances between their intended receivers and the base station. Among these two users, the near and far users are denoted as u_n^L and u_f^L with distances r_n^L and r_f^L , respectively, and $r_n^L \leq r_f^L$. Hence, we have $\alpha_f \geq \alpha_n$. As users u_1^L and u_2^L follow the same location distribution, they have the same probability (i.e., 0.5) of being the near or far user. For instance, if $r_1^L \leq r_2^L$, then we have $u_n^L = u_1^L$ and $u_f^L = u_2^L$. Otherwise, we have $u_n^L = u_2^L$ and $u_f^L = u_1^L$. The paired NOMA users having the same low priority are ordered based on their distances to the base station. Due to the uniform distribution of users u_1^L and u_2^L , according to [37], the PDF of the distance between far user u_f^L and base station S is given by

$$f(r_f^L) = 4(r_f^L)^3 / r^4, \quad 0 \leq r_f^L \leq r. \quad (43)$$

When NOMA is enabled, the probability of successful packet reception at user u_f^L is given by

$$\begin{aligned} q_{Lf|Ln}^{\text{ON}}(\alpha_f) &= \mathbb{P}\left(\Gamma_{Lf|Ln}(t, \alpha_f) \geq \Gamma_{\text{th}}^L\right) \\ &= \mathbb{E}_{x_f^L} \left[\exp\left(-N_2/\ell(x_f^L)\right) \right] \\ &\stackrel{(a)}{=} \frac{4}{r^4} \int_0^r \exp\left(-N_2(1+(r_f^L)^\beta)\right) (r_f^L)^3 dr_f^L \\ &= \frac{4}{r^4 \beta} N_2^{-4/\beta} \exp(-N_2) \gamma\left(\frac{4}{\beta}, N_2 r^\beta\right), \end{aligned} \quad (44)$$

if $\alpha_f^2 > \Gamma_{\text{th}}^L \alpha_n^2$, otherwise, $q_{Lf|Ln}^{\text{ON}}(\alpha_f) = 0$, where (a) follows by substituting (43).

Similarly, the PDF of the distance between near user u_n^L and base station S is given by

$$f(r_n^L) = 4 \frac{r_n^L}{r^2} \left(1 - \frac{(r_n^L)^2}{r^2}\right), \quad 0 \leq r_n^L \leq r_L. \quad (45)$$

When NOMA is enabled, the probability of successful packet reception at user u_n^L is given by

$$\begin{aligned} q_{Ln|Lf}^{\text{ON}}(\alpha_n) &= \mathbb{P}\left(\Gamma_{Lf \rightarrow Ln}(t, \alpha_f) \geq \Gamma_{\text{th}}^L, \Gamma_{Ln}(t, \alpha_n) \geq \Gamma_{\text{th}}^L\right) \\ &= \mathbb{E}_{x_n^L} \left[\exp\left(-N_3/\ell(x_n^L)\right) \right] \\ &\stackrel{(a)}{=} \frac{4}{r^2} \int_0^r \exp\left(-N_3(1+(r_n^L)^\beta)\right) \left(r_n^L - \frac{(r_n^L)^3}{r^2}\right) dr_n^L \\ &= \frac{4}{r^2 \beta} N_3^{-2/\beta} \exp(-N_3) \gamma\left(\frac{2}{\beta}, N_3 r^\beta\right) \\ &\quad - \frac{4}{r^4 \beta} N_3^{-4/\beta} \gamma\left(\frac{4}{\beta}, N_3 r^\beta\right), \end{aligned} \quad (46)$$

1046 if $\alpha_f^2 > \Gamma_{th}^L \alpha_n^2$, otherwise, $q_{L_n|L_f}^{ON}(\alpha_n) = 0$, where (a) follows
 1047 by substituting (45).

1048 Based on (44) and (46), we have

$$1049 \quad q_{L1|L2}^{ON} = q_{L_f|L_n}^{ON}(\alpha_f) + q_{L_n|L_f}^{ON}(\alpha_n). \quad (47)$$

1050 When OMA is enabled, the probability of successful packet
 1051 reception at user u_1^L is given by

$$1052 \quad q_{L1}^{OMA} = \mathbb{P}(\Gamma_{L1}(t, 1) \geq \Gamma_{th}^L) \\
 1053 \quad = \frac{2}{r^2 \beta} \rho_L^{-2/\beta} \exp(-\rho_L) \gamma\left(\frac{2}{\beta}, \rho_L r^\beta\right). \quad (48)$$

1054 By substituting (9), (47), and (48) into (7), the average
 1055 service rate of queue Q_L in dominant system Φ_1^{ON} can be
 1056 derived. By Loynes' theorem, queue Q_L is stable if (8) holds.

1057 APPENDIX C 1058 PROOF OF THEOREM 1

1059 Our proof is based on a similar technique as the proofs
 1060 in [16]–[18]. The dominant systems (i.e., Φ_1^{ON} and Φ_2^{ON})
 1061 are modifications of the original opportunistic NOMA system
 1062 Φ^{ON} . The queue lengths in the dominant systems are never
 1063 shorter than the queue lengths in the original opportunistic
 1064 NOMA system Φ^{ON} as an empty queue can contribute dummy
 1065 packets. The transmission of dummy packets reduces the prob-
 1066 ability of successful packet reception by generating co-channel
 1067 interference, but does not contribute to the throughput. Hence,
 1068 the stability condition obtained for the dominant systems
 1069 (i.e., Φ_1^{ON} and Φ_2^{ON}) is sufficient for the stability of the
 1070 original opportunistic NOMA system Φ^{ON} .

1071 As only two queues are considered, the stability condition of
 1072 the original opportunistic NOMA system Φ^{ON} is determined
 1073 by the two parallel dominant systems (i.e., Φ_1^{ON} and Φ_2^{ON}).
 1074 In particular, dominant systems Φ_1^{ON} and Φ_2^{ON} explore all
 1075 possible choices of the average arrival rates λ_L and λ_H that
 1076 can lead to a stable system, respectively. In dominant system
 1077 Φ_1^{ON} , some λ_L would cause queue Q_L to be always non-empty.
 1078 As long as queue Q_L always has packets to transmit, queue
 1079 Q_L does not contribute dummy packets and hence the behavior
 1080 of dominant system Φ_1^{ON} is identical to that of the original
 1081 opportunistic NOMA system Φ^{ON} . As a result, dominant
 1082 system Φ_1^{ON} and the original opportunistic NOMA system
 1083 Φ^{ON} are indistinguishable at the boundary of the stability

1084 region (i.e., line CD in Fig. 2). Similarly, dominant system
 1085 Φ_2^{ON} and the original opportunistic NOMA system Φ^{ON} are
 1086 also indistinguishable at the boundary of the stability region
 1087 (i.e., line AC in Fig. 2). Similar indistinguishability argu-
 1088 ments are used in [16]–[18]. Thereby, the stability condition
 1089 obtained for the dominant systems (i.e., Φ_1^{ON} and Φ_2^{ON}) is also
 1090 necessary for the stability of the original opportunistic NOMA
 1091 system Φ^{ON} . As a result, we have $\mathcal{R}^{ON} = \mathcal{R}_1^{ON} \cup \mathcal{R}_2^{ON}$.

1092 APPENDIX D 1093 PROOF OF LEMMA 3

1094 Due to the independence of events $\{\Gamma_{H1|L1}(t, \alpha_H) \geq \Gamma_{th}^H\}$
 1095 and $\{\Gamma_{H1 \rightarrow L1}^{FCN}(t, \alpha_H) < \Gamma_{th}^H\}$, the probability of successful
 1096 packet reception at user u_1^H when NOMA is enabled, denoted
 1097 as $q_{H1}^N(\alpha_H)$, is given by

$$1098 \quad q_{H1}^N(\alpha_H) \\
 1099 \quad = \mathbb{P}(\Gamma_{H1|L1}(t, \alpha_H) \geq \Gamma_{th}^H) \mathbb{P}(\Gamma_{H1 \rightarrow L1}^{FCN}(t, \alpha_H) < \Gamma_{th}^H) \\
 1100 \quad = \exp\left(-\frac{\rho_H(1+r_H^\beta)}{\alpha_H^2 - \Gamma_{th}^H \alpha_L^2}\right) \mathbb{E}_{x_1^L} \left[1 - \exp\left(-\frac{N_4}{\ell(x_1^L)}\right)\right] \\
 1101 \quad = \exp\left(-\frac{\rho_H(1+r_H^\beta)}{\alpha_H^2 - \Gamma_{th}^H \alpha_L^2}\right) \\
 1102 \quad \times \left(1 - \frac{2}{r^2 \beta} N_4^{-2/\beta} \exp(-N_4) \gamma\left(\frac{2}{\beta}, N_4 r^\beta\right)\right), \quad (49)$$

1103 if $\alpha_H^2 > \Gamma_{th}^H \alpha_L^2$, otherwise, $q_{H1}^N(\alpha_H) = 0$.

1104 The probability of successful packet reception at user u_1^H
 1105 when cooperative NOMA is enabled, denoted as $q_{H1}^{FCN}(\alpha_H)$,
 1106 can be expressed as

$$1107 \quad q_{H1}^{FCN}(\alpha_H) \stackrel{(a)}{=} \mathbb{E}_{x_1^L} \left[\mathbb{P}\left(\left(\alpha_H^2 - \Gamma_{th}^H \alpha_L^2\right) P_S |h_1^H(t)|^2 \ell(x_1^H) \right. \right. \\
 1108 \quad \left. \left. + P_L |g_{1,1}^{HL}(t)|^2 \ell(x_1^H - x_1^L) \geq \Gamma_{th}^H \sigma^2\right) \right. \\
 1109 \quad \left. \times \mathbb{P}\left(|h_1^L(t)|^2 \geq \frac{N_4}{\ell(x_1^L)}\right) \right], \quad (50)$$

$$\begin{aligned} & \mathbb{P}\left(Z |h_1^H(t)|^2 + P_L |g_{1,1}^{HL}(t)|^2 \ell(x_1^H - x_1^L) \geq \Gamma_{th}^H \sigma^2\right) \\ &= \int_0^{\frac{\Gamma_{th}^H \sigma^2}{Z}} \int_{\frac{\Gamma_{th}^H \sigma^2 - |h_1^H(t)|^2}{P_L \ell(x_1^H - x_1^L)}}^\infty \exp\left(-|h_1^H(t)|^2\right) \exp\left(-|g_{1,1}^{HL}(t)|^2\right) d|g_{1,1}^{HL}(t)|^2 d|h_1^H(t)|^2 \\ &+ \int_{\frac{\Gamma_{th}^H \sigma^2}{Z}}^\infty \int_0^\infty \exp\left(-|h_1^H(t)|^2\right) \exp\left(-|g_{1,1}^{HL}(t)|^2\right) d|g_{1,1}^{HL}(t)|^2 d|h_1^H(t)|^2 \\ &= \begin{cases} \frac{1}{1 - \frac{Z}{P_L \ell(x_1^H - x_1^L)}} \left(\exp\left(-\frac{\Gamma_{th}^H \sigma^2}{P_L \ell(x_1^H - x_1^L)}\right) - \exp\left(-\frac{\Gamma_{th}^H \sigma^2}{Z}\right) \right) + \exp\left(-\frac{\Gamma_{th}^H \sigma^2}{Z}\right), & \text{if } Z \neq P_L \ell(x_1^H - x_1^L), \\ \frac{\Gamma_{th}^H \sigma^2}{Z} \exp\left(-\frac{\Gamma_{th}^H \sigma^2}{P_L \ell(x_1^H - x_1^L)}\right), & \text{if } Z = P_L \ell(x_1^H - x_1^L). \end{cases} \end{aligned} \quad (51)$$

$$\begin{aligned}
 & \mathbb{P}(\Gamma_{\text{H1} \rightarrow \text{L1}}^{\text{FCN}}(t, \alpha_{\text{H}}) < \Gamma_{\text{th}}^{\text{H}}, \Gamma_{\text{H1} \rightarrow \text{L1}}(t, \alpha_{\text{H}}) \geq \Gamma_{\text{th}}^{\text{H}}, \Gamma_{\text{L1}}^{\text{FCN}}(t, \alpha_{\text{L}}) \geq \Gamma_{\text{th}}^{\text{L}}) \\
 &= \mathbb{P}\left(\max\left\{\frac{\Gamma_{\text{th}}^{\text{H}}}{\alpha_{\text{H}}^2 - \Gamma_{\text{th}}^{\text{H}} \alpha_{\text{L}}^2}, \frac{\Gamma_{\text{th}}^{\text{L}}}{\alpha_{\text{L}}^2}\right\} \frac{\sigma^2}{P_S \ell(x_1^{\text{L}})} \leq |h_1^{\text{L}}(t)|^2 < \frac{(\zeta P_{\text{L}} + \sigma^2) \Gamma_{\text{th}}^{\text{H}}}{(\alpha_{\text{H}}^2 - \Gamma_{\text{th}}^{\text{H}} \alpha_{\text{L}}^2) P_S \ell(x_1^{\text{L}})}\right) \\
 &= \mathbb{E}_{x_1^{\text{L}}}\left[\exp\left(-\max\left\{\frac{\rho_{\text{H}}}{\alpha_{\text{H}}^2 - \Gamma_{\text{th}}^{\text{H}} \alpha_{\text{L}}^2}, \frac{\rho_{\text{L}}}{\alpha_{\text{L}}^2}\right\} \frac{1}{\ell(x_1^{\text{L}})}\right) - \exp\left(-\frac{(\zeta P_{\text{L}} + \sigma^2) \Gamma_{\text{th}}^{\text{H}}}{(\alpha_{\text{H}}^2 - \Gamma_{\text{th}}^{\text{H}} \alpha_{\text{L}}^2) P_S \ell(x_1^{\text{L}})}\right)\right] \\
 &= \frac{2}{r^2 \beta} N_1^{-2/\beta} \exp(-N_1) \gamma\left(\frac{2}{\beta}, N_1 r^\beta\right) - \frac{2}{r^2 \beta} N_4^{-2/\beta} \exp(-N_4) \gamma\left(\frac{2}{\beta}, N_4 r^\beta\right), \tag{53}
 \end{aligned}$$

1110 where (a) follows from the independent channel fading
 1111 assumption across different links.

1112 By conditioning on location coordinate x_1^{L} , we obtain (51),
 1113 as shown at the bottom of previous page.

1114 By substituting (51) into (50), we have

$$\begin{aligned}
 & q_{\text{H1}}^{\text{FCN}}(\alpha_{\text{H}}) \\
 &= \mathbb{E}_{x_1^{\text{L}}}\left[\left(\frac{\left(\exp\left(-\frac{\Gamma_{\text{th}}^{\text{H}} \sigma^2}{P_{\text{L}} \ell(x_1^{\text{H}} - x_1^{\text{L}})}\right) - N_5\right)}{1 - \frac{Z}{P_{\text{L}} \ell(x_1^{\text{H}} - x_1^{\text{L}})}} + N_5\right)\right. \\
 &\quad \left. \times \exp\left(-\frac{N_4}{\ell(x_1^{\text{L}})}\right)\right] \\
 &= C(\alpha_{\text{H}}) + \frac{2N_5}{r^2 \beta} N_4^{-2/\beta} \exp(-N_4) \gamma\left(\frac{2}{\beta}, N_4 r^\beta\right), \tag{52}
 \end{aligned}$$

1119 where $C(\alpha_{\text{H}})$ is given in (21), which can be calculated
 1120 numerically using commercial software (e.g., Mathematica).
 1121 By substituting (49) and (52) into (18), we obtain the average
 1122 service rate of queue Q_{H} . Hence, queue Q_{H} is stable if
 1123 $\lambda_{\text{H}} < \mu_{\text{H}}^{\text{FCN1}} = q_{\text{H1}}^{\text{N}}(\alpha_{\text{H}}) + q_{\text{H1}}^{\text{FCN}}(\alpha_{\text{H}})$.

APPENDIX E PROOF OF LEMMA 4

1126 The probability of successful packet reception at user u_1^{L}
 1127 when cooperative NOMA is enabled, if $\alpha_{\text{H}}^2 > \Gamma_{\text{th}} \alpha_{\text{L}}^2$, can be
 1128 expressed as

$$\begin{aligned}
 & \mathbb{P}(\Gamma_{\text{H1} \rightarrow \text{L1}}^{\text{FCN}}(t, \alpha_{\text{H}}) \geq \Gamma_{\text{th}}^{\text{H}}, \Gamma_{\text{L1}}^{\text{FCN}}(t, \alpha_{\text{L}}) \geq \Gamma_{\text{th}}^{\text{L}}) \\
 &= \mathbb{E}_{x_1^{\text{L}}}\left[\exp\left(-\max\left\{\frac{\Gamma_{\text{th}}^{\text{H}}}{\alpha_{\text{H}}^2 - \Gamma_{\text{th}}^{\text{H}} \alpha_{\text{L}}^2}, \frac{\Gamma_{\text{th}}^{\text{L}}}{\alpha_{\text{L}}^2}\right\} \frac{(\zeta P_{\text{L}} + \sigma^2)}{P_S \ell(x_1^{\text{L}})}\right)\right] \\
 &= \frac{2}{r^2 \beta} N_6^{-2/\beta} \exp(-N_6) \gamma\left(\frac{2}{\beta}, N_6 r^\beta\right).
 \end{aligned}$$

1132 When NOMA is enabled, the probability of successful
 1133 packet reception at user u_1^{L} is given by (53), as shown
 1134 at the top of this page, where $\alpha_{\text{H}}^2 > \Gamma_{\text{th}} \alpha_{\text{L}}^2$. By substi-
 1135 tuting (24), (47), and (48) into (22), we can obtain the average
 1136 service rate of queue Q_{L} . Hence, queue Q_{L} is stable if
 1137 $\lambda_{\text{L}} < \mu_{\text{L}}^{\text{FCN1}} = \frac{\lambda_{\text{H}}}{\mu_{\text{FCN1}}^{\text{H}}} q_{\text{L1}|\text{H1}}^{\text{FCN}}(\alpha_{\text{L}}) + \left(1 - \frac{\lambda_{\text{H}}}{\mu_{\text{FCN1}}^{\text{H}}}\right) \eta$, where
 1138 η is given in (10).

APPENDIX F PROOF OF LEMMA 5

Given that there are K low-priority users, we have

$$\begin{aligned}
 & q_{\text{H1}}^{\text{FC}} = 1 - \mathbb{P}(\Gamma_{\text{H1}}(t, 1) < \Gamma_{\text{th}}^{\text{H}}) \\
 &\quad \times \mathbb{E}_{x_k^{\text{R}}}\left[\prod_{u_k^{\text{R}} \in \mathcal{U}^{\text{L}}} \left[1 - \mathbb{P}(\Gamma_{\text{H1} \rightarrow \text{Rk}}^{\text{FC}}(t) \geq \Gamma_{\text{th}}^{\text{H}})\right.\right. \\
 &\quad \left.\left. \times \mathbb{P}(\Gamma_{\text{H1Rk}}^{\text{FC}}(t) \geq \Gamma_{\text{th}}^{\text{H}} | \Gamma_{\text{H1}}(t, 1) < \Gamma_{\text{th}}^{\text{H}})\right]\right], \tag{54}
 \end{aligned}$$

1145 where $\mathbb{P}(\Gamma_{\text{H1}}(t, 1) < \Gamma_{\text{th}}^{\text{H}}) = 1 - \exp(-\rho_{\text{H}}(1 + r_{\text{H}}^\beta))$.
 1146 As the K low-priority users are uniformly distributed within
 1147 the network coverage area, we have

$$\begin{aligned}
 & \mathbb{E}_{x_k^{\text{R}}}\left[\prod_{u_k^{\text{R}} \in \mathcal{U}^{\text{L}}} \left[1 - \mathbb{P}(\Gamma_{\text{H1} \rightarrow \text{Rk}}^{\text{FC}}(t) \geq \Gamma_{\text{th}}^{\text{H}})\right.\right. \\
 &\quad \left.\left. \times \mathbb{P}(\Gamma_{\text{H1Rk}}^{\text{FC}}(t) \geq \Gamma_{\text{th}}^{\text{H}} | \Gamma_{\text{H1}}(t, 1) < \Gamma_{\text{th}}^{\text{H}})\right]\right] \\
 &\stackrel{(a)}{=} \left(\frac{1}{\pi r^2} \int_0^r \int_0^{2\pi} \left(1 - \mathbb{P}(\Gamma_{\text{H1} \rightarrow \text{Rk}}^{\text{FC}}(t) \geq \Gamma_{\text{th}}^{\text{H}})\right.\right. \\
 &\quad \left.\left. \times \mathbb{P}(\Gamma_{\text{H1Rk}}^{\text{FC}}(t) \geq \Gamma_{\text{th}}^{\text{H}} | \Gamma_{\text{H1}}(t, 1) < \Gamma_{\text{th}}^{\text{H}})\right) r_k^{\text{R}} d\tau_k^{\text{R}} dr_k^{\text{R}}\right)^K \\
 &\stackrel{(b)}{=} \left(\frac{1}{\pi r^2} \int_0^r \int_0^{2\pi} \left(1 - \mathbb{P}(\Gamma_{\text{H1} \rightarrow \text{Rk}}^{\text{FC}}(t) \geq \Gamma_{\text{th}}^{\text{H}})\right.\right. \\
 &\quad \left.\left. \times \frac{\mathbb{P}(\Gamma_{\text{H1Rk}}^{\text{FC}}(t) \geq \Gamma_{\text{th}}^{\text{H}}, \Gamma_{\text{H1}}(t, 1) < \Gamma_{\text{th}}^{\text{H}})}{\mathbb{P}(\Gamma_{\text{H1}}(t, 1) < \Gamma_{\text{th}}^{\text{H}})}\right) r_k^{\text{R}} d\tau_k^{\text{R}} dr_k^{\text{R}}\right)^K \\
 &\stackrel{(c)}{=} \left(1 - \frac{C(1)}{1 - \exp(-\rho_{\text{H}}(1 + r_{\text{H}}^\beta))}\right)^K, \tag{55}
 \end{aligned}$$

1155 where (a) follows from the probability generating func-
 1156 tional (PGFL) of the BPP [22], (b) follows from the definition
 1157 of conditional probability, (c) is obtained using similar steps
 1158 as for deriving (50), and $C(1)$ is given in Lemma 3 by setting
 1159 $\alpha_{\text{H}}^2 = 1$ and $\alpha_{\text{L}}^2 = 0$.

By substituting (55) into (54), we have

$$\begin{aligned}
 q_{H1}^{FC} &= 1 - \left(1 - \exp\left(-\rho_H \left(1 + r_H^\beta\right)\right) \right) \\
 &\times \left(1 - \frac{C(1)}{1 - \exp\left(-\rho_H \left(1 + r_H^\beta\right)\right)} \right)^K \\
 &\stackrel{(a)}{=} \exp\left(-\rho_H \left(1 + r_H^\beta\right)\right) + \sum_{j=1}^K \binom{K}{j} (-1)^{j+1} (C(1))^j \\
 &\times \left(1 - \exp\left(-\rho_H \left(1 + r_H^\beta\right)\right) \right)^{1-j}, \quad (56)
 \end{aligned}$$

where (a) follows from the binomial expansion. By substituting (56) into (27), we can derive the average service rate of queue Q_H in dominant system Φ_2^{FCN} . Hence, queue Q_H is stable if $\lambda_H < \mu_H^{FCN2} = (1 - \lambda_L/\mu_L^{FCN2}) q_{H1}^{FC} + \lambda_L/\mu_L^{FCN2} (q_{H1}^{CN}(\alpha_H) + q_{H1}^{FCN}(\alpha_H))$.

APPENDIX G PROOF OF PROPOSITION 1

Based on Fig. 2, in order for $\mathcal{R}^{OMA} \subset \mathcal{R}^{FCN}$ to hold, points D, E, and F have to be on the right side of line AB. To ensure that point E is on the right side of line AB, condition $q_{H1}^{FC} > \mu_H^{OMA}$ should hold, which is obtained based on the coordinates of points A and E. According to (56), condition $q_{H1}^{FC} > \mu_H^{OMA}$ always holds. In addition, to guarantee that point F is on the right side of line AB, the Y-coordinate of point F should be larger than the Y-coordinate of the point that is on line AB and has the same X-coordinate as point F. Hence, condition $q_{L1|H1}^{FCN}(\alpha_L) > q_{L1}^{OMA} \left(1 - \frac{q_{H1}^{CN}(\alpha_H) + q_{H1}^{FCN}(\alpha_H)}{\mu_H^{OMA}} \right)$ should hold, where $\alpha_L^2 = 1 - \alpha_H^2$. Similarly, to ensure that point D is on the right side of line AB, condition $\eta > q_{L1}^{OMA}$, equivalent to $q_{L1L2}^{ON} > q_{L1}^{OMA}$, should hold, where q_{L1L2}^{ON} is given in (47). As a result, cooperative NOMA with full-duplex relaying achieves a larger stable throughput region than OMA if both (36) and (37) hold.

REFERENCES

- [1] Y. Zhou and V. W.S. Wong, "Stable throughput region of downlink NOMA transmissions with limited CSI," in *Proc. IEEE ICC*, Paris, France, May 2017, pp. 1–7.
- [2] Y. Saito, A. Benjebbour, Y. Kishiyama, and T. Nakamura, "System-level performance evaluation of downlink non-orthogonal multiple access (NOMA)," in *Proc. IEEE PIMRC*, London, U.K., Sep. 2013, pp. 611–615.
- [3] A. Benjebbour, A. Li, Y. Saito, Y. Kishiyama, A. Harada, and T. Nakamura, "System-level performance of downlink NOMA for future LTE enhancements," in *Proc. IEEE Globecom Workshops*, Atlanta, GA, USA, Dec. 2013, pp. 66–70.
- [4] J. Choi, "On the power allocation for MIMO-NOMA systems with layered transmissions," *IEEE Trans. Wireless Commun.*, vol. 15, no. 5, pp. 3226–3237, May 2016.
- [5] B. Di, S. Bayat, L. Song, and Y. Li, "Radio resource allocation for downlink non-orthogonal multiple access (NOMA) networks using matching theory," in *Proc. IEEE Globecom*, San Diego, CA, USA, Dec. 2015, pp. 1–6.
- [6] Z. Ding, P. Fan, and H. V. Poor, "Impact of user pairing on 5G non-orthogonal multiple access downlink transmissions," *IEEE Trans. Veh. Technol.*, vol. 65, no. 8, pp. 6010–6023, Aug. 2016.
- [7] Z. Ding, R. Schober, and H. V. Poor, "A general MIMO framework for NOMA downlink and uplink transmission based on signal alignment," *IEEE Trans. Wireless Commun.*, vol. 15, no. 6, pp. 4438–4454, Jun. 2016.
- [8] Z. Ding, H. Dai, and H. V. Poor, "Relay selection for cooperative NOMA," *IEEE Wireless Commun. Lett.*, vol. 5, no. 4, pp. 416–419, Aug. 2016.
- [9] Z. Ding, P. Fan, and H. V. Poor, "Random beamforming in millimeter-wave NOMA networks," *IEEE Access*, vol. 5, pp. 7667–7681, 2017.
- [10] D. Kim, H. Lee, and D. Hong, "A survey of in-band full-duplex transmission: From the perspective of PHY and MAC layers," *IEEE Commun. Surveys Tuts.*, vol. 17, no. 4, pp. 2017–2046, 4th Quart., 2015.
- [11] M. Jain *et al.*, "Practical, real-time, full duplex wireless," in *Proc. ACM MobiCom*, Las Vegas, NV, USA, Sep. 2011, pp. 301–312.
- [12] W. Li, J. Lilleberg, and K. Rikkinen, "On rate region analysis of half- and full-duplex OFDM communication links," *IEEE J. Sel. Areas Commun.*, vol. 32, no. 9, pp. 1688–1698, Sep. 2014.
- [13] Y. Sun, D. W. K. Ng, Z. Ding, and R. Schober, "Optimal joint power and subcarrier allocation for full-duplex multicarrier non-orthogonal multiple access systems," *IEEE Trans. Commun.*, vol. 65, no. 3, pp. 1077–1091, Mar. 2017.
- [14] T. Riihonen, S. Werner, and R. Wichman, "Hybrid full-duplex/half-duplex relaying with transmit power adaptation," *IEEE Trans. Wireless Commun.*, vol. 10, no. 9, pp. 3074–3085, Sep. 2011.
- [15] A. Sharma, R. K. Ganti, and J. K. Milleth, "Joint backhaul-access analysis of full duplex self-backhauling heterogeneous networks," *IEEE Trans. Wireless Commun.*, vol. 16, no. 3, pp. 1727–1740, Mar. 2017.
- [16] O. Simeone, Y. Bar-Ness, and U. Spagnolini, "Stable throughput of cognitive radios with and without relaying capability," *IEEE Trans. Commun.*, vol. 55, no. 12, pp. 2351–2360, Dec. 2007.
- [17] C. Kam, S. Kompella, G. D. Nguyen, J. E. Wieselthier, and A. Ephremides, "Cognitive cooperative random access for multicast: Stability and throughput analysis," *IEEE Trans. Control Netw. Syst.*, vol. 1, no. 2, pp. 135–144, Jun. 2014.
- [18] S. Kompella, G. Nguyen, C. Kam, J. E. Wieselthier, and A. Ephremides, "Cooperation in cognitive underlay networks: Stable throughput trade-offs," *IEEE/ACM Trans. Netw.*, vol. 22, no. 6, pp. 1756–1768, Dec. 2014.
- [19] I. Krikidis, J. N. Laneman, J. Thompson, and S. McLaughlin, "Protocol design and throughput analysis for multi-user cognitive cooperative systems," *IEEE Trans. Wireless Commun.*, vol. 8, no. 9, pp. 4740–4751, Sep. 2009.
- [20] P. Xu, Y. Yuan, Z. Ding, X. Dai, and R. Schober, "On the outage performance of non-orthogonal multiple access with 1-bit feedback," *IEEE Trans. Wireless Commun.*, vol. 15, no. 10, pp. 6716–6730, Oct. 2016.
- [21] M. Haenggi and R. K. Ganti, *Interference in Large Wireless Networks*. Boston, MA, USA: Now, 2009.
- [22] M. Haenggi, *Stochastic Geometry for Wireless Networks*. Cambridge, U.K.: Cambridge Univ. Press, 2012.
- [23] P. H. J. Nardelli, M. Kountouris, P. Cardieri, and M. Latva-Aho, "Throughput optimization in wireless networks under stability and packet loss constraints," *IEEE Trans. Mobile Comput.*, vol. 13, no. 8, pp. 1883–1895, Aug. 2014.
- [24] Y. Zhou and W. Zhuang, "Performance analysis of cooperative communication in decentralized wireless networks with unsaturated traffic," *IEEE Trans. Wireless Commun.*, vol. 15, no. 5, pp. 3518–3530, May 2016.
- [25] K. Stamatiou and M. Haenggi, "Random-access Poisson networks: Stability and delay," *IEEE Commun. Lett.*, vol. 14, no. 11, pp. 1035–1037, Nov. 2010.
- [26] S. Weber and J. G. Andrews, "Transmission capacity of wireless networks," *Found. Trends Netw.*, vol. 5, nos. 2–3, pp. 109–281, 2012, doi: <http://dx.doi.org/10.1561/1300000032>.
- [27] R. M. Loynes, "The stability of a queue with non-independent inter-arrival and service times," *Proc. Cambridge Philos. Soc.*, vol. 58, no. 3, pp. 497–520, 1962.
- [28] *Study on Downlink Multiuser Superposition Transmission (MUST) for LTE*, document 3GPP TR 36.859, Jan. 2016.
- [29] M. Haenggi, "The local delay in Poisson networks," *IEEE Trans. Inf. Theory*, vol. 59, no. 3, pp. 1788–1802, Mar. 2013.
- [30] F. Baccelli and B. Błaszczyszyn, *Stochastic Geometry Wireless Networks, Part I Theory*. Boston, MA, USA: Now, 2009.
- [31] W. Luo and A. Ephremides, "Stability of N interacting queues in random-access systems," *IEEE Trans. Inf. Theory*, vol. 45, no. 5, pp. 1579–1587, Jul. 1999.
- [32] I. S. Gradshteyn and I. M. Ryzhik, *Table of Integrals, Series, and Products*. New York, NY, USA: Academic, 2014.
- [33] D. Bharadia and S. Katti, "FastForward: Fast and constructive full duplex relays," in *Proc. ACM SIGCOMM*, Chicago, IL, USA, Aug. 2014, pp. 1–12.

1289 [34] K.-C. Hsu, K. C.-J. Lin, and H.-Y. Wei, "Full-duplex delay-and-forward
1290 relaying," in *Proc. ACM MobiHoc*, Paderborn, Germany, Jul. 2016,
1291 pp. 221–230.

1292 [35] Z. Yang, Z. Ding, P. Fan, and G. K. Karagiannidis, "On the performance
1293 of non-orthogonal multiple access systems with partial channel informa-
1294 tion," *IEEE Trans. Commun.*, vol. 64, no. 2, pp. 654–667, Feb. 2016.

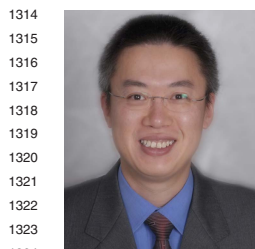
1295 [36] M. K. Tsatsanis, R. Zhang, and S. Banerjee, "Network-assisted diversity
1296 for random access wireless networks," *IEEE Trans. Signal Process.*,
1297 vol. 48, no. 3, pp. 702–711, Mar. 2000.

1298 [37] S. Srinivasa and M. Haenggi, "Distance distributions in finite uniformly
1299 random networks: Theory and applications," *IEEE Trans. Veh. Technol.*,
1300 vol. 59, no. 2, pp. 940–949, Feb. 2010.

the IEEE SmartGridComm 2013, the IEEE SmartGridComm 2017, and
the IEEE Globecom 2013. He is the Chair of the IEEE Vancouver Joint
Communications Chapter. He served as a Chair for the IEEE Communications
Society Emerging Technical Sub-Committee on Smart Grid Communications.



Yong Zhou (S'13–M'16) received the B.Sc. and M.Eng. degrees from Shandong University, Jinan, China, in 2008 and 2011, respectively, and the Ph.D. degree from the University of Waterloo, Waterloo, ON, Canada, in 2015. From 2015 to 2017, he was a Post-Doctoral Research Fellow with the Department of Electrical and Computer Engineering, The University of British Columbia, Vancouver, Canada. He is currently an Assistant Professor with the School of Information Science and Technology, ShanghaiTech University, China. His research interests include performance analysis and resource allocation of 5G networks. He served as a technical program committee member for several conferences.



Vincent W.S. Wong (S'94–M'00–SM'07–F'16) received the B.Sc. degree from the University of Manitoba, Winnipeg, MB, Canada, in 1994, the M.A.Sc. degree from the University of Waterloo, Waterloo, ON, Canada, in 1996, and the Ph.D. degree from The University of British Columbia (UBC), Vancouver, BC, Canada, in 2000. From 2000 to 2001, he was a Systems Engineer with Microsemi. He joined the Department of Electrical and Computer Engineering, UBC, in 2002, where he is currently a Professor. His research areas include protocol design, optimization, and resource management of communication networks, with applications to wireless networks, smart grid, mobile cloud computing, and Internet of Things. He received the 2014 UBC Killam Faculty Research Fellowship. He is an Editor of the IEEE TRANSACTIONS ON COMMUNICATIONS. He served as a Guest Editor for the IEEE JOURNAL ON SELECTED AREAS IN COMMUNICATIONS and the IEEE WIRELESS COMMUNICATIONS. He also served on the editorial boards for the IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY and the *Journal of Communications and Networks*. He was a Technical Program Co-Chair of the IEEE SmartGridComm 2014 and a Symposium Co-Chair of the IEEE ICC 2018,



Robert Schober (S'98–M'01–SM'08–F'10) received the Diploma and Ph.D. degrees in electrical engineering from the Friedrich-Alexander University of Erlangen–Nuremberg (FAU), Germany, in 1997 and 2000, respectively. From 2002 to 2011, he was a Professor and the Canada Research Chair with The University of British Columbia (UBC), Vancouver, Canada. Since 2012, he has been an Alexander von Humboldt Professor and the Chair for Digital Communication with FAU. His research interests fall into the broad areas

of communication theory, wireless communications, and statistical signal processing.

Dr. Schober is a fellow of The Canadian Academy of Engineering and The Engineering Institute of Canada. He received several awards for his research, including the 2002 Heinz Maier-Leibnitz Award of the German Science Foundation (DFG), the 2004 Innovations Award of the Vodafone Foundation for Research in Mobile Communications, the 2006 UBC Killam Research Prize, the 2007 Wilhelm Friedrich Bessel Research Award of the Alexander von Humboldt Foundation, the 2008 Charles McDowell Award for Excellence in Research from UBC, the 2011 Alexander von Humboldt Professorship, the 2012 NSERC E. W. R. Steacie Fellowship, and the 2017 Wireless Communications Recognition Award by the IEEE Wireless Communications Technical Committee. He is listed as a 2017 Highly Cited Researcher by the Web of Science and a Distinguished Lecturer by the IEEE Communications Society (ComSoc). From 2012 to 2015, he served as the Editor-in-Chief for the IEEE TRANSACTIONS ON COMMUNICATIONS. He is currently the Chair of the Steering Committee of the IEEE TRANSACTIONS ON MOLECULAR, BIOLOGICAL AND MULTI-SCALE COMMUNICATION, an Editorial Board Member of the PROCEEDINGS OF THE IEEE, a Member at Large of the Board of Governors of ComSoc, and the ComSoc Director of Journals.