



Hybrid NOMA/OMA with Buffer-Aided Relay Selection in Cooperative Networks

Nikolaos Nomikos, *Member, IEEE*, Themistoklis Charalambous, *Member, IEEE*, Demosthenes Vouyioukas, *Senior Member, IEEE*, George K. Karagiannidis, *Fellow, IEEE*, and Risto Wichman, *Member, IEEE*

Abstract—Non-orthogonal multiple access (NOMA) aims to increase the spectral efficiency of fifth generation (5G) networks by relaxing the orthogonal use of radio-resources. In this work, a network with multiple half-duplex (HD) buffer-aided (BA) relays is considered, where the source transmits with fixed rate towards two users. The users might demand the same rate by the source (e.g., two cellular users requiring the same service), or they could have different rate requirements (e.g., a cellular user coexisting with an Internet of Things (IoT) device). By deploying multiple BA relays, increased reliability and additional degrees of freedom are provided. Leveraging the spectral efficiency of NOMA and the increased diversity gain of BA relaying, two relay selection algorithms with broadcasting are proposed for power-domain (PD) NOMA and hybrid NOMA/OMA, namely BA-NOMA and BA-NOMA/OMA, respectively. BA-NOMA can improve the performance in terms of outage probability when the power allocation factor α is selected such that robustness against channel uncertainties due to, e.g., outdated channel state information (CSI), is provided. Moreover, BA-NOMA/OMA further improves the sum-rate by switching to OMA when the relays cannot serve the users through NOMA. For both cases, a theoretical analysis for the outage probability is conducted and the asymptotic performance is studied. Finally, numerical results and comparisons with other state-of-the-art algorithms are provided for the outage probability, average throughput and average delay.

Index Terms—NOMA, Hybrid NOMA/OMA, BA relays, relay selection, cooperative diversity, spectral efficiency.

I. INTRODUCTION

A. Background

The expected success of 5G networks relies, not only on traditional communication practices, but also on novel communication paradigms, thus deviating from the current trends. Among those paradigms, non-orthogonal multiple access (NOMA) differentiates from traditional orthogonal multiple access (OMA), allowing users to share the same resource (time/frequency/code) and exploit different channel power levels [1]. Under NOMA, at the transmitter, superposition coding is deployed [2], while at the receivers, successive

interference cancellation (SIC) is performed [3], [4]. The surveys in [5], [6] presented various challenges of NOMA, including the need for efficient user pairing, based on different channel quality for optimal NOMA, while the potential of NOMA for increasing the resource efficiency was underlined. Regarding resource allocation in NOMA, various works have provided efficient solutions. The authors in [7] studied joint sub-channel assignment and power allocation in the downlink of a NOMA network, aiming to maximize the weighted sum-rate and maintaining user fairness. In addition, spectral and power resource allocation optimization, towards increasing the energy efficiency of the network, was presented in [8]. As power allocation to multiple users is a key issue in NOMA networks, several surveys studied various power-domain (PD) NOMA scenarios, emphasizing the need for the development of novel cooperation techniques, facilitating NOMA when relays are available in multi-user networks [9]–[11]. Also, an important practical issue raised by the aforementioned works is the deteriorated performance of NOMA under imperfect/outdated channel state information (CSI) [9]. Furthermore, when only statistical CSI is available, NOMA was shown to outperform OMA when power allocation and decoding order are optimized under the max-min fairness criterion [12]. Finally, when multiple antennas are available at the source, precoding and the minorization-maximization algorithm were used to maximize the sum-rate of single-antenna destination [13], while for multi-antenna destinations, multiple-input multiple output (MIMO) NOMA was proved to provide improved performance compared to MIMO OMA [14].

Another imminent communication paradigm is cooperative relaying for mitigating the effects of path-loss, shadowing and multi-path fading. The seminal work of [15], triggered a furore of contributions on relay techniques. Among those techniques, opportunistic relay selection (ORS) [16], [17] and buffer-aided (BA) relaying (see, e.g., [18], [19] and references therein) are of significant importance, as relay-assisted cellular networks can enhance the coverage of 5G networks. In addition, BA-ORS reduces the outage probability of the network, especially in delay-tolerant applications. Towards this end, [20] proposed hybrid relay selection (HRS) based on the max-max relay selection (MMRS), allocating one time-slot to two relays, having the strongest source-relay ($\{S \rightarrow R\}$) and relay-destination ($\{R \rightarrow D\}$) links. Then, max-link algorithm with adaptive link selection was proposed in [21], providing increased diversity, by dedicating each time-slot to either an $\{S \rightarrow R\}$ or to an $\{R \rightarrow D\}$ transmission. A hybrid MMRS/max-link algorithm was presented in [22], while

N. Nomikos and D. Vouyioukas are with the Department of Information and Communication Systems Engineering, University of the Aegean, Samos, Greece. *Emails:* {nnomikos, dvouyiou}@aegean.gr.

T. Charalambous is with the Department of Electrical Engineering and Automation, School of Electrical Engineering, Aalto University, Espoo, Finland. *Email:* themistoklis.charalambous@aalto.fi.

G. K. Karagiannidis is with the Department of Electrical and Computer Engineering, Aristotle University of Thessaloniki, Thessaloniki, Greece. *Email:* geokarag@auth.gr.

R. Wichman is with the Department of Signal Processing and Acoustics, School of Electrical Engineering, Aalto University, Espoo, Finland. *Email:* risto.wichman@aalto.fi.

other works provided delay-aware algorithms (see, e.g., [23]–[28]), prioritizing $\{R \rightarrow D\}$ transmissions and/or considering buffer state information (BSI) for relay selection. Moreover, to further improve the diversity of BA relaying, various works imposed thresholds at the buffers, in terms of the number of packets, activating relays satisfying those thresholds [29], [30]. Recently, $\{S \rightarrow R\}$ broadcasting for BA-ORS was proposed in [31], aiming to reduce the delay of MMRS and max-link. By using broadcasting, more packets are available at the relays' buffers and lower delay is achieved. Then, the work in [32] proposed low-complexity (LoCo) link selection, focusing on practical challenges of BA-ORS, including asymmetric links, outdated CSI and feedback errors. LoCo-Link exploited broadcasting and $\{R \rightarrow D\}$ prioritization to improve the performance. As outdated CSI results in severe performance degradation, its effect on BA relaying was studied in [33], showing that a coding gain can still be achieved, while [34] revealed a trade-off among CSI acquisition overhead and the quality of CSI for throughput maximization.

In order to improve the spectral efficiency of 5G networks, relaying has been jointly studied with NOMA in several works. The authors in [35] exploited users with strong channels for decoding the messages of other users. Then, by forming user pairs, the decoded messages were relayed to the users with weak channels. In [36], a single half-duplex (HD) relay assisted in the transmission to two users. Compared to OMA, NOMA offered the same diversity and, at the same time, improved spectral efficiency and fairness. Next, the authors in [37] studied a network consisting of one source and two receivers. Since communication between the receivers was allowed, one of them operated as a full-duplex (FD) relay. Then, by forming a relay broadcast channel and employing noisy network coding and dirty paper coding, the achievable rate region of the network was improved. A similar topology, where each node is equipped with multiple antennas was investigated in [38], achieving to increase the rate of a cell-edge user through joint beamforming by the source and a user with better channel conditions towards the source. Next, the study in [39] used a dedicated FD relay in a two-user NOMA network. Compared to HD relaying, FD cooperative NOMA improved the ergodic sum capacity for low-to-medium signal-to-noise ratio (SNR), even when residual self-interference remained. Additionally, NOMA allowed two sources to communicate simultaneously with two destinations, by transmitting in parallel, towards a shared relay [40]. The proposed transmission strategy showed promising performance in terms of sum ergodic capacity, both for perfect and imperfect SIC. Furthermore, NOMA was studied in ORS networks in [41], where a two-stage ORS algorithm was developed. Performance evaluation showed that improved diversity gain was achieved, compared to conventional max-min ORS. In order to reduce CSI acquisition overhead, partial relay selection using the CSI at the transmitter (CSIT) of only the $\{S \rightarrow R\}$ or $\{R \rightarrow D\}$ links was analyzed in [42]. It was shown that significant performance gains can be provided when the number of relays increased from one to two, while for more relays, the performance gains were negligible. However, NOMA with outdated CSI has not been extensively studied until now, and

only a few studies have emerged, focusing on single hop [43] or single relay [44] topologies, without providing exact expressions for the power allocation coefficient.

Aiming to enhance the performance of NOMA, recently, BA relaying was adopted. The authors of [45] proposed adaptive mode selection in a single relay network. More specifically, transmission switched between NOMA, serving simultaneously two users or single user OMA when NOMA was infeasible. From the analysis, it was concluded that the system throughput increased compared to BA relaying without NOMA and conventional relaying with NOMA. For a similar topology, the work in [46] proposed adaptive link selection for cases with full CSIT and no CSIT in the $\{R \rightarrow D\}$ links. Expressions for the outage probability were derived, and results showed that a diversity gain equal to two can be achieved for buffer sizes larger than two. More recently, two BA-ORS algorithms were proposed in [47] for a NOMA relay network. The algorithms selected the best relay, by considering BSI and $\{R \rightarrow D\}$ prioritization, in order to reduce the average delay. Nonetheless, the performance of BA relaying for NOMA with outdated CSI has not been studied, while the majority of works consider a fixed power allocation coefficient. In addition, the potential of BA-ORS in NOMA and hybrid NOMA/OMA networks has not been investigated.

B. Contributions

In this work, we study the integration of broadcasting in BA ORS networks where NOMA is employed to provide communication of one source to multiple destinations. More specifically, we exploit broadcasting to enhance the diversity of the network, since more relays are able to decode and store the source's packets, providing more options for the $\{R \rightarrow D\}$ transmission. Furthermore, contrary to existing works on NOMA, a practical issue is studied: in cases where CSI might not be timely acquired, ORS with outdated CSI is obtained; our scheme accounts for outdated CSI. Finally, outage and diversity gain analysis, for both NOMA and hybrid NOMA/OMA are conducted. In greater detail, the following contributions are given:

- The process for choosing the power allocation factor α for NOMA is given in detail, assuming the existence of two destinations with (possibly) different rate requirements. The case of users with different rate requirements served in non-orthogonal channels is envisioned in several studies on the coexistence of cellular and IoT networks; see, e.g., [9], [48], [49]. More specifically, the selection of α is based on instantaneous CSI and a practical extension is provided to ensure robustness against possibly imperfect/outdated CSI.
- In order to enhance the feasibility of our approach, α is given for cases when CSI might be outdated. More specifically, a two-hop multi-relay topology is adopted, contrary to the NOMA networks with outdated CSI in [43], [44] and the impact of outdated CSI in relay selection and its interaction with the value of α is presented.
- Two BA-ORS algorithms are developed in order to improve the outage, sum-rate and average delay performance of the network. The first algorithm, namely BA-NOMA

prioritizes $\{R \rightarrow D\}$ transmission from the relay with the maximum buffer length and, if a feasible set of links to the two destinations does not exist, $\{S \rightarrow R\}$ broadcasting is performed. The second algorithm efficiently blends NOMA and OMA, resulting in hybrid BA-NOMA/OMA. So, when a NOMA transmission in the $\{R \rightarrow D\}$ links fails, OMA transmission towards a single destination is adopted, thus avoiding a complete outage.

- A theoretical outage analysis is conducted, by modeling the different states of the network as Markov Chains. Also, the asymptotic performance of the proposed algorithms is examined and the asymptotic outage probability and diversity gain are derived.

From the performance evaluation it is shown that both algorithms improve the performance of NOMA and hybrid NOMA/OMA BA relay networks, in terms of outage probability, average sum-rate and average delay.

C. Structure

The structure of this paper is as follows. Section II presents the system model and preliminaries, necessary for the development of our results. Section III provides the process for choosing the power allocation coefficient α . Then, Section IV gives in detail the BA-NOMA relay selection algorithm, while the case of hybrid BA-NOMA/OMA is presented in Section V. Theoretical analysis of the proposed algorithms is given in the corresponding section. In Section VI, we conduct a performance evaluation and, finally, in Section VII, we conclude and discuss possible future directions.

II. SYSTEM MODEL AND PRELIMINARIES

A. System model

We consider a relay-assisted network consisting of one source, S , two destinations, D_1 and D_2 , and a cluster \mathcal{C} of K HD decode-and-forward (DF) relays $R_k \in \mathcal{C}$ ($1 \leq k \leq K$). Due to severe fading, the direct links between the source and the destinations do not exist and we assume that communication is established via the relays only. Each relay R_k is equipped with a buffer of size L , where L denotes the maximum number of data elements that can be stored from the source's transmissions. The number of packets in the buffer of relay R_k is denoted by Q_k . Each buffer is allocated equally to both destinations i.e., the same amount of data elements for D_1 and D_2 can be stored at the relays. The system model is depicted in Fig. 1.

Time is divided into "slots" of one packet duration (e.g., fixed-size packets). At any arbitrary time-slot t , the quality of the wireless channels is degraded by additive white Gaussian noise (AWGN) and frequency non-selective Rayleigh block fading, according to a complex Gaussian distribution with zero mean and variance σ_{ij}^2 for the $\{i \rightarrow j\}$ link. For simplicity, the AWGN is assumed to be normalized with zero mean and unit variance. The complex channel coefficient for the $\{i \rightarrow j\}$ link is denoted by h_{ij} , and the channel gain, $g_{ij} \triangleq |h_{ij}|^2$, is assumed to be exponentially non-identically distributed, as is the case of asymmetric topology.

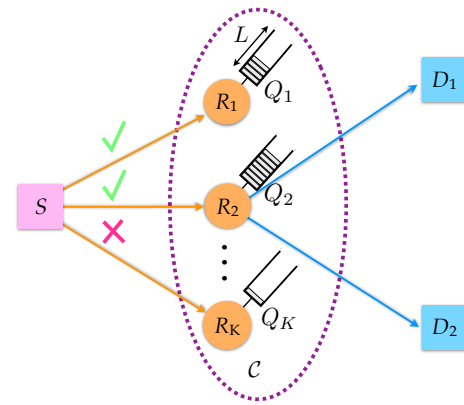


Fig. 1. BA-NOMA two-hop relay network.

The source node is assumed to be saturated (i.e., it always has data to transmit) and the required information rate, r_i , for successful reception at each destination, D_i , is fixed and may differ, depending on the application; for example, if D_1 is a cellular user and D_2 is an IoT device, the rate requirements differ and, hence, $r_1 \neq r_2$. Equivalently, a transmission from a transmitter i to its corresponding receiver j is successful if the SNR Γ_{ij} at the receiver is greater than or equal to a threshold γ_j , called the *capture ratio*. The value of γ_j depends on the modulation and coding characteristics of the application. The variance of thermal noise at relay R_k is denoted by σ_k^2 and it is assumed to be AWGN. At each time-slot, the source S or one of the relays R_k attempts to transmit a packet, using a fixed power level P_i , $i \in \{S, R_1, \dots, R_K\}$.

The retransmission process is based on an acknowledgement/negative-acknowledgement (ACK/NACK) mechanism, in which short-length error-free packets are broadcasted by the receivers over a separate narrow-band channel.

B. Transmission in the $\{S \rightarrow R\}$ link

In this topology, a time-slot is dedicated for a $\{S \rightarrow R\}$ or a $\{R \rightarrow D\}$ transmission. In the general case, each destination might demand a different rate r_j , $j \in \{1, 2\}$, and in order to avoid overflow or starvation of the buffers, the source transmits with rate $r_1 + r_2$ [46]. Therefore,

$$\Gamma_{SR_k}(P_S) \triangleq \frac{g_{SR_k} P_S}{\sigma_k^2} \geq 2^{r_1+r_2} - 1. \quad (1)$$

On the contrary, a $\{S \rightarrow R\}$ link is in outage if $\gamma_{R_k} < 2^{r_1+r_2} - 1$, and the probability of outage is given by

$$p_{out\{S \rightarrow R\}} \triangleq \mathbb{P} \left[g_{SR_k} < \frac{(2^{r_1+r_2} - 1) \sigma_k^2}{P_S} \right]. \quad (2)$$

Remark 1. It must be noted that a rate equal to $r_1 + r_2$ has been selected in order to avoid the $\{S \rightarrow R\}$ transmission being the bottleneck of the end-to-end transmission and due to the fact that no CSI is used due to the broadcasting nature of the $\{S \rightarrow R\}$ transmission. Nonetheless, as transmit SNR increases, for schemes with $\{R \rightarrow D\}$ prioritization like ours, it would be beneficial – if CSI is available – to allow the source

to transmit with adaptive rate, as relays might be selected for consecutive time-slots to transmit towards the users. Such a case is beyond the scope of this work, but could be an interesting extension for NOMA relay networks.

Let $b_{SR} \triangleq (b_{SR_1}, b_{SR_2}, \dots, b_{SR_K})$ be the binary representation of the feasible $\{S \rightarrow R\}$ links due to the fulfillment of eq. (1) (i.e., if transmission on link $\{S \rightarrow R_k\}$ is possible, then $b_{SR_k} = 1$). Moreover, let $q_{SR} \triangleq (q_{SR_1}, q_{SR_2}, \dots, q_{SR_K})$ be the binary representation of the feasible links due to the fulfillment of the queue conditions (i.e., for a $\{S \rightarrow R\}$ link the buffer in a relay is not full). By \mathcal{F}_{SR} , we denote the sets of $\{S \rightarrow R\}$ links that are feasible having a cardinality of F_{SR} .

C. Transmission in the $\{R \rightarrow D\}$ link

On the other hand, if a relay R_k is selected, the information symbols of D_1 and D_2 are superimposed and a NOMA transmission is performed. More specifically, the transmitted superimposed information symbol consisting of the information symbols x_1 and x_2 of each destination, is defined as,

$$x = \sqrt{\alpha}x_1 + \sqrt{1-\alpha}x_2 \quad (3)$$

with $\mathbb{E}[|x_1^2|] = \mathbb{E}[|x_2^2|] = 1$ and $0 \leq \alpha \leq 1$.

Then, D_1 will receive an information symbol y_1 containing the desired symbol, as well as the symbol of D_2 , i.e.,

$$y_1 = h_{R_k D_1} \sqrt{\alpha P_{R_k}} x_1 + h_{R_k D_1} \sqrt{(1-\alpha) P_{R_k}} x_2 + \eta_1; \quad (4)$$

equivalently, the received information symbol y_2 at D_2 is given by

$$y_2 = h_{R_k D_2} \sqrt{\alpha P_{R_k}} x_1 + h_{R_k D_2} \sqrt{(1-\alpha) P_{R_k}} x_2 + \eta_2, \quad (5)$$

where η_1 and η_2 denote the AWGN at each destination.

D. Outdated CSI

In practical systems, the CSI for the selection of a link might be different from the actual one, due to delays generated by the feedback mechanism. More specifically, CSI may be outdated, as the channel may vary in the period from the end of the estimation and the start of the actual transmission [50], or because it may not be fed back constantly, in order to reduce the coordination overhead [33].

From a practical perspective, wireless nodes with outdated CSI are also considered herein, and its effect on the proposed link selection algorithms is examined. In a system, where CSI feedback is delayed, the actual channel response h_{ij} conditioned on the channel response \hat{h}_{ij} that was estimated in the $\{i \rightarrow j\}$ link, during the selection period is given by [50]

$$h_{ij} | \hat{h}_{ij} \sim \mathcal{CN}(\rho_i \hat{h}_{ij}, 1 - \rho_i^2), \quad (6)$$

where $\rho_i \in [0, 1)$ denotes the correlation coefficient between h_{ij} and \hat{h}_{ij} . In relay selection networks, it has been shown that outdated CSI has a degrading effect on the diversity of the network. In [50], it was proven that a diversity order equal to one is obtained, independently of the number of the employed relays. Thus, the network performance degrades to that of single relay networks or to random relay selection, even when

$\rho_i \approx 1$ for asymptotically high SNR values. By adopting the Jakes' model [51], ρ_i is given by $\rho_i = J_0(2\pi f_{d_i} T_{D_i})$, where f_{d_i} is the Doppler frequency, T_{D_i} is the delay between link selection and the start of information transmission and $J_0(\cdot)$ is the zero-order Bessel function of the first kind.

III. CHOOSING α FOR NOMA

In this section, the selection of α determining the power allocation for each destination is presented. Since full CSIT is available at the relay, the power allocation coefficient α can be calculated in each time-slot. The value of α is chosen to ensure that the x_1 and x_2 are decoded successfully. In other words, in order to have SIC, the SINR, for at least one of the symbols, should be greater than or equal to a threshold at both destinations D_1 and D_2 . This process must be performed by each relay in order to define the set of relays that can perform NOMA transmissions in the $\{R \rightarrow D\}$ links.

For example, for decoding x_2 at both destinations,

$$\Gamma_{R_k D_j}(P_{R_k}) = \frac{(1-\alpha)P_{R_k}g_{R_k D_j}}{\alpha P_{R_k}g_{R_k D_j} + \sigma_{D_j}^2} \geq \gamma_j, \quad j \in \{1, 2\}. \quad (7)$$

Note that $\gamma_j \equiv 2^{r_j} - 1$. Then, once x_2 is successfully decoded, x_1 can be decoded interference-free at destination D_1 according to

$$\Gamma_{R_k D_1}(P_{R_k}) = \frac{\alpha P_{R_k}g_{R_k D_1}}{\sigma_{D_1}^2} \geq \gamma_1. \quad (8)$$

For this example, from inequalities (7) and (8) we get that

$$\alpha \leq \frac{P_{R_k}g_{R_k D_2} - \gamma_2 \sigma_{D_2}^2}{P_{R_k}g_{R_k D_2}(1 + \gamma_2)}, \quad (9a)$$

$$\alpha \leq \frac{P_{R_k}g_{R_k D_1} - \gamma_2 \sigma_{D_1}^2}{P_{R_k}g_{R_k D_1}(1 + \gamma_1)}, \quad (9b)$$

$$\alpha \geq \frac{\gamma_1 \sigma_{D_1}^2}{P_{R_k}g_{R_k D_1}}. \quad (9c)$$

Then, α can take values in the range

$$\alpha_{\min} \leq \alpha \leq \max\{0, \min\{1, \alpha_{\max}\}\}, \quad (10)$$

where

$$\alpha_{\min} \triangleq \frac{\gamma_1 \sigma_{D_1}^2}{P_{R_k}g_{R_k D_1}},$$

$$\alpha_{\max} \triangleq \min\left\{\frac{P_{R_k}g_{R_k D_1} - \gamma_2 \sigma_{D_1}^2}{P_{R_k}g_{R_k D_1}(1 + \gamma_2)}, \frac{P_{R_k}g_{R_k D_2} - \gamma_2 \sigma_{D_2}^2}{P_{R_k}g_{R_k D_2}(1 + \gamma_2)}\right\}.$$

Similarly, for decoding x_1 first at both destinations, we have

$$\alpha_{\min} \triangleq \max\left\{\frac{\gamma_1(P_{R_k}g_{R_k D_1} + \sigma_{D_1}^2)}{P_{R_k}g_{R_k D_1}(1 + \gamma_1)}, \frac{\gamma_1(P_{R_k}g_{R_k D_2} + \sigma_{D_2}^2)}{P_{R_k}g_{R_k D_2}(1 + \gamma_1)}\right\},$$

$$\alpha_{\max} \triangleq \frac{P_{R_k}g_{R_k D_2} - \gamma_2 \sigma_{D_2}^2}{P_{R_k}g_{R_k D_2}}.$$

The outage probability of NOMA is equal to

$$p_{\text{out}\{R \rightarrow D\}} = \mathbb{P}[\alpha_{\min} > \min\{1, \alpha_{\max}\}]. \quad (11)$$

Let $b_{RD} \triangleq (b_{R_1 D}, b_{R_2 D}, \dots, b_{R_K D})$ be the binary representation of the feasible $\{R \rightarrow D\}$ links due to the fulfillment of eqs. (7), (8) (i.e., if a NOMA transmission on the set of

links $\{R_k \rightarrow D_1\}$, $\{R_k \rightarrow D_2\}$ is possible, then $b_{R_k D} = 1$. Similarly, let $q_{RD} \triangleq (q_{R_1 D}, q_{R_2 D}, \dots, q_{R_K D})$ be the binary representation of the feasible links due to the fulfillment of the queue conditions. By \mathcal{F}_{RD} , we denote the sets $\{R \rightarrow D\}$ links that are feasible, having a cardinality of F_{RD} .

For practical implementation, the values α_{\min} and α_{\max} for each case can be computed at each R_k . Assuming that each R_k has the CSIT for D_1, D_2 the range of values of α given by eq. (10) is accurately calculated. Nonetheless, in practical systems channel estimation errors might occur. In this case, it is desired to choose α in a way that robustness against CSIT estimation errors is provided. Towards this end, we propose to choose α , so as

$$\alpha = \frac{\alpha_{\min} + \max\{0, \min\{1, \alpha_{\max}\}\}}{2}. \quad (12)$$

In addition, the selected relay must inform the two destinations on the decoding strategy that they must adopt. This is achieved by adding an extra bit to the packet's header. If the bit value is "0", D_1 will perform SIC, thus decoding its packet interference-free, while D_2 will decode its packet by considering the signal intended for D_1 , as interference. The reverse strategies are adopted when the bit value is "1". It must be noted that each R_k examines sequentially the possible decoding strategies and when the first strategy that fulfills the rate requirements is found, its decision is transmitted by the selected relay to D_1 and D_2 .

IV. BUFFER-AIDED NOMA

A. Description of the algorithm

The first selection algorithm aims at improving the performance of BA-NOMA relay networks. More specifically, the BA-NOMA relay selection algorithm aims to activate a relay to simultaneously serve D_1 and D_2 through NOMA. So, when multiple relays are available, in the $\{S \rightarrow R\}$ link, S broadcasts the combined signals of D_1 and D_2 , towards the K available relays, with a rate $r = r_1 + r_2$. Through broadcasting, more packets are available at the relays' buffers and more importantly, CSIT is not required at the source, thus significantly reducing the implementation complexity.

Then, in the $\{R \rightarrow D\}$ link, each relay determines the power allocation factor α for NOMA transmission, according to (12). Also, in order to reduce the average delay of NOMA, BA-NOMA prioritizes the selection of the relay with the maximum number of packets in its queue. Thus, if a feasible set of links to D_1 and D_2 exists, relay $R_i^* \in \mathcal{F}_{RD}$, having the maximum number of packets in its queue among the set \mathcal{F}_{RD} is activated for transmission. If the set \mathcal{F}_{RD} is empty, source broadcasting is performed. The BA-NOMA link selection algorithm for a single time-slot is summarized in Algorithm 1.

By exploiting broadcasting and selecting the relay with the maximum number of packets for $\{R \rightarrow D\}$ transmission, diversity can be maintained. Nonetheless, NOMA transmissions to D_1 and D_2 often require high SINR, especially for high data rate requirements, which leads to outages for both links, even though single link transmissions to one of the destinations would be feasible. Thus, it is necessary to devise a hybrid algorithm combining NOMA and OMA, efficiently switching

Algorithm 1 The BA-NOMA algorithm

```

1: input  $\mathcal{F}_{RD}, \alpha$ 
2: if  $\mathcal{F}_{RD} = \emptyset$  then
3:   The source broadcasts packets for  $D_1$  and  $D_2$ .
4:    $Q_j \leftarrow Q_j + 2, \quad \forall j \in \mathcal{F}_{SR}$ 
5: else
6:    $i' = \arg \max_{i \in \mathcal{F}_{RD}} Q_i$  ( $\{R \rightarrow D\}$  link)
7:   if more than one relays have the same maximum queue
       length then
8:      $i^*$  is chosen randomly among the set of relays in  $i'$ .
9:   else
10:     $i^* = i'$ .
11:   end if
12:    $Q_{i^*} \leftarrow Q_{i^*} - 2$ 
13: end if
14: Output Link  $\{R_{i^*} \rightarrow D\}$  is activated for NOMA transmission or the set of links in  $\mathcal{F}_{SR}$  receive a combined packet with rate  $r_1 + r_2$  from the source, if  $\mathcal{F}_{SR} \neq \emptyset$ .

```

between the two multiple-access schemes, in order to avoid a complete network outage.

B. Theoretical analysis of the BA-NOMA algorithm

In this algorithm, we observe that the data for destinations D_1 and D_2 that the source transmits to the relays remain together in the buffer and they are eventually transmitted simultaneously to the destinations. Hence, we can treat these data as one packet in the network. As a result, the theoretical analysis of this work is similar, *mutatis mutandis*, to those of [31], [32], and can be cast as a discrete time Markov chain (DTMC) in which each state represents a possible state of the buffers. The main difference in the BA-NOMA algorithm is that the values of the transition probabilities of the DTMC differ from the source to the relays, since the source transmits packets for two destinations, and the same holds for the transition probabilities from the relays to the destinations, as packets are transmitted only when both links are not in outage.

States of the DTMC. As aforementioned, the states of the DTMC represent all the possible states of the buffers. The state of the DTMC can be represented by $S_r \in \mathcal{S}$, where \mathcal{S} is the set of all available states, $r \in \mathbb{N}$, $1 \leq r \leq |\mathcal{S}|$, and $|\mathcal{S}|$ represents the cardinality of all the possible combinations of the buffer states.

Construction of the state transition matrix of the DTMC. Let $(X_t)_{t \geq 0}$ denote the discrete-time Markov random process capturing the evolution of the system. Also, let $\mathbf{A} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$ denote the state transition matrix of the DTMC, in which the entry $\mathbf{A}_{i,j} = \mathbb{P}(S_j \rightarrow S_i) \triangleq \mathbb{P}(X_{t+1} = S_i | X_t = S_j)$ is the transition probability to move from state S_j at time t to state S_i at time $(t + 1)$. At each time slot, the buffer state and, hence, the state of the system, can be modified as follows: (a) the number of elements of one or more buffers can be decreased by one, if a relay node is selected for transmission and the transmission is successful, and also other relays, that have the same packet in their buffer, drop it; (b) the number of elements of one or more buffers can be increased by one, if the

source node is selected for broadcasting and the transmission to at least one relay is successful; (c) the state of all buffers remains unchanged when there is an outage event (i.e., all the $\{S \rightarrow R\}$ and $\{R \rightarrow D\}$ links are in outage).

Properties of the DTMC. The DTMC depicts finite-size buffers and, hence, it is stationary, irreducible and aperiodic (SIA) [21], i.e., a steady state (also known as the distribution of the DTMC) λ exists such that $\mathbf{A}\lambda = \lambda$ and $\mathbf{b}^T \lambda = 1$, where $\mathbf{b} = [1 \ 1 \dots 1]^T$. Steady-state λ of the column stochastic matrix \mathbf{A} of the MC that models the states of the network system is given by [21, Lemma 1] $\lambda = (\mathbf{A} - \mathbf{I} + \mathbf{B})^{-1} \mathbf{b}$, where \mathbf{I} is the identity matrix of appropriate dimensions and \mathbf{B} is a square matrix of appropriate dimensions with $\mathbf{B}_{ij} = 1$ for all i, j . The analytical expressions for the outage probability, average throughput and average packet delay, for our proposed algorithm, and for *any policy defined for finite-length buffer-aided relay selection* based on the construction of the DTMC can be found in several papers; see, for example, [32, Equations (7), (8) and (11)].

C. Asymptotic performance of the BA-NOMA algorithm

When the *average* channel SNR, denoted herein by ϕ , goes to infinity (i.e., $\phi \rightarrow \infty$), the probability of a $\{S \rightarrow R\}$ or a $\{R \rightarrow D\}$ link being in outage goes to zero, provided the buffer is not full or empty, respectively. The prioritization of the $\{R \rightarrow D\}$ links at high SNR makes the scheme deterministic, as it is the case of the LoCo – Link scheme in [32]. The main reason is the following: since $\{R \rightarrow D\}$ links are prioritized for transmission, there will be a time, say t , at which all buffers will eventually become empty. At the next time-slot (i.e., $t+1$), since there are no packets in the buffers, the source broadcasts *the same* packet to all the relays. At the next time-slot (i.e., $t+2$), priority is again given to the $\{R \rightarrow D\}$ links and, hence, once a relay transmits the packet in its buffer, all other relays discard their packet and the relays become empty again. Since this procedure is repeated indefinitely, we can cast this as a two-state DTMC: one, denoted by S_e , corresponds to the case in which there are no packets in the buffers and the other, denoted by S_f , corresponds to the case in which there is one identical packet in all relay buffers. For simplicity of exposition of the analysis, we assume that all the channel gains are independent and identically distributed (i.i.d.) with rate parameter equal to 1, i.e., $g_{SR_k}, g_{R_k D_j} \sim \text{Exp}(1) \ \forall k \in \{1, \dots, K\}, j \in \{1, 2\}$.

Derivation of the Asymptotic Outage Probability. Since the scheme switches between two possible states, S_e and S_f , then as $\phi \rightarrow \infty$, the probability of the system being in outage is given by

$$p_{\text{out}}^{\infty} \triangleq \lim_{\phi \rightarrow \infty} p_{\text{out}} = \mathbb{P}(S_e) p_{\text{out}}^{S_e} + \mathbb{P}(S_f) p_{\text{out}}^{S_f}, \quad (13)$$

where $\mathbb{P}(S_e)$ and $\mathbb{P}(S_f)$ are the probabilities that the state of the buffers are in S_e and S_f , respectively, and $p_{\text{out}}^{S_e}$ and $p_{\text{out}}^{S_f}$ are the outage probabilities to the corresponding states. Since, states S_e and S_f are switching continuously, the probability of being at one of the states at any given time is equal to $1/2$,

i.e., $\mathbb{P}(S_e) = \mathbb{P}(S_f) = \frac{1}{2}$. When the state is S_e , there are K available $\{S \rightarrow R\}$ links only. Therefore,

$$p_{\text{out}}^{S_e} = \lim_{\phi \rightarrow \infty} \left(1 - e^{-\frac{\gamma_0}{\phi}}\right)^K \stackrel{(a)}{\approx} \lim_{\phi \rightarrow \infty} \left(\frac{\gamma_0}{\phi}\right)^K, \quad (14)$$

where $\gamma_0 = 2^{r_1+r_2} - 1$ and (a) stems from the fact that for very small x , $e^x \approx 1 + x$. When the state is S_f , there are K available *pairs* of $\{R \rightarrow D\}$ links and n ($n = 0$ for $L = 1$ and $n = K$ for $L \geq 2$) available $\{S \rightarrow R\}$ links (that are not used unless the $\{R \rightarrow D\}$ links are in outage). Suppose that x_2 is decoded first. Then, in the asymptotic case, since when $\phi \rightarrow \infty$ it is equivalent to $P_{R_k} \rightarrow \infty$, (7) becomes

$$\begin{aligned} \Gamma_{R_k D_j}^{\infty}(P_{R_k}) &\triangleq \lim_{P_{R_k} \rightarrow \infty} \frac{(1 - \alpha) P_{R_k} g_{R_k D_j}}{\alpha P_{R_k} g_{R_k D_j} + \sigma_{D_j}^2} \\ &\equiv \lim_{\phi \rightarrow \infty} \frac{(1 - \alpha) \phi g_{R_k D_j}}{\alpha \phi g_{R_k D_j} + 1} \\ &= \frac{(1 - \alpha)}{\alpha} \geq \gamma_j, \quad j \in \{1, 2\}. \end{aligned} \quad (15)$$

The inequality emerging from (15) can be written as

$$\alpha \leq \frac{1}{1 + \max\{\gamma_1, \gamma_2\}}, \quad (16)$$

and it is justified by inequalities (9a) and (9b) for $P_{R_k} \rightarrow \infty$. Condition (16) can be guaranteed from our choice of parameter α , since it can be easily shown that if inequalities (9a) and (9b) hold, then (16) holds.

In the following proposition, we provide the probability for each of the $\{R \rightarrow D\}$ links to be in outage.

Proposition 1. *The outage probability at D_2 in which the data to D_1 (i.e., x_1) are not required to be decoded first, is given by*

$$p_{\text{out}, D_2} = 1 - e^{-\frac{\gamma_2}{[1 - \alpha(1 + \gamma_2)]\phi}}. \quad (17)$$

The probability of D_1 being in outage in BA-NOMA algorithm, given that it decodes x_2 first, is given by

$$p_{\text{out}, D_1} = \left(1 - e^{-\frac{\gamma_1}{\min\{[1 - \alpha(1 + \gamma_1)], \alpha\}\phi}}\right). \quad (18)$$

Proof. See Appendix A. □

Note that, if one of the destinations is not able to successfully receive its data, the whole packet (with x_1 and x_2) is dropped. In the following proposition, we provide the probability of the $\{R \rightarrow D\}$ transmission to be in outage.

Proposition 2. *The probability of any of the $\{R \rightarrow D\}$ links being in outage is given by*

$$p_{\text{out}, D} = 1 - e^{-\left(\frac{\gamma_1}{\alpha_1 \phi} + \frac{\gamma_2}{\alpha_2 \phi}\right)}, \quad (19)$$

where $\alpha_1 \triangleq \min\{[1 - \alpha(1 + \gamma_1)], \alpha\}$ and $\alpha_2 \triangleq 1 - \alpha(1 + \gamma_2)$.

Proof. See Appendix B. □

From Proposition 2,

$$p_{\text{out}}^{S_f} = \lim_{\phi \rightarrow \infty} \left(1 - e^{-\left(\frac{\gamma_1}{\alpha_1 \phi} + \frac{\gamma_2}{\alpha_2 \phi}\right)} \right)^K \left(1 - e^{-\frac{\gamma_0}{\phi}} \right)^n$$

$$\stackrel{(a)}{\approx} \lim_{\phi \rightarrow \infty} \left(\frac{\gamma_1}{\alpha_1 \phi} + \frac{\gamma_2}{\alpha_2 \phi} \right)^K \left(\frac{\gamma_0}{\phi} \right)^n \quad (20)$$

$$= \lim_{\phi \rightarrow \infty} \left(\frac{\gamma_{1,2}}{\phi} \right)^K \left(\frac{\gamma_0}{\phi} \right)^n, \quad (21)$$

where (a), as before, stems from the fact that for very small x , $e^x \approx 1 + x$ and $\gamma_{1,2} \triangleq \frac{\gamma_1}{\alpha_1} + \frac{\gamma_2}{\alpha_2}$. Hence, (13) becomes

$$p_{\text{out}}^\infty = \frac{1}{2} \left(p_{\text{out}}^{S_e} + p_{\text{out}}^{S_f} \right)$$

$$= \frac{1}{2} \lim_{\phi \rightarrow \infty} \left[\left(\frac{\gamma_0}{\phi} \right)^K + \left(\frac{\gamma_{1,2}}{\phi} \right)^K \left(\frac{\gamma_0}{\phi} \right)^n \right]. \quad (22)$$

Derivation of the diversity gain. The diversity gain of the proposed relay selection scheme can be computed as follows:

$$d = - \lim_{\phi \rightarrow \infty} \frac{\log p_{\text{out}}}{\log \phi}$$

$$= - \lim_{\phi \rightarrow \infty} \frac{\log \left[\frac{1}{2} \left(\left(\frac{\gamma_0}{\phi} \right)^K + \left(\frac{\gamma_{1,2}}{\phi} \right)^K \left(\frac{\gamma_0}{\phi} \right)^n \right) \right]}{\log \phi}. \quad (23)$$

For buffer size $L = 1$, then $n = 0$, and (23) becomes

$$d_{(L=1)} = - \lim_{\phi \rightarrow \infty} \frac{\log \left[\frac{1}{2} \left(\left(\frac{\gamma_0}{\phi} \right)^K + \left(\frac{\gamma_{1,2}}{\phi} \right)^K \right) \right]}{\log \phi}$$

$$= - \lim_{\phi \rightarrow \infty} \frac{\log \left(\frac{\gamma_0}{\phi} \right)^K + \log \left[\left(1 + \left(\frac{\gamma_{1,2}}{\gamma_0} \right)^K \right) \right]}{\log \phi}$$

$$= K - \lim_{\phi \rightarrow \infty} \frac{\log \left[\left(1 + \left(\frac{\gamma_{1,2}}{\gamma_0} \right)^K \right) \right]}{\log \phi} = K.$$

For buffer size $L \geq 2$, then $n = K$, and therefore (23) becomes

$$d_{(L \geq 2)} = - \lim_{\phi \rightarrow \infty} \frac{\log \left[\frac{1}{2} \left(\left(\frac{\gamma_0}{\phi} \right)^K + \left(\frac{\gamma_{1,2}}{\phi} \right)^K \left(\frac{\gamma_0}{\phi} \right)^K \right) \right]}{\log \phi}$$

$$= - \lim_{\phi \rightarrow \infty} \frac{\log \left(\frac{\gamma_0}{\phi} \right)^K + \log \left[\left(1 + \left(\frac{\gamma_{1,2}}{\phi} \right)^K \right) \right]}{\log \phi}$$

$$\stackrel{(a)}{=} K - \lim_{\phi \rightarrow \infty} \frac{\left(\frac{\gamma_{1,2}}{\phi} \right)^K}{\log \phi} = K,$$

where (a) stems from the fact that for small x , $\log(1+x) \approx x$.

Remark 2. The diversity gain for $L \geq 2$ converges much faster to K since the term that goes to zero decays with exponential rate, i.e., $O(\phi^K \log \phi)$, whereas for $L = 1$ the diversity gain converges with logarithmic rate only, i.e., $O(\log \phi)$.

V. HYBRID BUFFER-AIDED NOMA/OMA

The second selection algorithm, namely hybrid BA-NOMA/OMA combines the two multiple-access schemes, maintaining the sum-rate of the network. More specifically, BA-NOMA/OMA allows the source to broadcast a packet with rate $r_1 + r_2$ towards the K relays. Then, the relays derive α according to (12). However, it is possible that a NOMA transmission might not be successful, especially for low available transmit SNR and high rate requirements. In this case, the network should not experience a complete outage, since, at least, one destination D_j might be able to receive its packet from the relays belonging in \mathcal{F}_{RD_j} , $j \in \{1, 2\}$. Thus, set \mathcal{F}_{RD_j} includes links that can achieve an OMA transmission towards D_j according to

$$\Gamma_{R_k D_j}(P_{R_k}) = \frac{P_{R_k} g_{R_k D_j}}{\sigma_{D_j}^2} \geq \gamma_j. \quad (24)$$

In order to maintain user fairness, at odd time-slots, first, D_1 is examined on whether or not it can receive its packet from $R_k \in \mathcal{F}_{RD_1}$, otherwise D_2 is examined on whether or not it can receive its packet from $R_k \in \mathcal{F}_{RD_2}$. The order with which, each destination is examined for successful reception is reversed at even time-slots. The hybrid BA-NOMA/OMA link selection algorithm for a single time-slot t is summarized in Algorithm 2.

For clarity, the case of having relays with equal queue lengths in the $\{R \rightarrow D\}$ links is omitted in the description of Algorithm 2. Nonetheless, in such cases, the relay is selected in a random manner, as described in lines 7 to 11 of Algorithm 1. BA-NOMA/OMA constitutes an efficient selection algorithm combining the merits of both multiple-access schemes. It is obvious that switching among NOMA and OMA allows the network to avoid outages and sustain its performance independently of the available transmit SNR. In the performance evaluation at Section VI, BA-NOMA/OMA is evaluated for various scenarios with different data rate requirements, showing promising performance gains compared to either NOMA or OMA algorithms.

Regarding the practical implementation of BA-NOMA/OMA, a 2-bit signaling scheme is devised. More specifically, the first bit of the signaling message denotes whether or not, a NOMA transmission is performed. So, when the first bit is “1” NOMA is feasible, while a “0” value triggers an OMA transmission. Then, the second bit informs each destination to perform a specific action. If NOMA is performed and the second bit is “0”, D_1 performs SIC, while a value equal to “1” leads D_2 to carry out SIC. However, for OMA, a value equal to “0” for the second bit flags that D_1 is scheduled, and a value equal to “1” for the second bit denotes that D_2 is scheduled. Table I includes all the possible cases for the 2-bit signaling, necessary for the operation of BA-NOMA/OMA.

A. Theoretical analysis of the hybrid BA-NOMA/OMA algorithm

In the hybrid BA-NOMA/OMA algorithm, if NOMA is not feasible, the algorithm adopts OMA in which one of the two

Algorithm 2 The hybrid BA–NOMA/OMA algorithm

```

1: input  $\mathcal{F}_{RD}, \mathcal{F}_{RD_1}, \mathcal{F}_{RD_2}, \alpha, t$ 
2: if  $\mathcal{F}_{RD} = \emptyset$  then
3:   if  $\text{mod}(t, 2) = 1$  then
4:     if  $\mathcal{F}_{RD_1} \neq \emptyset$  then
5:        $i' = \arg \max_{i \in \mathcal{F}_{RD_1}} Q_i$    ( $\{R \rightarrow D\}$  link)
6:        $Q_{i'} \leftarrow Q_{i'} - 1$ 
7:     else
8:       if  $\mathcal{F}_{RD_2} \neq \emptyset$  then
9:          $i' = \arg \max_{i \in \mathcal{F}_{RD_2}} Q_i$    ( $\{R \rightarrow D\}$  link)
10:         $Q_{i'} \leftarrow Q_{i'} - 1$ 
11:      else
12:        The source broadcasts its value.
13:         $Q_j \leftarrow Q_j + 2, \quad \forall j \in \mathcal{F}_{SR}$ 
14:      end if
15:    end if
16:  else
17:    if  $\mathcal{F}_{RD_2} \neq \emptyset$  then
18:       $i' = \arg \max_{i \in \mathcal{F}_{RD_2}} Q_i$    ( $\{R \rightarrow D\}$  link)
19:       $Q_{i'} \leftarrow Q_{i'} - 1$ 
20:    else
21:      if  $\mathcal{F}_{RD_1} \neq \emptyset$  then
22:         $i' = \arg \max_{i \in \mathcal{F}_{RD_1}} Q_i$    ( $\{R \rightarrow D\}$  link)
23:         $Q_{i'} \leftarrow Q_{i'} - 1$ 
24:      else
25:        The source broadcasts its value.
26:         $Q_j \leftarrow Q_j + 2, \quad \forall j \in \mathcal{F}_{SR}$ 
27:      end if
28:    end if
29:  end if
30: else
31:    $i' = \arg \max_{i \in \mathcal{F}_{RD}} Q_i$    ( $\{R \rightarrow D\}$  link)
32:    $Q_{i'} \leftarrow Q_{i'} - 2$ 
33: end if
34: Output Link  $\{R_{i'} \rightarrow D\}$  is activated for NOMA transmission or link  $\{R_{i'} \rightarrow D\}$  is activated for OMA transmission towards  $D_1$  or  $D_2$ , or the set of links in  $\mathcal{F}_{SR}$  receive a packet from the source, if  $\mathcal{F}_{SR} \neq \emptyset$ .

```

TABLE I
THE 2-BIT SIGNALING SCHEME OF BA–NOMA/OMA

| Bit sequence | Multiple access scheme | Strategy |
|--------------|------------------------|--------------------|
| 00 | OMA | D_1 is scheduled |
| 01 | OMA | D_2 is scheduled |
| 10 | NOMA | D_1 performs SIC |
| 11 | NOMA | D_2 performs SIC |

destinations receives its data. As mentioned in the description of the algorithm, a fairness approach is used in which odd time slots are used for transmitting to D_1 , if NOMA is not possible, and even time slots for transmitting to D_2 . Note however that if in an odd/even time slot in which NOMA is not feasible D_1/D_2 is in outage, then the slot is used for transmitting to D_2/D_1 , if the link is not in outage. As a result, we cannot treat these data as one packet in the network any more, as we

did for the BA–NOMA algorithm in Section IV, because at a relay, the data might get split because of OMA and only the data for one of the destinations is transmitted. As a result, the theoretical analysis of this part differs, even though it can be cast as a DTMC in which each state represents a possible state of the buffers.

States of the DTMC. The state of the DTMC is given by the state of the buffers, as in Section IV-B.

Construction of the state transition matrix of the DTMC. At each time slot, the buffer state and, hence, the state of the system, can be modified as follows: (a) the number of elements of one or more buffers can be decreased by two, if a relay node is selected for transmission and the transmission is successful in NOMA, and also the other relays, that have the same packet in their buffer, drop it; (b) the number of elements of one or more buffers can be decreased by one only, if a relay node is selected for transmission and the transmission is successful in OMA (i.e., NOMA is not feasible, but OMA is), (c) the number of elements of one or more buffers can be increased by two, if the source node is selected for broadcasting and the transmission to at least one relay is successful; (d) the state of all buffers remains unchanged when there is an outage event (i.e., all the $\{S \rightarrow R\}$ and all the $\{R \rightarrow D\}$ links are in outage). Unlike BA–NOMA, from the relays you may have either one or two packets transmitted. We also assume, for simplicity, that the packet sizes for each destination are of the same size, but in principle they can be of different sizes as well.

Properties of the DTMC. Despite the differences in the construction of the state transition matrix of the DTMC, once the DTMC is constructed, the analytical expressions for the outage probability, average throughput and average packet delay do not change. Constructing examples for our DTMC is trivial and beyond the scope of this paper.

B. Asymptotic performance of the hybrid BA–NOMA/OMA algorithm

The advantage of the hybrid BA–NOMA/OMA algorithm is that when the BA–NOMA algorithm is in outage because, either of the destinations has failed to decode successfully, the algorithm switches to OMA, thus allowing one of the destinations to receive data. Hence, the outage probability is given by

$$\begin{aligned}
 p_{\text{out}}^{S_f} &= \lim_{\phi \rightarrow \infty} \left[\left(1 - e^{-\frac{\gamma_1}{\alpha_1 \phi}} \right) \left(1 - e^{-\frac{\gamma_2}{\alpha_2 \phi}} \right) \right]^K \left(1 - e^{-\frac{\gamma_0}{\phi}} \right)^n \\
 &\stackrel{(a)}{\approx} \lim_{\phi \rightarrow \infty} \left(\frac{\gamma_1}{\alpha_1 \phi} \right)^K \left(\frac{\gamma_2}{\alpha_2 \phi} \right)^K \left(\frac{\gamma_0}{\phi} \right)^n, \quad (25)
 \end{aligned}$$

where (a) stems from the fact that for very small x , $e^x \approx 1+x$. Hence, (13) now becomes

$$\begin{aligned}
 p_{\text{out}}^\infty &= \frac{1}{2} \left(p_{\text{out}}^{S_e} + p_{\text{out}}^{S_f} \right) \\
 &= \frac{1}{2} \lim_{\phi \rightarrow \infty} \left[\left(\frac{\gamma_0}{\phi} \right)^K + \left(\frac{\gamma_1}{\alpha_1 \phi} \right)^K \left(\frac{\gamma_2}{\alpha_2 \phi} \right)^K \left(\frac{\gamma_0}{\phi} \right)^n \right].
 \end{aligned}$$

Derivation of the diversity gain. The diversity gain of the proposed relay selection scheme can be computed as follows:

$$d = - \lim_{\phi \rightarrow \infty} \frac{\log p_{\text{out}}}{\log \phi} = - \lim_{\phi \rightarrow \infty} \frac{\log \left[\frac{1}{2} \left(\left(\frac{\gamma_0}{\phi} \right)^K + \left(\frac{\gamma_1}{\alpha_1 \phi} \right)^K \left(\frac{\gamma_2}{\alpha_2 \phi} \right)^K \left(\frac{\gamma_0}{\phi} \right)^n \right) \right]}{\log \phi}. \quad (26)$$

For buffer size $L = 1$, then $n = 0$, and (26) becomes

$$d_{(L=1)} = - \lim_{\phi \rightarrow \infty} \frac{\log \left[\frac{1}{2} \left(\left(\frac{\gamma_0}{\phi} \right)^K + \left(\frac{\gamma_1}{\alpha_1 \phi} \right)^K \left(\frac{\gamma_2}{\alpha_2 \phi} \right)^K \right) \right]}{\log \phi} = - \lim_{\phi \rightarrow \infty} \frac{\log \left(\frac{\gamma_0}{\phi} \right)^K + \log \left[\left(1 + \left(\frac{\gamma_1 \gamma_2}{\alpha_1 \alpha_2 \gamma_0 \phi} \right)^K \right) \right]}{\log \phi} \stackrel{(a)}{=} K - \lim_{\phi \rightarrow \infty} \frac{\left(\frac{\gamma_1 \gamma_2}{\alpha_1 \alpha_2 \gamma_0 \phi} \right)^K}{\log \phi} = K.$$

where (a) stems from the fact that for small x , $\log(1+x) \approx x$. For buffer size $L \geq 2$, then $n = K$, and therefore (26) becomes

$$d_{(L \geq 2)} = - \lim_{\phi \rightarrow \infty} \frac{\log \left[\frac{1}{2} \left(\left(\frac{\gamma_0}{\phi} \right)^K + \left(\frac{\gamma_1}{\alpha_1 \phi} \right)^K \left(\frac{\gamma_2}{\alpha_2 \phi} \right)^K \left(\frac{\gamma_0}{\phi} \right)^K \right) \right]}{\log \phi} = - \lim_{\phi \rightarrow \infty} \frac{\log \left(\frac{\gamma_0}{\phi} \right)^K + \log \left[\left(1 + \left(\frac{\gamma_1}{\alpha_1 \phi} \right)^K \left(\frac{\gamma_2}{\alpha_2 \phi} \right)^K \right) \right]}{\log \phi} \stackrel{(a)}{=} K - \lim_{\phi \rightarrow \infty} \frac{\left(\frac{\gamma_1}{\alpha_1 \phi} \right)^K \left(\frac{\gamma_2}{\alpha_2 \phi} \right)^K}{\log \phi} = K,$$

where (a) stems from the fact that for small x , $\log(1+x) \approx x$.

Remark 3. The diversity gain for $L \geq 2$ converges much faster to K since the term that goes to zero decays with exponential rate, i.e., $O(\phi^{2K} \log \phi)$, whereas for $L = 1$ the diversity gain converges with rate as that for $L \geq 2$ for the BA-NOMA, i.e., $O(\phi^K \log \phi)$.

VI. NUMERICAL RESULTS

Next, performance evaluation results are given for BA-NOMA and BA-NOMA/OMA. The topology consists of a two-hop network with two users, a varying number of relays K and a buffer size $L = 4$ packets. To provide additional insight on the performance of the two selection algorithms, as well as a more realistic scenario, the channels of the two users are considered to be independent non identically distributed (i.n.i.d). Regarding the required rate by each user, two cases are evaluated, i.e. equal rate and different rate. The case of equal rate can be mapped to two cellular users requiring the same service from a base station, while the case of a high-rate user and a low-rate user might correspond to a cellular user and an IoT device sharing the same wireless

channel. The considered OMA schemes include the LoCo-Link algorithm of [32], the $\{R \rightarrow D\}$ prioritization algorithm of [24] and the delay-aware algorithm of [26]. Also, apart from the two proposed algorithms, NOMA versions of the algorithms of [24] and [26] are evaluated, as well as the delay-aware NOMA algorithm of [47]. For OMA, a fixed scheduler allocates the wireless channel to the first user in the odd time-slots and to the second user in the even time-slots, while outages occur when at two consecutive time-slots there is no relay that can perform a successful transmission/reception. Also, for NOMA, outages occur if at a given time-slot there is no relay that can perform a successful transmission/reception. For a fair comparison, the hybrid BA-NOMA/OMA algorithm is not included in the outage results, as it cannot be directly compared to NOMA and OMA algorithms, similarly to [46]. Outages occur according to eq. (2) and eq. (11).

A. Equal rate requirements

In the first scenario, the two users require equal rates $r_1 = r_2 = r = 3$ BPCU and in this case, the channel asymmetry is defined as $\sigma_{R_k D_1}^2 = 4\sigma_{R_k D_2}^2$. For a fair comparison, in the OMA algorithms, the transmission in each hop should satisfy twice the rate requirement, since the overall communication demands two times the time-slots of NOMA. For the hybrid NOMA/OMA, when NOMA fails, the OMA transmission in the RD link selects one user and transmits with rate r .

1) *Outage probability:* The outage probability for NOMA and OMA transmissions is depicted in Fig. 2 for a network where $K = 4$ and $L = 4$. As each user has the same rate requirements, possible gains for NOMA derive from channel asymmetries. As a result, it can be seen that due to D_1 having stronger channels towards the relays than D_2 , NOMA algorithms outperform their OMA versions. For both multiple-access cases, for low transmit SNR, the network is in outage as rate requirements are increased, while after 15 dB the network's performance improves. Focusing on the NOMA algorithms, equal diversity performance is observed, since all of them perform $\{R \rightarrow D\}$ prioritization. The same observation applies for the OMA algorithms, where at each time-slot, if an $\{R \rightarrow D\}$ link is not in outage and there are packets in the relays' buffers, it is scheduled for transmission.

2) *Average sum-rate:* Next, average sum-rate results are given in Fig. 3, for NOMA, OMA and hybrid NOMA/OMA algorithms when $K = 4$ and $L = 4$. It can be seen that the hybrid NOMA/OMA provides superior performance for all the SNR range. It is obvious that switching to OMA transmission in the $\{R \rightarrow D\}$ link, when NOMA cannot be performed, allows at least one user to be served at its required rate. As a result, the additional signaling complexity of hybrid NOMA/OMA is justified in order to achieve better performance. Then, for SNR values between 10 dB and 15 dB, OMA outperforms NOMA, as the interference between the users' signals cannot guarantee that both signals are successfully decoded at each destination. Nonetheless, when increased transmit SNR is available, due to the robust selection of α , NOMA offers improved performance compared to OMA, as the users' power allocation is based on channel

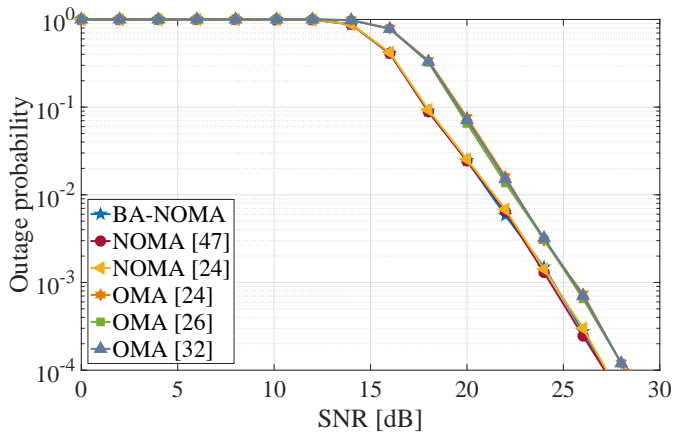


Fig. 2. Outage probability for $K = 4$, $L = 4$ and various algorithms.

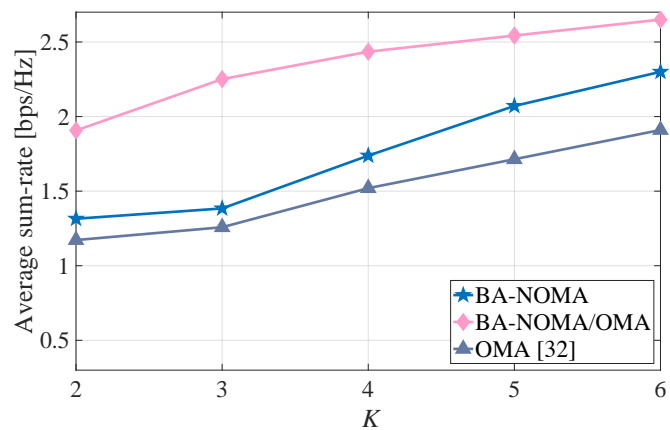


Fig. 4. Average sum-rate for varying K , $L = 4$ and $\text{SNR} = 16$ dB.

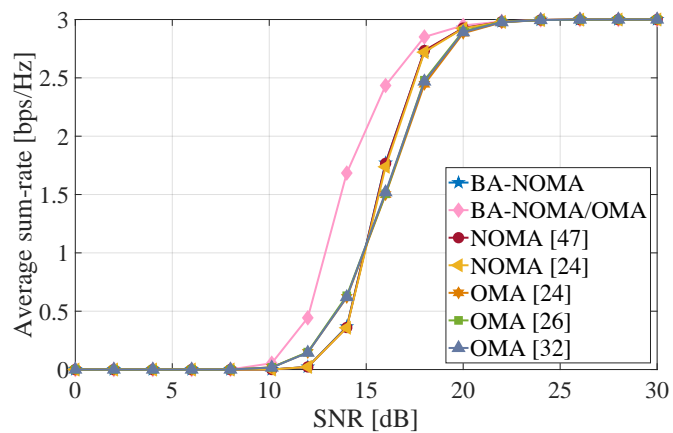


Fig. 3. Average sum-rate for $K = 4$, $L = 4$ and various algorithms.

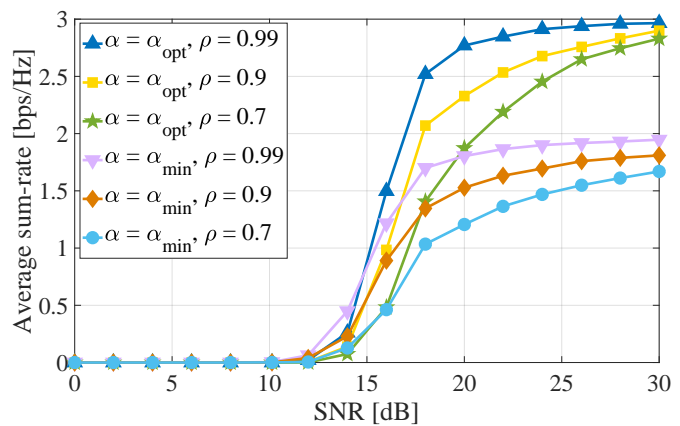


Fig. 5. Average sum-rate for varying ρ , $K = 4$, $L = 4$ and different α values.

asymmetries and the required rate. Again, it can be seen that no differentiation is observed within each multiple-access family, due to $\{R \rightarrow D\}$ prioritization.

The next comparison in Fig. 4 investigates the average sum-rate performance of the BA-NOMA, hybrid BA-NOMA/OMA and the OMA of [32] for increasing number of relays K and fixed transmit SNR equal to 16 dB. All the algorithms are based on broadcasting in the $\{S \rightarrow R\}$ links and $\{R \rightarrow D\}$ prioritization. It can be seen that the hybrid BA-NOMA/OMA outperforms the other multiple-access cases, due to its flexible scheduling in the $\{R \rightarrow D\}$ links, switching among NOMA and OMA. In addition, with each relay that is added to the network, the sum-rate increases, due to increased diversity. For the other two multiple access cases, it can be seen that for a small number of relays similar performance is achieved. However, when more relays are available, BA-NOMA provides better performance, as both users can be served simultaneously and the increased diversity is better exploited, compared to the case of OMA.

Then, Fig. 5 depicts the sum-rate of the network when α is determined using outdated CSI. So, various correlation coefficient ρ values are examined in order to illustrate the impact of outdated CSI on BA-NOMA depending on the adopted α value. More specifically, for all ρ cases, when

$\alpha = \alpha_{min}$, the sum-rate of BA-NOMA is significantly reduced, compared to the case of employing an α value using eq. (12), denoted as α_{opt} . It can be observed that independently of ρ , α_{opt} allows BA-NOMA to reach the upper bound for the sum-rate for this topology where fixed rated transmissions are performed.

3) *Average delay*: Then, Fig. 6 includes average delay results, where the performance of each user is shown when $K = 4$ and $L = 4$, for the hybrid BA-NOMA/OMA and OMA algorithms. On the contrary, for the NOMA algorithms, the users are not distinguished, as the same delay is experienced, since they are served simultaneously at each time-slot. For lower SNR values, it can be seen that hybrid BA-NOMA/OMA provides the best delay performance to both users, while D_1 experiences reduced delay than D_2 . As the transmit SNR increases, hybrid NOMA/OMA and NOMA algorithms provide almost equivalent delay performance, with some differences, due to instances of switching for the hybrid algorithm, resulting in cases where packets for one user might not exist and thus, less NOMA transmissions are performed. Nonetheless, after 22 dB, hybrid BA-NOMA/OMA relies on NOMA transmissions and so, delay performance is identical to NOMA algorithms. Moreover, there are some differences between OMA algorithms, with the algorithm of [32] providing

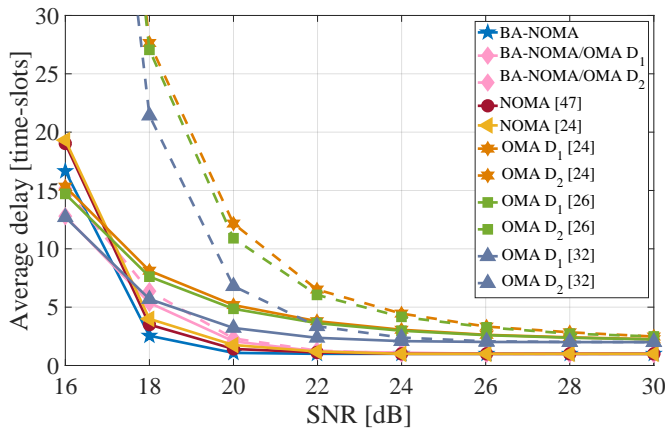


Fig. 6. Average delay for $K = 4$, $L = 4$ and various algorithms.

slightly better performance, since more packets are available due to broadcasting in the $\{S \rightarrow R\}$ links. For the OMA algorithms, it must be noted that increased delay is experienced by the weak user D_2 , due to excessive outages when this user is scheduled. However, when transmit conditions improve, both users achieve similar delay performance. Overall, after 22 dB and when delay performance is stabilized, the adoption of hybrid BA-NOMA/OMA and NOMA algorithms is able to halve the delay in comparison to OMA.

B. Different rate requirements

In the second scenario, the two users require different rates $r_1 = 4$ BPCU and $r_2 = 1$ BPCU. It is well known that NOMA is superior to OMA when increased asymmetry between the users' channels is experienced. In this comparison, the first user D_1 is assumed to have stronger channels towards the relays than D_2 , an asymmetry described by $\sigma_{R_k D_1}^2 = 2\sigma_{R_k D_2}^2$. More importantly, as increased asymmetry in rate requirements is considered, this scenario might correspond to a network where D_1 is a cellular user demanding a high rate service, while D_2 is an IoT device.

1) *Outage probability*: The outage performance of the network for $K = 4$, $L = 4$ and different rate requirements for each user is illustrated at Fig. 7. It can be observed that due to the increased asymmetry stemming from different channel conditions and rate requirements, NOMA algorithms offer significantly reduced outage probability compared to OMA algorithms. Also, all the NOMA algorithms provide similar behavior with a slight improvement offered by BA-NOMA and NOMA of [47], where broadcasting in the $\{S \rightarrow R\}$ links and the selection of the relay with the maximum buffer size are performed, respectively. However, due to $\{R \rightarrow D\}$ prioritization, diversity gain is equal for all NOMA and OMA algorithms. It must be noted that, as transmit SNR increases, there is almost a 10 dB coding gain between NOMA and OMA.

2) *Average sum-rate*: In the next comparison, the average sum-rate performance is evaluated in Fig. 8 for $K = 4$ and $L = 4$. In this case, it can be observed that for low SNR, OMA is able to serve D_2 with the required low rate. On

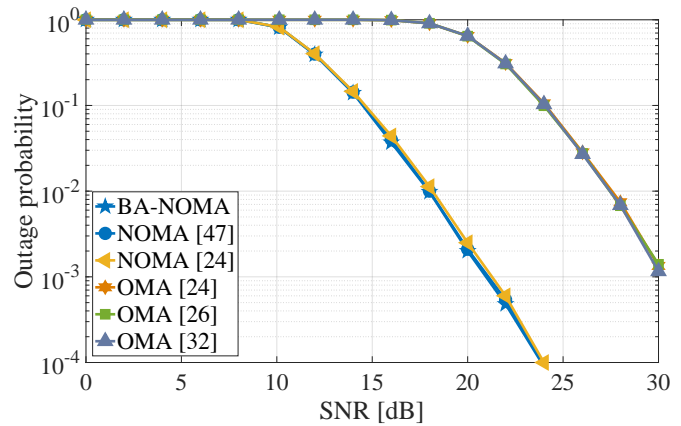


Fig. 7. Outage probability for $K = 4$, $L = 4$ and various algorithms.

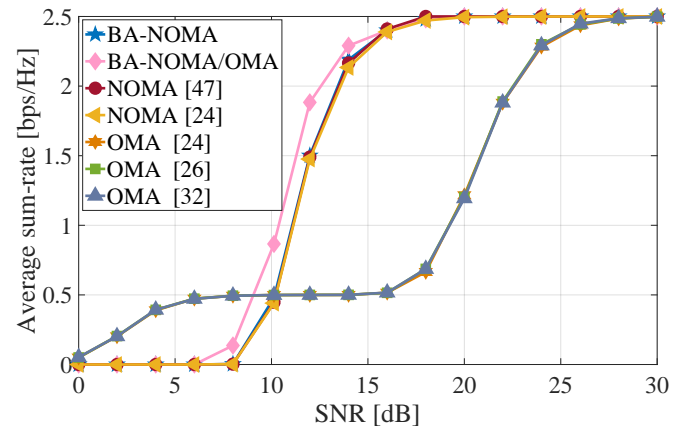


Fig. 8. Average sum-rate for $K = 4$, $L = 4$ and various algorithms.

the other hand, hybrid BA-NOMA/OMA cannot serve the users until the transmit SNR reaches a value of 8 dB. This can be explained by the rate requirement in the $\{S \rightarrow R\}$ links where $r_1 + r_2 = 5$ BPCU must be achieved. As a result, outages occur for this algorithm. In addition, BA-NOMA requires a transmit SNR above 10 dB to operate. However, as transmit SNR increases hybrid BA-NOMA/OMA and NOMA algorithms provide the best performance due to the robust selection of α , considering both the channels and the rates of the users. More importantly, hybrid BA-NOMA/OMA provides increased performance gains, as it overcomes NOMA outages in the $\{R \rightarrow D\}$ links by switching to OMA.

Another performance comparison for the average sum-rate is depicted in Fig. 9. Here, the performance is evaluated versus the increasing number of relays and fixed transmit SNR equal to 16 dB. Starting from the hybrid BA-NOMA/OMA, it can be observed that independently of K , it provides the highest sum-rate compared to BA-NOMA and OMA. Nonetheless, for this transmit SNR, BA-NOMA operates efficiently and is able to provide a high sum-rate that reaches that of hybrid BA-NOMA/OMA as K increases, thus exploiting the increased diversity. On the contrary, for this transmit SNR value, OMA cannot improve its performance and only the low rate user is served.

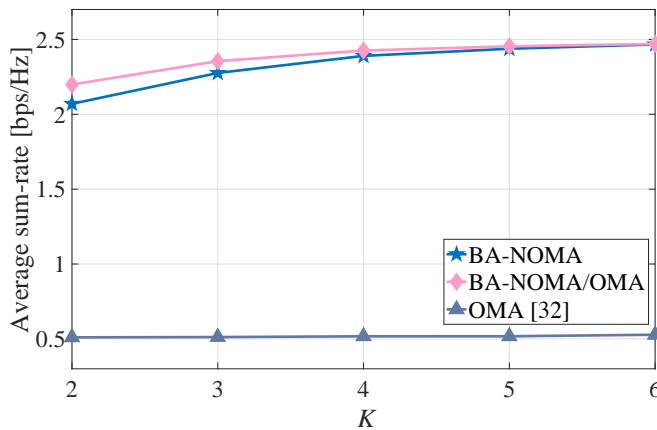


Fig. 9. Average sum-rate for varying K , $L = 4$ and $\text{SNR} = 16$ dB.

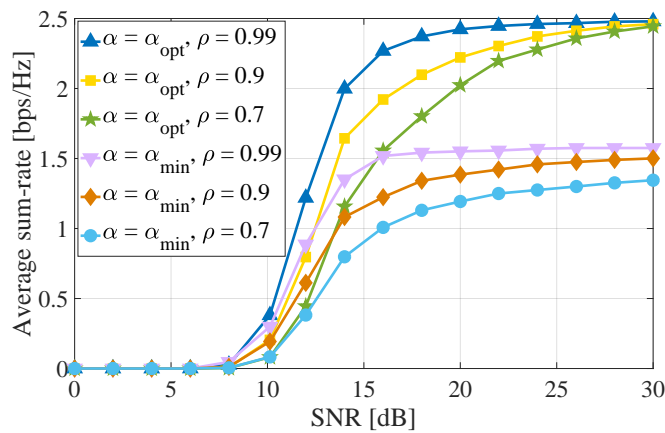


Fig. 10. Average sum-rate for varying ρ , $K = 4$, $L = 4$ and different α values.

After, Fig. 10 shows the sum-rate of the network under various correlation coefficient ρ values, when outdated CSI is used for determining α for BA-NOMA. It is clear that for all the values of ρ , $\alpha = \alpha_{min}$ degrades the performance of BA-NOMA. On the contrary, choosing α based on eq. (12) provides robustness against delays in CSI acquisition. As a result, the fixed end-to-end sum-rate of 2.5 bps/Hz can be achieved and both users can be served.

3) *Average delay*: For this highly asymmetric scenario in terms of rate requirements, delay results are included in Fig. 11. The OMA algorithms, due to the fixed scheduling and the absence of interference are able to support the low rate user D_2 , even for low transmit SNR values. Nonetheless, the high rate user D_1 experiences outages and its delay performance improves after 20 dB. The hybrid BA-NOMA/OMA algorithm is able to adequately serve both users at lower SNR and as conditions improve, it achieves significantly improved delay performance compared to OMA algorithms. Finally, NOMA algorithms demand high transmit SNR in order to reduce the average delay of the users and especially BA-NOMA is able to improve the delay performance at a lower SNR by exploiting the existence of more packets, due to broadcasting in the $\{S \rightarrow R\}$ links. For asymptotically high transmit SNR, hybrid BA-NOMA/OMA and NOMA algorithms offer half the delay

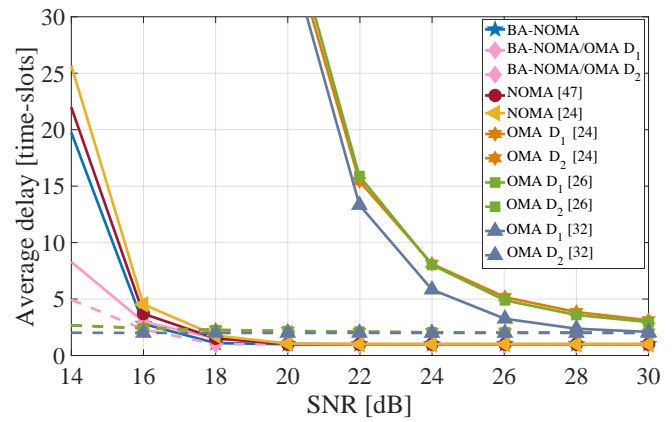


Fig. 11. Average delay for $K = 4$, $L = 4$ and various algorithms.

compared to the OMA algorithms, as users are simultaneously served at each time-slot.

VII. CONCLUSIONS AND FUTURE DIRECTIONS

In this paper, a two-hop topology where two users communicate with a source through NOMA was examined. Also, the process of determining the power allocation coefficient α was presented. Then, two BA relay selection algorithms were proposed for NOMA and hybrid NOMA/OMA relying on broadcasting from the source in order to maintain the diversity of the network and reduce the implementation complexity. Finally, performance evaluation was conducted for users with equal and different rate requirements and asymmetric channel conditions. It was shown that the proposed BA-ORS algorithms can improve the performance of NOMA, and especially hybrid NOMA/OMA networks, in terms of average sum-rate and delay due to efficient scheduling and $\{R \rightarrow D\}$ prioritization.

Possible future directions consist of developing selection algorithms, aiming to maximize the sum-rate of multi-source networks and the consideration of BA full-duplex relaying. Towards this end, scenarios with both direct $\{S \rightarrow D\}$ connectivity and relay selection can provide additional degrees of freedom to the transmission [52], thus improving the performance of NOMA. Also, practical issues, such as outdated CSI [9]–[11] and imperfect SIC [53] should be considered when selecting the relay in NOMA and hybrid NOMA/OMA networks.

APPENDIX A PROOF OF PROPOSITION 1

The probability of D_2 being in outage in BA-NOMA algorithm is given by

$$p_{\text{out}, D_2} \triangleq \mathbb{P} \left(\frac{(1 - \alpha)\phi g_{R_k D_2}}{\alpha\phi g_{R_k D_2} + 1} < \gamma_2 \right). \quad (27)$$

Solving in the parenthesis with respect to $g_{R_k D_2}$ we obtain

$$p_{\text{out}, D_2} = \mathbb{P} \left(g_{R_k D_2} < \frac{\gamma_2}{[1 - \alpha(1 + \gamma_2)]\phi} \right) \quad (28)$$

Let $F_W(w)$ denote the cumulative distribution function (cdf) of a random variable W ; e.g., for the exponential distribution $F_W(w) = 1 - \exp(-\lambda w)$. Then, from (28), p_{out,D_1} is given by

$$p_{\text{out},D_2} = 1 - e^{-\frac{\gamma_2}{[1-\alpha(1+\gamma_2)]\phi}}.$$

The probability of D_1 being in outage in BA-NOMA algorithm is that either the decoding of x_2 has failed (event A) or the decoding of x_1 has failed (event B). Since these two events are not independent, the probability of D_1 being in outage can be shown to be

$$\begin{aligned} p_{\text{out},D_1} &= \mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B) \\ &= 1 - e^{-\frac{\gamma_1}{[1-\alpha(1+\gamma_1)]\phi}} + 1 - e^{-\frac{\gamma_2}{\alpha\phi}} \\ &\quad + \min \left\{ 1 - e^{-\frac{\gamma_1}{[1-\alpha(1+\gamma_1)]\phi}}, 1 - e^{-\frac{\gamma_2}{\alpha\phi}} \right\} \\ &= 1 - e^{-\frac{\gamma_1}{\min\{[1-\alpha(1+\gamma_1)], \alpha\}\phi}}. \end{aligned}$$

The proof is complete. \square

APPENDIX B

PROOF OF PROPOSITION 2

As in the proof of Proposition (1) in Appendix (A), for D_1 the event of failing to decode x_2 is denoted by A , and for D_2 the event of failing to decode of x_w is denoted by B . In this case, the events are independent (since the channels are assumed to be independently distributed) and, therefore,

$$\begin{aligned} p_{\text{out},D} &= \mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A)\mathbb{P}(B) \\ &= 1 - e^{-\frac{\gamma_1}{\alpha_1\phi}} + 1 - e^{-\frac{\gamma_2}{\alpha_2\phi}} - \left(1 - e^{-\frac{\gamma_1}{\alpha_1\phi}}\right) \left(1 - e^{-\frac{\gamma_2}{\alpha_2\phi}}\right) \\ &= 1 - e^{-\left(\frac{\gamma_1}{\alpha_1\phi} + \frac{\gamma_2}{\alpha_2\phi}\right)}, \end{aligned}$$

where $\alpha_1 \triangleq \min\{[1-\alpha(1+\gamma_1)], \alpha\}$ and $\alpha_2 \triangleq 1-\alpha(1+\gamma_2)$. The proof is complete. \square

REFERENCES

- [1] L. Lei, D. Yuan, C. K. Ho and S. Sun, "Power and channel allocation for non-orthogonal multiple access in 5G systems: Tractability and computation," *IEEE Trans. on Wireless Commun.*, vol. 15, no. 12, pp. 8580–8594, Dec. 2016.
- [2] L. Dai, B. Wang, Y. Yuan, S. Han, C. I. I and Z. Wang, "Non-orthogonal multiple access for 5G: Solutions, challenges, opportunities, and future research trends," *IEEE Commun. Mag.*, vol. 53, no. 9, pp. 74–81, Sept. 2015.
- [3] Kenichi Higuchi and Anass Benjebbour, "Non-orthogonal multiple access (NOMA) with successive interference cancellation for future radio access" *IEICE Trans. on Commun.*, vol. E98.B, no. 3, pp. 403–414, Mar. 2015.
- [4] H. Hacı, H. Zhu and J. Wang, "Performance of non-orthogonal multiple access with a novel asynchronous interference cancellation technique," *IEEE Trans. on Commun.*, vol. 65, no. 3, pp. 1319–1335, March 2017.
- [5] S. M. R. Islam, N. Avazov, O. A. Dobre and K. S. Kwak, "Power-domain non-orthogonal multiple access (NOMA) in 5G systems: Potentials and challenges," *IEEE Commun. Surv. & Tut.*, vol. 19, no. 2, pp. 721–742, Secondquarter 2017.
- [6] S. M. R. Islam, M. Zeng, O. A. Dobre and K. Kwak, "Resource allocation for downlink NOMA systems: Key techniques and open issues," *IEEE Wireless Commun.*, vol. 25, no. 2, pp. 40–47, April 2018.
- [7] B. Di, L. Song and Y. Li, "Sub-channel assignment, power allocation, and user scheduling for non-orthogonal multiple access networks," *IEEE Trans. on Wireless Commun.*, vol. 15, no. 11, pp. 7686–7698, Nov. 2016.
- [8] F. Fang, H. Zhang, J. Cheng and V. C. M. Leung, "Energy-efficient resource allocation for downlink non-orthogonal multiple access network," *IEEE Trans. on Commun.*, vol. 64, no. 9, pp. 3722–3732, Sept. 2016.
- [9] Z. Ding, X. Lei, G. K. Karagiannidis, R. Schober, J. Yuan and V. K. Bhargava, "A survey on non-orthogonal multiple access for 5G networks: Research challenges and future trends," *IEEE J. on Select. Areas Commun.*, vol. 35, no. 10, pp. 2181–2195, Oct. 2017.
- [10] Y. Liu, Z. Qin, M. El-kashlan, Z. Ding, A. Nallanathan and L. Hanzo, "Nonorthogonal multiple access for 5G and beyond," *Proc. IEEE*, vol. 105, no. 12, pp. 2347–2381, Dec. 2017.
- [11] Y. Wang, B. Ren, S. Sun, S. Kang and X. Yue, "Analysis of non-orthogonal multiple access for 5G," *China Commun.*, vol. 13, no. Supplement2, pp. 52–66, 2016.
- [12] S. Shi, L. Yang and H. Zhu, "Outage balancing in downlink nonorthogonal multiple access with statistical channel state information," *IEEE Trans. on Wireless Commun.*, vol. 15, no. 7, pp. 4718–4731, July 2016.
- [13] M. F. Hanif, Z. Ding, T. Ratnarajah and G. K. Karagiannidis, "A minorization-maximization method for optimizing sum rate in the downlink of non-orthogonal multiple access systems," *IEEE Trans. on Signal Process.*, vol. 64, no. 1, pp. 76–88, Jan. 2016.
- [14] M. Zeng, A. Yadav, O. A. Dobre, G. I. Tsiropoulos and H. V. Poor, "Capacity comparison between MIMO-NOMA and MIMO-OMA with multiple users in a cluster," *IEEE J. on Select. Areas Commun.*, vol. 35, no. 10, pp. 2413–2424, Oct. 2017.
- [15] J. N. Laneman, D. N. C. Tse and G. W. Wornell, "Cooperative diversity in wireless networks: Efficient protocols and outage behavior," *IEEE Trans. Inform. Theory*, vol. 50, pp. 3062–3080, Dec. 2004.
- [16] A. Bletsas, A. Khisti, D. Reed and A. Lippman, "A simple cooperative diversity method based on network path selection," *IEEE J. Select. Areas Commun.*, vol. 24, pp. 659–672, March 2006.
- [17] D. S. Michalopoulos and G. K. Karagiannidis, "Performance analysis of single relay selection in Rayleigh fading," *IEEE Trans. Wireless Commun.*, vol. 7, pp. 3718–3724, Oct. 2008.
- [18] N. Nomikos, T. Charalambous, I. Krikidis, D. N. Skoutas, D. Vouyioukas, M. Johansson, C. Skianis, "A survey on buffer-aided relay selection," *IEEE Commun. Surv. & Tut.*, vol. 18, no. 2, pp. 1073–1097, Secondquarter 2016.
- [19] N. Zlatanov, A. Ikhlef, T. Islam and R. Schober, "Buffer-aided cooperative communications: opportunities and challenges," *IEEE Commun. Mag.*, vol. 52, no. 4, pp. 146–153, April 2014.
- [20] A. Ikhlef, D. S. Michalopoulos and R. Schober, "Max-max relay selection for relays with buffers," *IEEE Trans. Wireless Commun.*, vol. 11, pp. 1124–1135, March 2012.
- [21] I. Krikidis, T. Charalambous and J. S. Thompson, "Buffer-aided relay selection for cooperative diversity systems without delay constraints," *IEEE Trans. Wireless Commun.*, vol. 11, pp. 1957–1967, May 2012.
- [22] M. Oiwa, C. Tosa and S. Sugiura, "Theoretical analysis of hybrid buffer-aided cooperative protocol based on max-max and max-link relay selections," *IEEE Trans. on Vehic. Tech.*, vol. 65, no. 11, pp. 9236–9246, Nov. 2016.
- [23] D. Poulimeneas, T. Charalambous, N. Nomikos, I. Krikidis, D. Vouyioukas, M. Johansson, "A delay-aware hybrid relay selection policy," *IEEE Int. Conf. on Telecomm., (ICT)*, May 2016.
- [24] Z. Tian, Y. Gong, G. Chen and J. Chambers, "Buffer-aided relay selection with reduced packet delay in cooperative networks," *IEEE Trans on Vehic. Tech.*, vol. 66, no. 3, pp. 2567–2575, March 2017.
- [25] S. Luo and K. C. Teh, "Buffer state based relay selection for buffer-aided cooperative relaying systems," *IEEE Trans. on Wireless Commun.*, vol. 14, no. 10, pp. 5430–5439, Oct. 2015.
- [26] N. Nomikos, D. Poulimeneas, T. Charalambous, I. Krikidis, D. Vouyioukas and M. Johansson, "Delay- and diversity-aware buffer-aided relay selection policies in cooperative networks," *IEEE Access*, Nov. 2018.
- [27] S. L. Lin and K. H. Liu, "Relay selection for cooperative relaying networks with small buffers," *IEEE Trans. on Vehic. Tech.*, vol. 65, no. 8, pp. 6562–6572, Aug. 2016.
- [28] W. Wicke, N. Zlatanov, V. Jamali and R. Schober, "Buffer-aided relaying with discrete transmission rates for the two-hop half-duplex relay network," *IEEE Trans. on Wireless Commun.*, vol. 16, no. 2, pp. 967–981, Feb. 2017.
- [29] B. Zhou, Y. Cui and M. Tao, "Stochastic throughput optimization for two-hop systems with finite relay buffers," *IEEE Trans. on Signal Proc.*, vol. 63, no. 20, pp. 5546–5560, Oct. 2015.
- [30] M. Oiwa, R. Nakai, S. Sugiura, "Buffer-state-and-thresholding-based amplify-and-forward cooperative networks," *IEEE Wireless Commun. Lett.*, vol. 6, no. 5, pp. 674–677, Oct. 2017.

- [31] M. Oiwa and S. Sugiura, "Reduced-packet-delay generalized buffer-aided relaying protocol: Simultaneous activation of multiple source-to-relay links," *IEEE Access*, vol. 4, no., pp. 3632–3646, June 2016.
- [32] N. Nomikos, T. Charalambous, D. Vouyioukas and G. K. Karagiannidis, "Low-complexity buffer-aided link selection with outdated CSI and feedback errors," *IEEE Trans. on Commun.*, vol. 66, no. 8, pp. 3694–3706, March 2018.
- [33] T. Islam, D. S. Michalopoulos, R. Schober and V. K. Bhargava, "Buffer-aided relaying with outdated CSI," *IEEE Trans. Wireless Commun.*, vol. 15, no. 3, pp. 1979–1997, March 2016.
- [34] V. Jamali, N. Waly, N. Zlatanov and R. Schober, "Optimal buffer-aided relaying with imperfect CSI," *IEEE Commun. Lett.*, vol. 20, no. 7, pp. 1309–1312, July 2016.
- [35] Z. Ding, M. Peng and H. V. Poor, "Cooperative non-orthogonal multiple access in 5G systems," *IEEE Commun. Lett.*, vol. 19, no. 8, pp. 1462–1465, Aug. 2015.
- [36] J. Men and J. Ge, "Performance analysis of non-orthogonal multiple access in downlink cooperative network," *IET Commun.*, vol. 9, no. 18, pp. 2267–2273, Dec. 2015.
- [37] J. So and Y. Sung, "Improving non-orthogonal multiple access by forming relaying broadcast channels," *IEEE Commun. Lett.*, vol. 20, no. 9, pp. 1816–1819, Sept. 2016.
- [38] Y. Li, M. Jiang, Q. Zhang, Q. Li and J. Qin, "Cooperative non-orthogonal multiple access in multiple-input-multiple-output channels," *IEEE Trans. on Wireless Commun.*, vol. 17, no. 3, pp. 2068–2079, March 2018.
- [39] C. Zhong and Z. Zhang, "Non-orthogonal multiple access with cooperative full-duplex relaying," *IEEE Commun. Lett.*, vol. 20, no. 12, pp. 2478–2481, Dec. 2016.
- [40] M. F. Kader, M. B. Shahab and S. Y. Shin, "Exploiting non-orthogonal multiple access in cooperative relay sharing," *IEEE Commun. Lett.*, vol. 21, no. 5, pp. 1159–1162, May 2017.
- [41] 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.
- [42] S. Lee, D. Benevides da Costa and T. Q. Duong, "Outage probability of non-orthogonal multiple access schemes with partial relay selection," *IEEE Int. Symp. on Pers., Ind., and Mob. Radio Commun. (PIMRC)*, pp. 1–6, Dec. 2016.
- [43] F. Fang, H. Zhang, J. Cheng, S. Roy and V. C. M. Leung, "Joint user scheduling and power allocation optimization for energy-efficient NOMA systems with imperfect CSI," *IEEE J. Select. Areas Commun.*, vol. 35, no. 12, pp. 2874–2885, Dec. 2017.
- [44] J. Men, J. Ge and C. Zhang, "Performance analysis for downlink relaying aided non-orthogonal multiple access networks with imperfect CSI over Nakagami- m fading," *IEEE Access*, vol. 5, pp. 998–1004, March 2017.
- [45] S. Luo and K. C. Teh, "Adaptive transmission for cooperative NOMA system with buffer-aided relaying," *IEEE Commun. Lett.*, vol. 21, no. 4, pp. 937–940, April 2017.
- [46] Q. Zhang, Z. Liang, Q. Li and J. Qin, "Buffer-aided non-orthogonal multiple access relaying systems in Rayleigh fading channels," *IEEE Trans. on Commun.*, vol. 65, no. 1, pp. 95–106, Jan. 2017.
- [47] N. Nomikos, T. Charalambous, D. Vouyioukas, G. K. Karagiannidis and R. Wichman, "Relay selection for buffer-aided non-orthogonal multiple access networks," *IEEE GLOBECOM Works.*, Dec. 2017.
- [48] M. Shirvanimoghaddam, M. Dohler and S. J. Johnson, "Massive non-orthogonal multiple access for cellular IoT: Potentials and limitations," *IEEE Commun. Mag.*, vol. 55, no. 9, pp. 55–61, Sept. 2017.
- [49] M. Shirvanimoghaddam, M. Condoluci, M. Dohler and S. J. Johnson, "On the fundamental limits of random non-orthogonal multiple access in cellular massive IoT," *IEEE J. Select. Areas Commun.*, vol. 35, no. 10, pp. 2238–2252, Oct. 2017.
- [50] J. L. Vicario, A. Bel, J. A. Lopez-Salcedo and G. Seco, "Opportunistic relay selection with outdated CSI: Outage probability and diversity analysis," *IEEE Trans. Wireless Commun.*, vol. 8, pp. 2872–2876, June 2009.
- [51] A. Goldsmith, "Wireless communications", *Cambridge University Press*, August 2005.
- [52] T. Charalambous, N. Nomikos, I. Krikidis, D. Vouyioukas and M. Johansson, "Modeling buffer-aided relay selection in networks with direct transmission capability," *IEEE Commun. Lett.*, vol. 19, no. 4, pp. 649–652, April 2015.
- [53] T. Manglayev, R. C. Kizilirmak, Y. H. Kho, N. Bazhayev and I. Lebedev, "NOMA with imperfect SIC implementation," *IEEE Intern. Conf. on Smart Tech. (EUROCON)*, July 2017.