

Low Latency Driven Effective Capacity Analysis for Non-orthogonal and Orthogonal Spectrum Access

Wenjuan Yu¹, Leila Musavian², Atta ul Quddus¹, Qiang Ni³ and Pei Xiao¹

¹5GIC, Institute of Communication Systems, University of Surrey, Guildford, GU2 7XH, UK

²School of Computer Science and Electronic Engineering, University of Essex, UK

³School of Computing and Communications, InfoLab21, Lancaster University, UK

Emails: {w.yu, a.quddus, p.xiao}@surrey.ac.uk, leila.musavian@essex.ac.uk, q.ni@lancaster.ac.uk

Abstract—In this paper, we theoretically investigate the performance of non-orthogonal and orthogonal spectrum access protocols (more generically known as NOMA) in supporting ultra-reliable low-latency communications (URLLC). The theory of effective capacity (EC) is adopted as a suitable delay-guaranteed capacity metric to flexibly represent the users' delay requirements. Then, the total EC difference between a downlink user-paired NOMA network and a downlink orthogonal multiple access (OMA) network is analytically studied. Exact closed-form expressions and the approximated closed-forms at high signal-to-noise ratios (SNRs) are derived for both networks and validated through simulation results. It is shown that for a user pair in which two users with the most distinct channel conditions are paired together, NOMA still achieves higher total EC (compared to OMA) in high SNR regime as the user group size becomes larger, although the EC performance of both NOMA and OMA reduces with the increase in group size. It is expected that the derived analytical framework can serve as a useful reference and practical guideline for designing favourable orthogonal and non-orthogonal spectrum access schemes in supporting low-latency services.

Index Terms—URLLC, NOMA, effective capacity, delay violation probability, exact closed-form expressions.

I. INTRODUCTION

With the explosive growth of delay-sensitive wireless communication applications and services such as vehicular communications, tactile Internet and virtual reality, the ability of supporting low-latency services becomes more and more important. The provision of ultra-reliable low-latency communications (URLLC) has also emphasized the importance of reliability and low-latency transmissions in the 5th generation (5G) cellular networks and beyond [2]. In this paper, we adopt the concept of effective capacity (EC) to describe the maximum constant arrival rate that can be served, while guaranteeing a statistical delay requirement [3]. From the introduction about EC given in Section III, it will become evident that EC is not only a suitable, but also a flexible, metric which can represent various delay requirements. This allows us to investigate the delay-constrained data rates that can be served by a proposed communication network, for different latency requirements.

Furthermore, due to the spectrum scarcity issue in future communication networks, it is important to study and investigate the performance of spectrum-efficient multiple access

protocols. In recent years, non-orthogonal multiple access (NOMA)¹ was introduced as a promising multiple access (MA) technique for 5G communications and it has been shown to be spectrally efficient since it allows multiple users to transmit with different power levels, but using the same subcarriers and time slots. Further, NOMA has been shown to have some advantages over conventional orthogonal multiple access (OMA) in different communication systems [4]–[9]. For a downlink transmission with NOMA applied, the base station (BS) will send a superimposed mixture containing all users' signals, then the users with stronger channel conditions can obtain the information of the users with weaker channel conditions in accordance with NOMA principle. This can be exploited to improve the weaker user's reception reliability [7]. Moreover, a NOMA network with fixed power allocation (F-NOMA) was compared with the conventional OMA in [5], which shows that F-NOMA can achieve a larger sum rate than OMA and if two users with very different channel conditions are paired, the performance of NOMA can be further improved. In [7], a non-orthogonal relaying scheme was proposed and studied, which improves the ergodic secrecy rate, compared to the conventional orthogonal relaying schemes.

Although the aforementioned studies were conducted to explore the applications of NOMA in different communication networks, they are more suitable for delay-insensitive services [10]. Focusing on delay-sensitive applications, the sum EC was formulated and maximized for a downlink NOMA network in [11], which was then solved with a suboptimal power control policy. In [12], the max-min EC problem was first proved to be quasi-concave and then solved using a bisection-based power allocation algorithm. However, [11] and [12] focused on designing efficient power allocation algorithms, rather than analytically analyzing the performance of NOMA in supporting delay-sensitive services. Considering the users with and without delay requirements, the exact analytical closed-form expressions for the ECs were proposed in [1] for a two-user NOMA network only. However, when *multiple* NOMA pairs are considered, the exact and the approximated closed-form expressions at high transmit signal-to-noise ratios (SNRs) for EC are not provided in [1].

Part of this conference paper was presented in [1].

¹We only consider power-domain NOMA in this work.

To investigate orthogonal and non-orthogonal multiple access protocols for low-latency services, in this paper we conduct a comparative study on a downlink user-paired NOMA network and a classical OMA network. The derived analytical framework can serve as a useful reference and practical guideline for designing favourable spectrum access protocols in supporting delay-sensitive communication scenarios. Specifically, we summarize the primary contributions as follows:

- For a downlink NOMA network with multiple pairs and a downlink OMA network, the exact analytical closed-form expressions are derived respectively, followed by the approximated closed-forms for EC at high SNRs. The accuracy of these obtained closed-form expressions is validated through simulation results.
- In Lemma 1, we prove that the performance gain of NOMA over OMA, in supporting delay-sensitive services, becomes stable in the high SNR regime. Simulation results confirm the theoretical analysis and further reveal the impact of delay requirements, power coefficient settings and user group size on the total EC difference between NOMA and OMA schemes.

II. SYSTEM MODELS

Consider a classical downlink transmission with one BS and K mobile users with single antenna employed. Without loss of generality, the K users' channels are assumed to be ordered as $|h_1|^2 \leq |h_2|^2 \leq \dots \leq |h_K|^2$, where h_k is the channel gain from the BS to the k^{th} user, following Rayleigh fading distribution with unit variance. Block fading channel models are assumed, i.e., the channel gain stays fixed within one fading-block and independently varies for the next fading-block. The fading-block length, T_f , is assumed to be equal to one frame size, which is an integer multiple of the symbol length. Two downlink networks with different multiple access protocols are described as follows.

A. A Downlink NOMA Network with Multiple User Pairs

For the downlink NOMA network, all K users are grouped into $K/2$ NOMA pairs², with the j^{th} pair described as $\phi_j = \{(u_j, w_j) \mid u_j \neq w_j, |h_{u_j}|^2 \leq |h_{w_j}|^2\}$, $\forall j \in \{1, 2, \dots, K/2\}$. Within each pair, NOMA is applied and for the inter-pair MA, the conventional OMA is utilized. For the j^{th} pair, the BS will send $\sqrt{\alpha_{u_j} P} s_{u_j} + \sqrt{\alpha_{w_j} P} s_{w_j}$ to both users on the same subcarrier/time-slot in accordance with the NOMA protocol. Here, α_{u_j} and α_{w_j} are the power coefficients for the users u_j and w_j , P is the total transmission power at the BS, s_{u_j} and s_{w_j} are the messages for the users u_j and w_j . By employing successive interference cancellation (SIC), the stronger user w_j with higher channel gains can detect and eliminate the message for the weaker user u_j , before it decodes its own. The user u_j with smaller channel gains will decode its information by considering the user w_j 's message as noise. It is assumed that $R_{u_j \rightarrow w_j} \geq \tilde{R}_{u_j}$ [5], so that SIC is guaranteed to be successfully applied at the user w_j . Here, $R_{u_j \rightarrow w_j}$ denotes the user w_j 's data rate to decode the user u_j 's message and \tilde{R}_{u_j} is

the target data rate for the user u_j . Assume that the target rate \tilde{R}_{u_j} equals to its achievable rate when the user u_j decodes its own message. Hence, it can be found that the requirement $R_{u_j \rightarrow w_j} \geq \tilde{R}_{u_j}$ always holds since $|h_{u_j}|^2 \leq |h_{w_j}|^2$ [5].

Fixed power coefficients are considered for all NOMA pairs in this paper, which are non-adaptive and follow NOMA principle, i.e., $\alpha_j = \{(\alpha_{u_j}, \alpha_{w_j}) \mid \alpha_{u_j} \geq \alpha_{w_j}, \alpha_{u_j} + \alpha_{w_j} = 1\}$, $\forall j \in \{1, 2, \dots, K/2\}$. Different non-adaptive transmit power policies will be simulated and discussed in Section IV, for the purpose of providing guidance on designing practical power allocation algorithms in future studies.

Then, for the j^{th} NOMA pair, the achievable data rates for both users can be respectively given by

$$R_{u_j} = \frac{2}{K} \log_2 \left(1 + \frac{\rho \alpha_{u_j} |h_{u_j}|^2}{\rho \alpha_{w_j} |h_{u_j}|^2 + 1} \right), \quad (\text{b/s/Hz}) \quad (1a)$$

$$R_{w_j} = \frac{2}{K} \log_2 (1 + \rho \alpha_{w_j} |h_{w_j}|^2), \quad (\text{b/s/Hz}) \quad (1b)$$

where ρ is the transmit SNR, i.e., $\rho = \frac{P}{N_0 B}$. Here, $N_0 B$ is the noise power with N_0 indicating the one-sided noise spectral density and B denoting the channel bandwidth.

B. A Conventional Downlink OMA Network

For a classical downlink OMA network applying time division multiple access (TDMA), each user is assumed to equally occupy $1/K$ of orthogonal resources, which is a typical allocation strategy. Although in this case, the available radio resources for each user become less, the transmit power for each user can be higher. Then, the achievable rate for the user k , in b/s/Hz, is given by

$$\bar{R}_k = \frac{1}{K} \log_2 (1 + \rho |h_k|^2), \quad k \in \{1, 2, \dots, K\}. \quad (2)$$

III. EFFECTIVE CAPACITY

Assume that there exists a virtual buffer for the k^{th} user at the BS, with an infinite buffer size. Define $D_k(t)$ as the delay experienced by a packet arriving at time t . According to [13], [14], the probability of the delay $D_k(t)$ exceeding a maximum delay bound D_{\max}^k can be estimated as

$$P_{\text{delay}}^{\text{out}} = \Pr\{D_k(t) > D_{\max}^k\} \approx \Pr\{Q(t) > 0\} e^{-\theta_k \mu D_{\max}^k}, \quad (3)$$

where $P_{\text{delay}}^{\text{out}}$ denotes the delay violation probability limit for the k^{th} user, $\Pr\{Q(t) > 0\}$ is the probability of a non-empty buffer at time t with $Q(t)$ indicating the buffer length, D_{\max}^k is the given delay bound in the unit of symbol duration, and θ_k ($\theta_k > 0$) represents the k^{th} user's exponential decay rate. It was proved that the constant arrival rate has to be limited to the value of μ , which equals to the EC, so that a target delay violation probability limit $P_{\text{delay}}^{\text{out}}$ can be met.

Assume that the given wireless link can support $C_k(t)$ packets per unit of time, which can be modeled as a stationary and ergodic random process with nonnegative values. Assuming that the service process satisfies Gärtner-Ellis theorem [15], the EC for the k^{th} user on a block-fading channel is defined as

$$E_c^k = -\frac{1}{\theta_k T_f B} \ln (\mathbb{E} [e^{-\theta_k T_f B C_k}]), \quad (\text{b/s/Hz}) \quad (4)$$

²Here, it is required that K is an even number.

where $\mathbb{E}[\cdot]$ is the expectation over the user k 's channel. Set $x_k = \rho|h_k|^2$. Since all K users' channels are assumed to be already ordered, the probability density function (PDF) of the ordered x_k , $\forall k \in \{1, \dots, K\}$, denoted by $f^{(k)}(x_k)$, follows the theory of order statistics and is given by [16]

$$f^{(k)}(x_k) = \psi_k f(x_k) F(x_k)^{k-1} (1 - F(x_k))^{K-k}, \quad (5)$$

where $\psi_k = 1/B(k, K - k + 1)$, in which the beta function $B(k, K - k + 1)$ equals to $\frac{k!(K - k + 1)!}{(K + 1)!}$ [17]. In (5), $f(x_k)$ and $F(x_k)$ are the PDF and the cumulative distribution function (CDF) of the unordered random variable x_k .

With different multiple access protocols applied, the available service rate for each user can be different. In the following, we study the total maximum achievable delay-constrained rate, i.e., the total EC, for the downlink transmission with NOMA applied and also the conventional OMA principle.

A. For A Downlink NOMA With Multiple User Pairs

By replacing C_k in (4) with the achievable rates (1a) and (1b), the individual EC for the users u_j and w_j in the j^{th} NOMA pair can be respectively given by

$$E_c^{u_j} = -\frac{1}{\theta_{u_j} T_f B} \ln \left(\mathbb{E} \left[\left(\frac{\rho |h_{u_j}|^2 + 1}{\rho \alpha_{w_j} |h_{u_j}|^2 + 1} \right)^{\frac{4}{K} \beta_{u_j}} \right] \right), \quad (6a)$$

$$E_c^{w_j} = -\frac{1}{\theta_{w_j} T_f B} \ln \left(\mathbb{E} \left[(1 + \rho \alpha_{w_j} |h_{w_j}|^2)^{\frac{4}{K} \beta_{w_j}} \right] \right), \quad (6b)$$

where $\beta_{u_j} = -\frac{\theta_{u_j} T_f B}{2 \ln 2}$ and $\beta_{w_j} = -\frac{\theta_{w_j} T_f B}{2 \ln 2}$. The total EC for multiple NOMA pairs, denoted by E_c^N , equals to $\sum_{j=1}^{K/2} (E_c^{u_j} + E_c^{w_j})$. We then provide the exact analytical closed-forms and approximations at high SNRs, for $E_c^{u_j}$ and $E_c^{w_j}$ in the downlink transmission with NOMA.

Theorem 1: The exact closed-forms for $E_c^{u_j}$ and $E_c^{w_j}$ are respectively given in (7a) and (7b) (shown on next page). At high SNRs, $\lim_{\rho \rightarrow \infty} E_c^{u_j} = \frac{2}{K} \log_2 \left(\frac{1}{\alpha_{w_j}} \right)$ and $E_c^{w_j}$ can be approximated by (7c).

Proof: The proof for the exact closed-forms is similar to that of two-user case in [1], so it is omitted here due to page limit. The proof for the approximated closed-forms at high SNRs is given in Appendix. ■

B. For A Conventional Downlink OMA Network

By inserting (2) into (4), the individual EC for the user k in a conventional OMA network can be given by

$$\bar{E}_c^k = -\frac{1}{\theta_k T_f B} \ln \left(\mathbb{E} \left[(1 + \rho |h_k|^2)^{\frac{2}{K} \beta_k} \right] \right). \quad (8)$$

The total EC in a conventional downlink OMA network, denoted by E_c^O , equals to $\sum_{k=1}^K \bar{E}_c^k$. The following two theorems derive the exact analytical closed-forms and also the closed-forms at high SNRs for \bar{E}_c^k , $\forall k \in \{1, \dots, K\}$.

Theorem 2: The exact closed-form for \bar{E}_c^k for the user k is given in (7d). At high SNRs, the approximated closed-form for \bar{E}_c^k is given in (7e).

Proof: The proof for the exact closed-forms is similar to that of two-user case in [1], hence it is omitted due to

page limit. The proof for the approximated closed-form at high SNRs is given in Appendix. ■

C. Comparison of Total Effective Capacity

The total EC difference between a downlink user-paired NOMA and a conventional OMA can be expressed as

$$\mathbf{E}_c^N - \mathbf{E}_c^O = \sum_{j=1}^{K/2} (E_c^{u_j} + E_c^{w_j}) - \sum_{k=1}^K \bar{E}_c^k, \quad (9)$$

and its exact closed-form and approximation at high SNRs can be obtained, by simply applying the analytical results given in (7a)-(7e).

Lemma 1: At extreme SNR values, $\lim_{\rho \rightarrow 0} \mathbf{E}_c^N - \mathbf{E}_c^O = 0$ and $\lim_{\rho \rightarrow \infty} \mathbf{E}_c^N - \mathbf{E}_c^O$ approaches a constant, given in (10).

$$\lim_{\rho \rightarrow \infty} \mathbf{E}_c^N - \mathbf{E}_c^O = \sum_{j=1}^{K/2} -\frac{1}{\theta_{u_j} T_f B} \ln \left(\frac{(\alpha_{w_j})^{-\frac{4}{K} \beta_{u_j}}}{\mathbb{E} \left[(|h_{u_j}|^2)^{\frac{2}{K} \beta_{u_j}} \right]} \right) - \frac{1}{\theta_{w_j} T_f B} \ln \left(\frac{(\alpha_{w_j})^{\frac{4}{K} \beta_{w_j}} \mathbb{E} \left[(|h_{w_j}|^2)^{\frac{4}{K} \beta_{w_j}} \right]}{\mathbb{E} \left[(|h_{w_j}|^2)^{\frac{2}{K} \beta_{w_j}} \right]} \right). \quad (10)$$

Proof: Please refer to Appendix G in [1]. ■

Numerical results in Section IV show that the downlink user-paired NOMA supports higher total EC values than the downlink OMA network, in high SNR regime. From Lemma 1, we can then conclude that the performance gain of NOMA over OMA approaches a constant value at very high SNRs.

Furthermore, from (10), we notice that the difference of the total ECs, i.e., $\mathbf{E}_c^N - \mathbf{E}_c^O$, depends on the user pairing setting ϕ_j and also the power coefficients α_j , $\forall j \in \{1, 2, \dots, K/2\}$. To investigate the influence of power coefficients, two different power settings are utilized and simulated in Section IV, in order to provide more comprehensive EC comparison results. The optimal user pairing algorithms are out of the scope of this paper, but will be considered as a future research topic.

IV. NUMERICAL RESULTS

In this section, we first validate the accuracy of the proposed exact analytical closed-forms and the approximated closed-forms at high SNRs for $E_c^{u_j}$, $E_c^{w_j}$ and \bar{E}_c^k . The theoretical insights given in Lemma 1 will be numerically confirmed as well. Then, the total EC comparison results between the multiple NOMA pairs and the downlink OMA network will be numerically studied. In our simulations, it is assumed that there are 6 mobile users, i.e., $K = 6$, and 3 NOMA pairs in total. Two different power coefficients settings are utilized: 1) a fixed power allocation; 2) a varied power setting according to [9], i.e., $\alpha_i = \frac{\mathcal{M}-i+1}{\delta}$ for all \mathcal{M} users sharing one radio resource. Here, δ is to ensure $\sum_{i=1}^{\mathcal{M}} \alpha_i = 1$. The delay exponents are assumed to be the same for all users and are set as 0.01, unless otherwise indicated. Further, it is assumed that $T_f = 0.5$ ms and $B = 180$ kHz.

Assume that the 1st user and the 4th user are paired together in the downlink NOMA network, i.e., $u_j = 1$ and $w_j = 4$. The power coefficients for this pair are given as $\alpha_{u_j} = 0.8$

$$E_c^{u_j} = -\frac{1}{\theta_{u_j} T_f B} \ln \left(\frac{(\alpha_{w_j})^{-\frac{4}{K} \beta_{u_j}} \psi_{u_j}}{\rho} \left(\sum_{s=0}^{u_j-1} \binom{u_j-1}{s} (-1)^s \frac{\rho}{K-u_j+1+s} - \frac{4}{K} \beta_{u_j} \frac{\alpha_{w_j}-1}{\alpha_{w_j}} \sum_{s=0}^{u_j-1} \binom{u_j-1}{s} (-1)^s e^{\frac{M-u_j+1+s}{\rho \alpha_{w_j}}} \right) \right. \\ \times E_i \left(-\frac{K-u_j+1+s}{\rho \alpha_{w_j}} \right) + \sum_{q=2}^{\infty} \left(\frac{4}{K} \beta_{u_j} \right)^q \left(\frac{\alpha_{w_j}-1}{\alpha_{w_j}} \right)^q \sum_{s=0}^{u_j-1} \binom{u_j-1}{s} (-1)^s \left(\frac{\sum_{i=1}^{q-1} \frac{(i-1)!}{(\alpha_{w_j})^{-i}} \left(-\frac{K-u_j+1+s}{\rho} \right)^{q-i-1}}{(q-1)!} \right) \\ \left. - \frac{\left(-\frac{K-u_j+1+s}{\rho} \right)^{q-1}}{(q-1)!} e^{\frac{K-u_j+1+s}{\rho \alpha_{w_j}}} E_i \left(-\frac{K-u_j+1+s}{\rho \alpha_{w_j}} \right) \right) \right), \quad (7a)$$

$$E_c^{w_j} = -\frac{1}{\theta_{w_j} T_f B} \ln \left(\frac{\psi_{w_j}}{\rho \alpha_{w_j}} \sum_{s=0}^{w_j-1} \binom{w_j-1}{s} (-1)^s U \left(1, 2 + \frac{4}{K} \beta_{w_j}, \frac{K-w_j+1+s}{\rho \alpha_{w_j}} \right) \right), \quad (7b)$$

$$E_c^{w_j} \approx -\frac{1}{\theta_{w_j} T_f B} \ln \left(\frac{\psi_{w_j}}{\rho} (\alpha_{w_j})^{\frac{4}{K} \beta_{w_j}} \sum_{s=0}^{w_j-1} \binom{w_j-1}{s} (-1)^s \left(\frac{K-w_j+1+s}{\rho} \right)^{-\frac{4}{K} \beta_{w_j}-1} \Gamma \left(\frac{4}{K} \beta_{w_j} + 1 \right) \right), \quad (7c)$$

$$\bar{E}_c^k = -\frac{1}{\theta_k T_f B} \ln \left(\frac{\psi_k}{\rho} \sum_{s=0}^{k-1} \binom{k-1}{s} (-1)^s U \left(1, 2 + \frac{2}{K} \beta_k, \frac{K-k+1+s}{\rho} \right) \right), \quad (7d)$$

$$\bar{E}_c^k \approx -\frac{1}{\theta_k T_f B} \ln \left(\frac{\psi_k}{\rho} \sum_{s=0}^{k-1} \binom{k-1}{s} (-1)^s \left(\frac{K-k+1+s}{\rho} \right)^{-\frac{2}{K} \beta_k-1} \Gamma \left(\frac{2}{K} \beta_k + 1 \right) \right). \quad (7e)$$

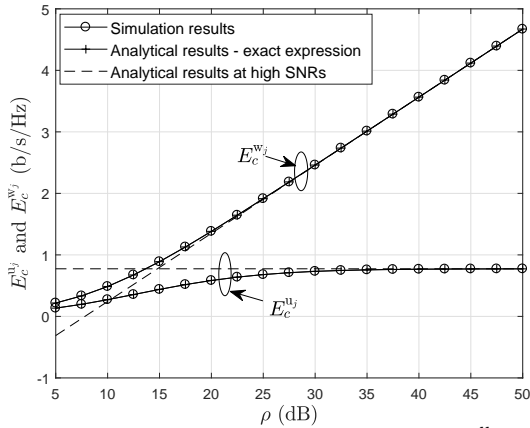


Fig. 1: Individual EC in one NOMA pair, i.e., $E_c^{u_j}$ and $E_c^{w_j}$, vs. the transmit SNR ρ .

and $\alpha_{w_j} = 0.2$. Fig. 1 shows the values of $E_c^{u_j}$ and $E_c^{w_j}$, calculated from the exact analytical closed-forms using (7a)-(7b), the closed-forms at high SNRs given in Theorem 1, and also Monte carlo results. From Fig. 1, we can conclude that the exact analytical closed-forms are accurate for $E_c^{u_j}$ and $E_c^{w_j}$ in one pair. At high SNRs, the proposed approximations for $E_c^{u_j}$ and $E_c^{w_j}$ match with the calculated exact values. Furthermore, Fig. 1 also indicates that the curve of $E_c^{u_j}$ for the weaker user approaches a constant in high SNR regime, while the $E_c^{w_j}$ curve for the stronger user shows a monotonically increasing

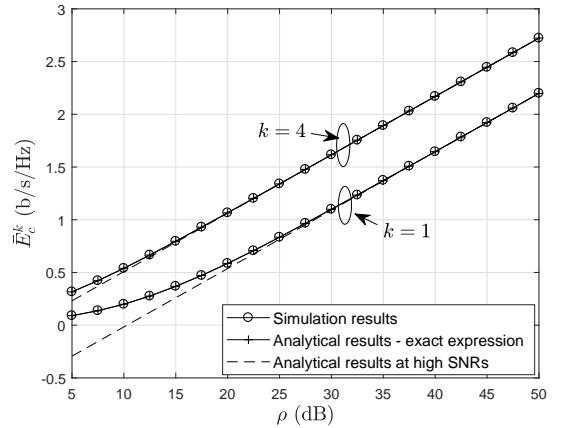


Fig. 2: Individual EC (\bar{E}_c^k) in OMA vs. the transmit SNR ρ .

trend with ρ . The theoretical proof for this phenomenon is omitted in this paper. However, we refer the interested readers to [1] which provides a theoretical proof for the trends of $E_c^{u_j}$ and $E_c^{w_j}$ versus the transmit SNR ρ .

Correspondingly, Fig. 2 is plotted which includes the curves of \bar{E}_c^k for the 1st user and the 4th user in a downlink OMA network. The results are calculated using the exact analytical closed-form given in (7d), the closed-form at high SNRs given in (7e), and also the Monte carlo results. From Fig. 2, one can notice that the derived exact analytical expression is accurate for \bar{E}_c^k in a downlink OMA network. In high SNR regime,

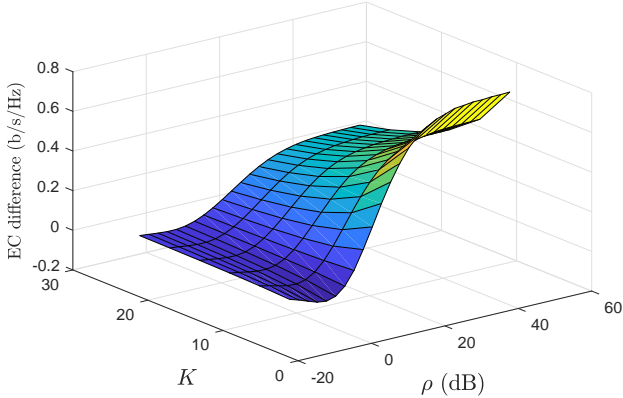


Fig. 3: $(E_c^1 + E_c^K) - (\bar{E}_c^1 + \bar{E}_c^K)$ vs. the user group size K and the transmit SNR ρ .

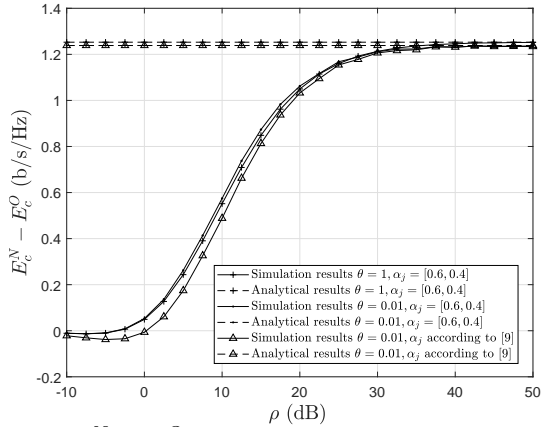


Fig. 4: $(E_c^N - E_c^O)$ vs. the transmit SNR ρ , for different power coefficient and θ settings.

the approximated closed-form matches with the exact values calculated with the Monte carlo method. Fig. 2 also shows that the 4th user always achieves higher EC values. This is because that all channels have been assumed to be ordered, hence the 4th user always has higher channel gains than the 1st user.

It was proved in [5] that the performance gain of NOMA can be improved by pairing two users with very different channel conditions. Inspired by [5], we assume that the user with the strongest channel conditions, i.e., the user K , is paired with the user with the weakest channel conditions, i.e., the user 1. Then, we plot Fig. 3 to study the total EC difference for the pair $(1, K)$, i.e., $(E_c^1 + E_c^K) - (\bar{E}_c^1 + \bar{E}_c^K)$, with respect to the user group size and the transmit SNR. From Fig. 3, it is evident that although the EC difference for the pair $(1, K)$ reduces with the group size K , it remains positive even for large group sizes. This indicates that when the user group size becomes larger, although the performance of NOMA and OMA reduces, the spectrum sharing gain by applying NOMA achieves higher total EC at high SNRs than OMA.

Assume that all 6 NOMA users are paired as follows: $\phi_1 = (1, 5)$, $\phi_2 = (2, 4)$ and $\phi_3 = (3, 6)$. All NOMA pairs have the same power coefficient settings, i.e., α_j is

the same for $j \in \{1, 2, 3\}$. Fig. 4 illustrates the curves of $E_c^N - E_c^O$ versus ρ , for different delay requirement scenarios and power coefficient settings within one pair. The solid lines are calculated with Monte carlo results and the dashed lines are plotted using the derived analytical closed-forms at high SNRs. Firstly, Fig. 4 shows that when all users have a slightly stringent delay requirement, i.e., when delay exponent θ changes from 0.01 to 1, the values of $E_c^N - E_c^O$ do not vary very much. But since θ indicates the exponential decay rate of the delay outage probability, with a slight increase of θ , the delay outage probability decreases in an exponential way. Secondly, two different power settings are included in Fig. 4 for the two users within one pair, which are a fixed power allocation $\alpha_j = (0.6, 0.4)$, $\forall j \in \{1, 2, 3\}$, and a varied power coefficient setting according to [9]. Fig. 4 shows that the two employed power settings achieve different $E_c^N - E_c^O$ values, but observe the same trend of $E_c^N - E_c^O$ versus ρ . Finally, Fig. 4 also indicates that the $E_c^N - E_c^O$ values at high SNRs converge to a constant, which confirms the theoretical analysis in Lemma 1.

V. CONCLUSIONS

In order to investigate the performance of different multiple access protocols for delay-sensitive applications, we focused on the downlink transmission which utilizes either user-paired NOMA principle or conventional OMA technique. It was proved that the performance gain of NOMA over OMA on the total delay-constrained rate converges at very high SNRs. This means that the further increase in SNR will lead to diminishing returns in the total delay-constrained rate. Further, for a special user pair with the two users having the most distinctive channel conditions, it was shown that at high SNRs, NOMA always achieves higher total EC even for a very large group size.

VI. APPENDIX

Proof for Theorem 1: In high SNR regime, $E_c^{u_j}$ for the weaker user in the j th NOMA pair, given in (6a), can be approximated as follows:

$$\lim_{\rho \rightarrow \infty} E_c^{u_j} = -\frac{1}{\theta_{u_j} T_f B} \ln \left(\frac{1}{\alpha_{w_j}} \right) \frac{4}{K} \beta_{u_j} = \frac{2}{K} \log_2 \frac{1}{\alpha_{w_j}}. \quad (11)$$

At high SNRs, $E_c^{w_j}$ for the stronger user in the j th NOMA pair, given in (6b), can be approximated and written as

$$E_c^{w_j} \approx -\frac{1}{\theta_{w_j} T_f B} \ln \left(\mathbb{E} \left[(\rho \alpha_{w_j} |h_{w_j}|^2)^{\frac{4}{K} \beta_{w_j}} \right] \right). \quad (12)$$

By setting $x_{w_j} = \rho |h_{w_j}|^2$ and inserting $f(x_{w_j}) = \frac{1}{\rho} e^{-\frac{x_{w_j}}{\rho}}$,

$$F(x_{w_j}) = 1 - e^{-\frac{x_{w_j}}{\rho}} \quad \text{into (5), (12) can be expanded as}$$

$$E_c^{w_j} \approx -\frac{1}{\theta_{w_j} T_f B} \ln \left(\frac{\psi_{w_j}}{\rho} (\alpha_{w_j}) \frac{4}{K} \beta_{w_j} \int_0^\infty (x_{w_j}) \frac{4}{K} \beta_{w_j} \right. \\ \left. \times e^{-\frac{(K-w_j+1)x_{w_j}}{\rho}} \left(1 - e^{-\frac{x_{w_j}}{\rho}} \right)^{w_j-1} dx_{w_j} \right) \quad (13)$$

$$= -\frac{1}{\theta_{w_j} T_f B} \ln \left(\frac{\psi_{w_j}}{\rho} (\alpha_{w_j}) \frac{4}{K} \beta_{w_j} \sum_{s=0}^{w_j-1} \binom{w_j-1}{s} (-1)^s \right. \\ \left. \times \int_0^\infty (x_{w_j}) \frac{4}{K} \beta_{w_j} e^{-\frac{(K-w_j+1+s)x_{w_j}}{\rho}} dx_{w_j} \right) \quad (14)$$

$$= -\frac{1}{\theta_{w_j} T_f B} \ln \left(\frac{\psi_{w_j}}{\rho} (\alpha_{w_j}) \frac{4}{K} \beta_{w_j} \sum_{s=0}^{w_j-1} \binom{w_j-1}{s} (-1)^s \right. \\ \left. \times \left(\frac{K-w_j+1+s}{\rho} \right)^{-\frac{4}{K} \beta_{w_j-1}} \Gamma \left(\frac{4}{K} \beta_{w_j} + 1 \right) \right). \quad (15)$$

From (13) to (14), it is obtained by replacing

$\left(1 - e^{-\frac{x_{w_j}}{\rho}}\right)^{w_j-1}$ with its binomial expansion

$\sum_{s=0}^{w_j-1} \binom{w_j-1}{s} (-1)^s e^{-\frac{x_{w_j}}{\rho} s}$. Then, from (3.382.4) in [18], we have that

$$\int_0^\infty (x + \beta)^v e^{-ux} dx = u^{-v-1} e^{\beta u} \Gamma(v+1, \beta u) \\ \text{for } |\arg \beta| < \pi, \operatorname{Re} u > 0. \quad (16)$$

By applying (16) to (14) and converting the incomplete gamma function $\Gamma\left(\frac{4}{K}\beta_{w_j} + 1, 0\right)$ to $\Gamma\left(\frac{4}{K}\beta_{w_j} + 1\right)$, (15) can be finally derived.

Proof for Theorem 2: At high SNRs, \bar{E}_c^k for the k^{th} user in a conventional OMA network, given in (8), can be approximated and written as

$$\bar{E}_c^k \approx -\frac{1}{\theta_k T_f B} \ln \left(\mathbb{E} \left[(\rho |h_k|^2) \frac{2}{K} \beta_k \right] \right). \quad (17)$$

By setting $y_k = \rho |h_k|^2$ and inserting the PDF $f_{(k)}(y_k)$, (17) can be expanded as

$$\bar{E}_c^k \approx -\frac{1}{\theta_k T_f B} \ln \left(\frac{\psi_k}{\rho} \int_0^\infty (y_k) \frac{2}{K} \beta_k e^{-\frac{(K-k+1)y_k}{\rho}} \right. \\ \left. \times \left(1 - e^{-\frac{y_k}{\rho}}\right)^{k-1} dy_k \right) \quad (18) \\ = -\frac{1}{\theta_k T_f B} \ln \left(\frac{\psi_k}{\rho} \sum_{s=0}^{k-1} \binom{k-1}{s} (-1)^s \right. \\ \left. \times \left(\frac{K-k+1+s}{\rho} \right)^{-\frac{2}{K} \beta_k-1} \Gamma \left(\frac{2}{K} \beta_k + 1 \right) \right). \quad (19)$$

From (18) to (19), it is achieved by first substituting

$\left(1 - e^{-\frac{y_k}{\rho}}\right)^{k-1}$ with its binomial expansion and then ap-

plying (16).

VII. ACKNOWLEDGEMENT

This work was partly funded by grant N62909-17-1-2114 from the US Office of Naval Research Global, the Royal Society project IEC170324, the EPSRC project EP/K011693/1, and EU FP7 CROWN project PIRSES-GA-2013-610524. In addition, the authors would also like to acknowledge the support of the University of Surrey 5GIC (<http://www.surrey.ac.uk/5gic>) members for this work.

REFERENCES

- [1] W. Yu, L. Musavian, and Q. Ni, "Link-layer capacity of NOMA under statistical delay QoS guarantees," *IEEE Trans. Commun.*, in press.
- [2] C. She, C. Yang, and T. Q. S. Quek, "Cross-layer optimization for ultra-reliable and low-latency radio access networks," *IEEE Trans. Wirel. Commun.*, vol. 17, no. 1, pp. 127–141, Jan. 2018.
- [3] W. Yu, L. Musavian, and Q. Ni, "Tradeoff analysis and joint optimization of link-layer energy efficiency and effective capacity toward green communications," *IEEE Trans. Wirel. Commun.*, vol. 15, no. 5, pp. 3339–3353, Jan. 2016.
- [4] Y. Liu, Z. Qin, M. ElKashlan, 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.
- [5] Z. Ding, P. Fan, and H. V. Poor, "Impact of user pairing on 5G nonorthogonal multiple-access downlink transmissions," *IEEE Trans. Veh. Technol.*, vol. 65, no. 8, pp. 6010–6023, Aug. 2016.
- [6] P. Xu and K. Cumanan, "Optimal power allocation scheme for non-orthogonal multiple access with α -fairness," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 10, pp. 2357–2369, Oct. 2017.
- [7] L. Lv, J. Chen, Q. Ni, and Z. Ding, "Design of cooperative non-orthogonal multicast cognitive multiple access for 5G systems: User scheduling and performance analysis," *IEEE Trans. Commun.*, vol. 65, no. 6, pp. 2641–2656, Jun. 2017.
- [8] Y. Liu, Z. Qin, M. ElKashlan, Y. Gao, and L. Hanzo, "Enhancing the physical layer security of non-orthogonal multiple access in large-scale networks," *IEEE Trans. Wirel. Commun.*, vol. 16, no. 3, pp. 1656–1672, Jan. 2017.
- [9] Z. Ding, Z. Yang, P. Fan, and H. V. Poor, "On the performance of non-orthogonal multiple access in 5G systems with randomly deployed users," *IEEE Signal Process. Lett.*, vol. 21, no. 12, pp. 1501–1505, Dec. 2014.
- [10] M. Gursory, D. Qiao, and S. Velipasalar, "Analysis of energy efficiency in fading channels under QoS constraints," *IEEE Trans. Wirel. Commun.*, vol. 8, no. 8, pp. 4252–4263, Aug. 2009.
- [11] J. Choi, "Effective capacity of NOMA and a suboptimal power control policy with delay QoS," *IEEE Trans. Commun.*, vol. 65, no. 4, pp. 1849–1858, Jan. 2017.
- [12] G. Liu, Z. Ma, X. Chen, Z. Ding, F. R. Yu, and P. Fan, "Cross-layer power allocation in nonorthogonal multiple access systems for statistical qos provisioning," *IEEE Trans. Veh. Technol.*, vol. 66, no. 12, pp. 11 388–11 393, Dec. 2017.
- [13] C. S. Chang, "Stability, queue length and delay, part ii: Stochastic queueing networks," in *Proc. IEEE Conf. Decision Contr.*, vol. 1, Tucson, Arizona, USA, Dec. 1992, pp. 1005–1010.
- [14] J. T. Lewis and R. Russell, *An Introduction to Large Deviations for Teletraffic Engineers*, 1997.
- [15] J. A. Bucklew, *Introduction to Rare Event Simulation*. Springer-Verlag New York Inc., 2004.
- [16] H. A. David and H. N. Nagaraja, *Order Statistics*. John Wiley, New York, 3rd ed., 2003.
- [17] M. Abramowitz and I. A. Stegun, *Handbook of Mathematical Functions*. New York: Dover, 1965.
- [18] I. S. Gradshteyn and I. M. Ryzhik, *Table of Integrals, Series, and Products*. New York: Academic Press, 6th ed., 2000.