

# Network Slicing for URLLC and eMBB with Max-Matching Diversity Channel Allocation

Elço João dos Santos Jr, Richard Demo Souza, *Senior Member, IEEE*,  
João Luiz Rebelatto, *Senior Member, IEEE*, and Hirley Alves, *Member, IEEE*

## Abstract

This work considers the problem of radio resource sharing between enhanced mobile broadband (eMBB) and ultra-reliable and low latency communications (URLLC), two heterogeneous 5G services. More specifically, we propose the use of a max-matching diversity (MMD) algorithm to properly allocate the channels to the eMBB users, considering both heterogeneous orthogonal multiple access (H-OMA) and heterogeneous non-orthogonal multiple access (H-NOMA) network slicing strategies. Our results indicate that MMD can simultaneously improve the eMBB achievable rate and the URLLC reliability regardless the network slicing strategy adopted.

## Index Terms

5G, eMBB, network slicing, URLLC, channel allocation, non-orthogonal multiple access.

## I. INTRODUCTION

5G technology aims at three heterogeneous use cases: enhanced mobile broadband (eMBB), ultra-reliable and low latency communications (URLLC) and massive machine type communications (mMTC) [1], [2]. The performance target of eMBB is to achieve high data rates with packet error rates (PER) around  $10^{-3}$ , while such kind of traffic is stable and tolerates a certain amount of latency. URLLC is an innovative service supported by 5G, which aims at much lower PER

Elço João dos Santos Jr. and Richard Demo Souza, EEL, UFSC, Florianópolis-SC, Brazil. e.joaocr@gmail.com, richard.demo@ufsc.br

João Luiz Rebelatto, CPGEI, UTFPR, Curitiba-PR, Brazil. jlrebelatto@utfpr.edu.br

Hirley Alves, CWC, University of Oulu, Oulu, Finland. hirley.alves@oulu.fi

This work has been supported in Brazil by CNPq and Print CAPES-UFSC “Automation 4.0”, and in Finland by Academy of Finland 6Genesis Flagship (Grant no318927) and EE-IoT (no319008).

with strict latency constraints. The required data rate of URLLC is typically low compared to eMBB, while traffic is intermittent. Finally, mMTC must provide connectivity to a large number of devices, sporadically transmitting short packets, where the target PER is much larger than for eMBB and URLLC.

The diverse requirements of the different 5G use cases pose great challenges for the network design. One of the enablers of 5G is network slicing [3], a service-oriented point of view to model a system that can support multiple types of users in a common physical infrastructure, which adapts the network shape and allocates resources according to the services in use [3]. The original slicing concept assumes an orthogonal approach, where the network resources of one service are isolated from others. Recently, focusing on the physical layer (PHY) and in the uplink, Popovski *et al* discuss heterogeneous orthogonal multiple access (H-OMA) and introduce the paradigm of heterogeneous non-orthogonal multiple access (H-NOMA), where the term heterogeneous refers to the heterogeneity of services (eMBB, URLLC, mMTC) [2]. In H-OMA the slicing of PHY resources is orthogonal among the services, while in H-NOMA heterogeneous services may share different PHY slices, making use of successive interference cancellation (SIC) decoding. Such concept is remarkably different than regular NOMA, in which radio resources are shared by devices with the same requirements [4].

By means of a communication theoretic approach, the trade-offs between H-OMA and H-NOMA in the slicing of eMBB and URLLC, as well as of eMBB and mMTC, have been analyzed in [2]. For instance, results indicate that H-NOMA slicing between eMBB and URLLC can achieve significant improvements in the system performance, a consequence of exploiting the concept of reliability diversity. Since URLLC traffic is much more reliable, and therefore has a high probability of being successfully decoded under the interference of eMBB traffic, sharing eMBB resources with URLLC brings considerable performance gains specially at high eMBB target data rates. Nevertheless, depending on the channel conditions and on the URLLC target data rate, H-OMA can be a better solution than H-NOMA [2]. Interesting trade-offs also arise in the case of eMBB and mMTC slicing, bringing new optimization opportunities for the system designer.

In [5], the H-OMA and H-NOMA strategies for eMBB and URLLC slicing from [2] are applied to a multi-cell scenario, where the URLLC traffic is decoded at the base station (BS) to meet latency requirements, while eMBB traffic is forwarded to a cloud server. The H-NOMA approach leads to improvements for both eMBB and URLLC when the activation probability

of the last is small and the edge has sufficient capacity. The authors of [6] recently investigated the slicing between URLLC and eMBB considering a minimum mean square error receiver with multiple antennas, comparing the performance with and without SIC. Results reveal that H-NOMA with SIC brings improvements in high SNR or with low URLLC loads, while H-OMA can achieve higher URLLC reliability.

It is interesting to note that, as eMBB traffic is scheduled only after a radio access and contention resolution phase, it is practical to consider that channel state information (CSI) for eMBB users is known at the BS, as implemented in current wireless standards [7], and assumed in [2]. As a consequence, supposing the use of orthogonal frequency division multiple access (OFDMA), different eMBB users can be adequately allocated to the frequency resources as to maximize the diversity and meet the reliability target. In this sense, in [8], a maximum-matching diversity (MMD) method based on random bipartite graph theory is considered for allocating users to channels in an OMA scenario with homogeneous types of users. Through MMD, frequency diversity for each user is maximized, equaling the number of independent channels.

In this work, we modify the communication theoretic analysis of [2] in order to investigate the impact of eMBB channel allocation in the performance of H-OMA and H-NOMA slicing between eMBB and URLLC uplink traffics. The main contributions of this work are summarized as follows:

- We evaluate the achievable rate of eMBB when adopting the MMD method from [8] to allocate the channels;
- We evaluate the performance of the network slicing between URLLC and eMBB with channel allocation. The increased frequency diversity achieved by the proposed MMD-aided scheme is beneficial to both H-OMA and H-NOMA, being capable of improving the eMBB achievable rate and the URLLC reliability simultaneously.

The rest of this paper is organized as follows. Section II presents the system model, while Section III brings the communication theoretic performance formulation of H-OMA and H-NOMA between eMBB and URLLC traffics, considering MMD channel allocation. Numerical results are discussed in Section IV, while Section V concludes the paper.

## II. SYSTEM MODEL

We consider the uplink of a network with eMBB and URLLC devices transmitting to a common BS. The bandwidth is divided into  $F$  channels of index  $f \in \{1, \dots, F\}$ , where each channel is subject to independent and identically distributed (i.i.d.) Rayleigh fading, which is assumed to be constant during one time slot (TS). The channel coefficient of user  $i \in \{B, U\}$  in channel  $f$  is thus  $H_{i,f} \sim \mathcal{CN}(0, \bar{\gamma}_i)$ , where  $\bar{\gamma}_i$  corresponds to the average signal-to-noise ratio (SNR), being  $G_{i,f} \triangleq |H_{i,f}|^2$  the channel gain, and where subscripts  $B$  and  $U$  refer to eMBB and URLLC devices, respectively. The number of channels allocated to user  $i$  is  $F_i \leq F$ , with  $i \in \{B, U\}$ . Moreover, each TS is divided into  $S$  minislots.

In accordance to [2], we adopt the following approaches regarding eMBB and URLLC transmissions:

- An URLLC device transmits in a pre-assigned minislot<sup>1</sup> (to meet latency requirements), in grant-free fashion, and spreads the transmission over  $F_U$  channels (to increase reliability). The activation probability of the device is  $a_U$ .
- An eMBB user transmits in a single channel  $f$  among the  $F_B$  available channels, but during the entire TS.

This time-frequency grid is illustrated in Fig. 1, considering H-OMA in Fig. 1(a) and H-NOMA in Fig. 1(b). It is worthy mentioning that, as in [2], we do not aim at evaluating the influence of the sharing of wireless resources among devices of the same type. Thus, we assume that radio access and competition among eMBB devices have been resolved prior to the considered time slot, *i.e.*, the number of eMBB devices able to transmit in such time slot is equal to the number of channels  $F_B$ . We also do not model collisions among URLLC devices, by assuming that a single URLLC device is active in a given pre-assigned minislot with some probability  $a_U$ .

Moreover, for eMBB we assume that the BS has CSI before transmission, obtained during a scheduling phase, as in [2]. However, differently from [2], in this work we consider MMD channel assignment for eMBB users as follows.

<sup>1</sup>One could relax this restriction and allow a URLLC device to use more than one minislot, establishing a trade-off between delay and performance. However, this is beyond the scope of this paper and is left as future work.

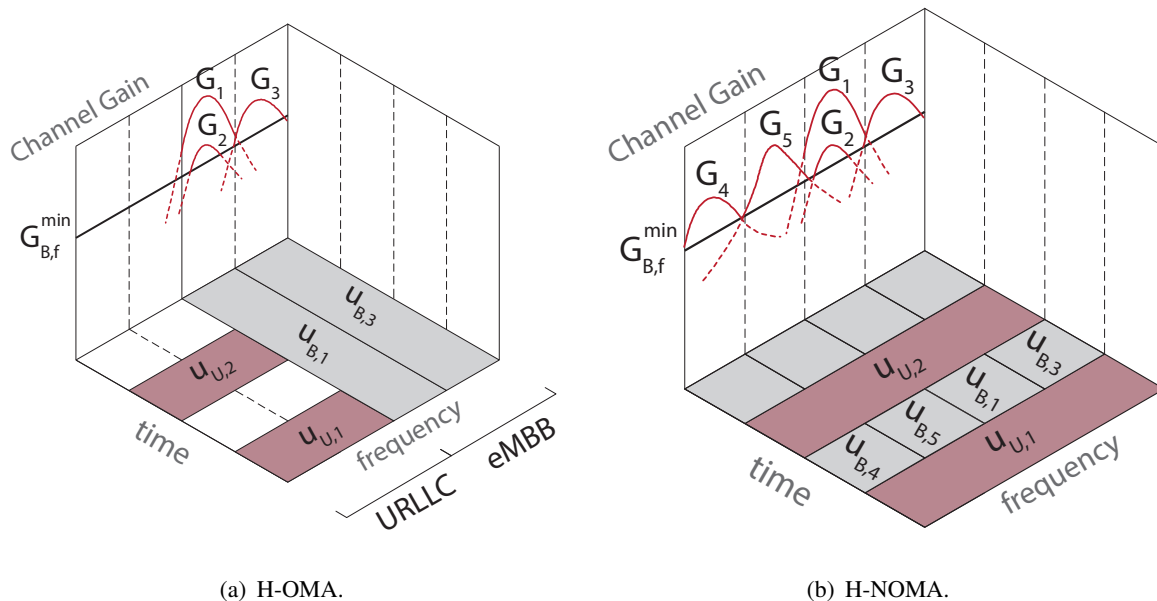


Fig. 1: System model with  $F = 4$  channels and  $S = 4$  subslots.

#### A. Max-Matching Diversity (MMD) Channel Allocation

The MMD channel allocation from [8] increases the frequency diversity by means of a proper allocation of users to channels. Through the random graph theory, herein eMBB users and channels are seen as opposite parts of the vertex set of a bipartite graph, while a user is connected to a channel if it is not in outage. The MMD algorithm then aims at minimizing the number of users in outage, by means of channel allocation.

The outage probability of the MMD scheme, in a scenario with  $F$  channels subjected to independent fading and with average SNR  $\bar{\gamma}$  per channel, is [8]

$$\mathcal{P}_{\text{MMD}}(\bar{\gamma}, F) = 2\mathcal{P}_s(\bar{\gamma})^F + O(\mathcal{P}_s(\bar{\gamma})^F) \approx 2\mathcal{P}_s(\bar{\gamma})^F, \quad (1)$$

where  $\mathcal{P}_s(\bar{\gamma})$  is the outage probability of a single channel and  $O(\cdot)$  is the higher order infinitesimal. Note in (1) that the MMD scheme from [8] achieves optimal frequency diversity, which equals the number of independent channels.

### III. SLICING FOR URLLC AND eMBB WITH MMD

#### A. eMBB with MMD

Let us consider a scenario where a radio resource  $f \in \{1, \dots, F\}$  is allocated exclusively to an eMBB user, following the MMD algorithm from [8]. As for eMBB it is reasonable to

assume CSI before data transmission, then transmit power can be adapted. The main objective of eMBB is to maximize its data rate, subject to the reliability requirement  $\epsilon_B$  and the average power constraint  $P_B = 1$ .

**Theorem 1.** *The MMD-aided eMBB rate is*

$$r_B^{\text{MMD}} = \log_2 (1 + G_{B,f}^{\text{tar}}), \quad [\text{bits/symbol}] \quad (2)$$

where

$$G_{B,f}^{\text{tar}} = \frac{\bar{\gamma}'_B}{\sum_{k=1}^{F_B} (-1)^{k-1} \binom{F_B}{k} k \Gamma \left( 0, \frac{k G_{B,f}^{\text{min}}}{\bar{\gamma}'_B} \right)}, \quad (3)$$

$$G_{B,f}^{\text{min}} = -\bar{\gamma}'_B \ln \left( 1 - \epsilon_B^{1/F_B} \right) \text{ and } \bar{\gamma}'_B \triangleq 2^{-(1/F_B)} \bar{\gamma}_B.$$

*Proof:* Please refer to Appendix A. □

## B. URLLC

Differently from eMBB users, it is assumed that URLLC devices do not have CSI, due to latency requirements. The URLLC device transmits data in all the  $F_U$  i.i.d. channels of a minislot, such that the outage probability, in the absence of interference from other services, is<sup>2</sup> [2]

$$\mathcal{P}_U(G_{U,f}) = \Pr \left( \frac{1}{F_U} \sum_{f=1}^{F_U} \log_2(1 + G_{U,f}) < r_U \right). \quad (4)$$

The target rate  $r_U$  is obtained by imposing the requirement  $\mathcal{P}_U(G_{U,f}) \leq \epsilon_U$  to (4), which we refer to as  $r_U^{\text{OMA}}$ .

## C. H-OMA

In H-OMA, orthogonal slicing is achieved by allocating  $F_B \leq F$  channels exclusively to the eMBB users, while the remaining  $F_U = F - F_B$  are allocated to URLLC. Thus, the rate region  $(r_B^{\text{OMA-MMD}}, r_U^{\text{OMA}})$  is obtained by considering  $r_B^{\text{OMA-MMD}}$  as the sum-rate of the active eMBB users,

$$r_B^{\text{OMA-MMD}} = F_B r_B^{\text{MMD}}, \quad (5)$$

where  $r_B^{\text{MMD}}$  comes from (2) and  $r_U^{\text{OMA}}$  is computed from (4).

<sup>2</sup>We follow [2] and assume that the block length utilized in the URLLC protocol is long enough so that the finite block length formulation can be well approximated by the asymptotic outage formulation [9].

#### D. H-NOMA

In H-NOMA, all the  $F$  channels are available simultaneously for both eMBB and URLLC ( $F_B = F_U = F$ ), such that the interference from URLLC transmissions into eMBB (and vice-versa) needs to be considered. To this end, we consider that the BS performs SIC<sup>3</sup>. Due to latency and reliability constraints, we assume that the SIC decoder always attempts to decode the URLLC transmission first, while treating the eMBB traffic as interference. In case of a successful decoding, the URLLC signal is then removed from the superimposed signal before attempting to decode the eMBB traffic. As a consequence, URLLC would only interfere with eMBB when the SIC decoder is not able to decode the URLLC messages.

In this scenario, we have the following Lemma regarding the achievable rate of the eMBB device.

**Lemma 1.** *The achievable rate of a MMD-aided eMBB device in H-NOMA is*

$$r_B^{NOMA-MMD} = F \log_2 (1 + G_{B,f}^{tar}), \quad (6)$$

$G_{B,f}^{tar}$  is upper bounded by (3) and the threshold SNR is

$$G_{B,f}^{min} \leq -\tilde{\gamma}'_B \ln \left( \frac{1 - \epsilon_B^{1/F_B}}{1 - \epsilon_U(1 - (1 - a_U)^S)} \right). \quad (7)$$

*Proof:* Following [2], the eMBB outage probability under H-NOMA can be bounded by the law of total probability as

$$\mathcal{P}_s(\tilde{\gamma}_B) = 1 - e^{-G_{B,f}^{min}/\tilde{\gamma}_B} \leq \frac{1 - \epsilon_B}{1 - \epsilon_U(1 - (1 - a_U)^S)}, \quad (8)$$

which accounts for the fact that the eMBB user is in outage when the SIC decoder does not decode the URLLC signal. Hence, by resorting to the fact that  $\mathcal{P}_B = \mathcal{P}_s(\tilde{\gamma}'_B)^{F_B} = \epsilon_B$  due to MMD channel allocation, one can isolate the threshold SNR from (8) as presented in (7), concluding the proof.  $\square$

The threshold from (7) indicates that the impact of URLLC transmissions in the eMBB decoding should be minimal, due to the fact that, by definition,  $\epsilon_U \ll \epsilon_B$ . On the other hand,

<sup>3</sup>Note that, as presented in [2], SIC outperforms other techniques of multi-user detection, such as puncturing and erasure decoders.

the eMBB interference in the URLLC traffic is supposed to be more critical, since URLLC is decoded prior to eMBB. As in [2] the outage probability of URLLC under H-NOMA is

$$\mathcal{P}_U^{\text{NOMA}}(G_{U,f}) = \Pr \left( \frac{1}{F_U} \sum_{f=1}^{F_U} \log_2 \left( 1 + \frac{G_{U,f}}{G_{B,f}^{\text{tar}}} \right) < r_U \right), \quad (9)$$

where we assume that the interference of eMBB is always present in the URLLC decoding, due to their long period activation. The URLLC achievable rate  $r_U^{\text{NOMA}}$  is then obtained by imposing the reliability constraint  $\mathcal{P}_U^{\text{NOMA}}(G_{U,f}) \leq \epsilon_U$ .

For comparison purposes, the eMBB achievable rate from [2], without channel allocation, can be obtained by setting  $F_B = 1$  in (3) (OMA) and in (7) (NOMA).

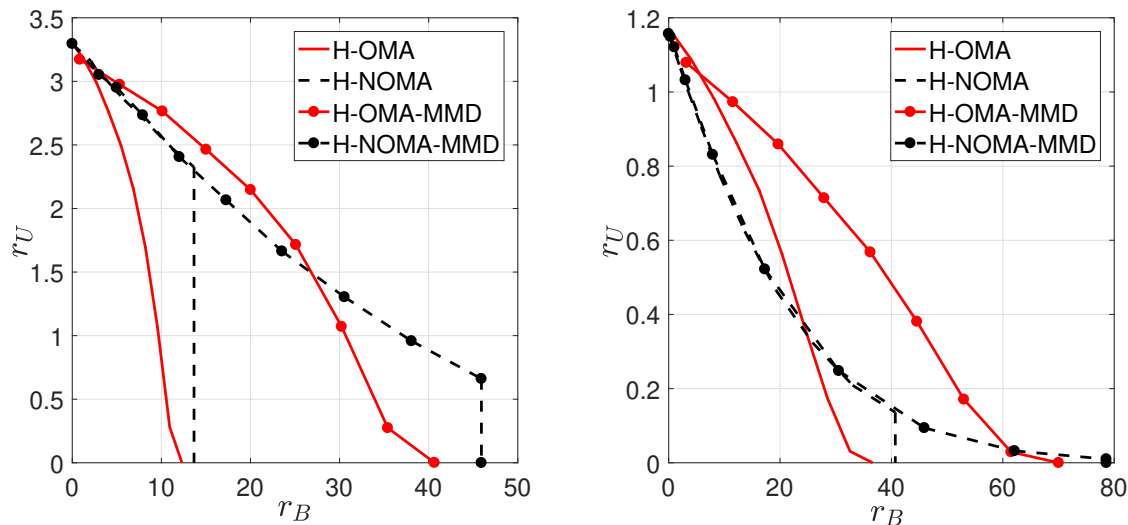
#### IV. NUMERICAL RESULTS

We resort to Monte Carlo simulations to evaluate the performance of the MMD-aided proposed schemes. Figs. 2(a) and 2(b) present the rate region  $(r_B, r_U)$ , respectively for the cases  $\bar{\gamma}_U > \bar{\gamma}_B$  and  $\bar{\gamma}_B > \bar{\gamma}_U$ . As expected, H-NOMA tends to achieve higher rates when  $\bar{\gamma}_U > \bar{\gamma}_B$  (due to the increased probability of recovering the URLLC in the SIC process), while H-OMA may be more advantageous when  $\bar{\gamma}_B > \bar{\gamma}_U$ . The beneficial impact of MMD is more evident in Fig. 2(a), where its additional frequency diversity compensates the worst average channel condition of eMBB. In this scenario, the maximum  $r_B^{\text{NOMA}}$  goes from  $\approx 13.7$  bits/symbol to  $r_B^{\text{NOMA-MMD}} \approx 46$  bits/symbol. Significant improvements are also noticed in Fig. 2(b), mainly to OMA.

Tab. I evaluates  $r_B$  for  $\epsilon_U \in \{10^{-5}, 10^{-6}, 10^{-7}\}$ ,  $r_U \in \{1, 2\}$  bits/symbol,  $\bar{\gamma}_B = 10$  dB, and  $\bar{\gamma}_U = 20$  dB. One can see that the MMD-aided schemes achieve larger  $r_B$ , while operating at more stringent values of  $\epsilon_U$ , simultaneously improving the rate of eMBB and reliability of URLLC. For example, while  $r_B^{\text{OMA-MMD}} = 17$  bits/symbol for  $\epsilon = 10^{-6}$ , the rate achieved at a higher target outage probability  $\epsilon = 10^{-5}$  is only  $r_B^{\text{OMA}} = 7.5$  bits/symbol, for a fixed  $r_U = 2$  bits/symbol. This shows that the gains provided by MMD do not benefit only the eMBB user: the URLLC user is indirectly benefited from the eMBB channel allocation, since eMBB can operate at lower transmission power levels, achieving the same (or even improved) performance, while decreasing the interference on URLLC devices. Moreover, in general, NOMA-MMD outperforms OMA-MMD for small values of  $r_U$ , and OMA-MMD becomes more advantageous when  $r_U$  increases.

Finally, since in practice eMBB CSI is already available at the BS for scheduling purposes, the main drawback of the proposed scheme resides on the complexity increase to execute the



(a)  $\bar{\gamma}_B = 10$  dB,  $\bar{\gamma}_U = 20$  dB.(b)  $\bar{\gamma}_B = 20$  dB,  $\bar{\gamma}_U = 10$  dB.Fig. 2:  $S=5$ ,  $a_U=0.1$ ,  $F=10$ ,  $\epsilon_B=10^{-3}$ ,  $\epsilon_U=10^{-5}$ .TABLE I:  $r_B$  vs  $\epsilon_U$ , for  $r_U \in \{1, 2\}$  bits/symbol.

$\epsilon_U$	$10^{-5}$		$10^{-6}$		$10^{-7}$	
	1	2	1	2	1	2
$r_B^{\text{OMA}}$	9.6	<b>7.5</b>	8.8	6.1	7.3	4.0
$r_B^{\text{NOMA}}$	13.7	13.7	13.7	13.7	13.7	10.0
$r_B^{\text{OMA-MMD}}$	30.8	21.5	26.8	<b>17.0</b>	25.2	9.0
$r_B^{\text{NOMA-MMD}}$	36.9	18.8	34.1	15.0	30.0	10.0

channel allocation, which scales with  $\mathcal{O}(F_B^{2.5})$  [8], and a very slight increase (few bits) in control traffic to eMBB users.

## V. FINAL COMMENTS

We considered network radio resource slicing between eMBB and URLLC users in a 5G system. By resorting to the MMD approach to properly allocate channels to the eMBB users, we showed that the eMBB achievable rate and the URLLC reliability can be improved simultaneously, under both H-OMA and H-NOMA network slicing strategies.

APPENDIX A  
PROOF OF THEOREM 1

Assuming that the eMBB devices have CSI [2], a transmission only occurs when  $G_{B,f}$ , the instantaneous channel gain, is greater than a threshold SNR  $G_{B,f}^{\min}$ . The outage probability of a point-to-point (single channel) communication is then

$$\begin{aligned} \mathcal{P}_s(\bar{\gamma}_B) &= \Pr[G_{B,f} < G_{B,f}^{\min}] = \int_0^{G_{B,f}^{\min}} p_{G_{B,f}}(x) dx \\ &= \int_0^{G_{B,f}^{\min}} \frac{e^{-x/\bar{\gamma}_B}}{\bar{\gamma}_B} dx = 1 - e^{-G_{B,f}^{\min}/\bar{\gamma}_B}. \end{aligned} \quad (10)$$

When considering MMD channel allocation from [8] with  $F_B$  independent channels, the outage probability of an eMBB user, in a scenario where the number of users is equal to  $F_B$ , can be approximated following (1) as

$$\mathcal{P}_B \approx 2\mathcal{P}_s(\bar{\gamma}_B)^{F_B}. \quad (11)$$

For Rayleigh fading, it can be shown from (10) that  $2\mathcal{P}_s(\bar{\gamma}_B)^{F_B} \approx \mathcal{P}_s(\bar{\gamma}'_B)^{F_B}$ , with  $\bar{\gamma}'_B \triangleq 2^{-(1/F_B)}\bar{\gamma}_B$ . After imposing the reliability constraint  $\mathcal{P}_B = \epsilon_B$ , we obtain the threshold SNR from (10) and (11) as

$$G_{B,f}^{\min} = -\bar{\gamma}'_B \ln \left( 1 - \epsilon_B^{1/F_B} \right). \quad (12)$$

Based on power inversion, the instantaneous power is chosen as  $P_B(G_{B,f}^{\text{MMD}}) = G_{B,f}^{\text{tar}}/G_{B,f}^{\text{MMD}}$  when  $G_{B,f}^{\text{MMD}} \geq G_{B,f}^{\min}$ . Otherwise,  $P_B(G_{B,f}^{\text{MMD}})$  is set to zero. The target SNR  $G_{B,f}^{\text{tar}}$  is then obtained by imposing the average power constraint

$$1 = \mathbb{E} [P_B(G_{B,f}^{\text{MMD}})] = \int_{G_{B,f}^{\min}}^{\infty} p_{G_{B,f}^{\text{MMD}}}(x) P_B(x) dx, \quad (13)$$

where  $p_{G_{B,f}^{\text{MMD}}}(x)$  is the probability density function (pdf) of the instantaneous channel gain after MMD,  $G_{B,f}^{\text{MMD}}$ , obtained from the cumulative density function (cdf)  $\mathcal{P}_s(\bar{\gamma}'_B)^{F_B}$  as

$$\begin{aligned} p_{G_{B,f}^{\text{MMD}}}(x) &= \frac{d}{dt} \left[ (1 - e^{-G_{B,f}^{\min}/\bar{\gamma}'_B})^{F_B} \right] \\ &= \frac{F_B}{\bar{\gamma}'_B} \left[ 1 - e^{-x/\bar{\gamma}'_B} \right]^{F_B-1} e^{-x/\bar{\gamma}'_B} \\ &= \sum_{k=1}^{F_B} (-1)^{k-1} \binom{F_B}{k} \frac{k e^{-kx/\bar{\gamma}'_B}}{\bar{\gamma}'_B}. \end{aligned} \quad (14)$$

The summation in (14) is obtained applying the binomial theorem [10], so that the pdf of the SNR of  $F_B$  channels can be expressed as a linear combination of  $F_B$  exponential pdfs. After replacing (14) in (13), we have:

$$\begin{aligned} 1 &= \int_{G_{B,f}^{\min}}^{\infty} \sum_{k=1}^{F_B} (-1)^{k-1} \binom{F_B}{k} \frac{k e^{-kx/\bar{\gamma}'_B} G_{B,f}^{\text{tar}}}{\bar{\gamma}'_B x} dx \\ &= G_{B,f}^{\text{tar}} \sum_{k=1}^{F_B} (-1)^{k-1} \binom{F_B}{k} \frac{k}{\bar{\gamma}'_B} \Gamma \left( 0, \frac{k G_{B,f}^{\min}}{\bar{\gamma}'_B} \right), \end{aligned} \quad (15)$$

where  $\Gamma(\cdot, \cdot)$  is the upper incomplete gamma function. Then, one can isolate  $G_{B,f}^{\text{tar}}$  from (15), concluding the proof.

## REFERENCES

- [1] M. Shafi *et al.*, “5G: A tutorial overview of standards, trials, challenges, deployment, and practice,” *IEEE J. Sel. Areas Commun.*, vol. 35, no. 6, pp. 1201–1221, June 2017.
- [2] P. Popovski, K. F. Trillingsgaard, O. Simeone, and G. Durisi, “5G wireless network slicing for eMBB, URLLC, and mMTC: A communication-theoretic view,” *IEEE Access*, vol. 6, pp. 55 765–55 779, 2018.
- [3] X. Foukas, G. Patounas, A. Elmokashfi, and M. K. Marina, “Network slicing in 5G: Survey and challenges,” *IEEE Commun. Mag.*, vol. 55, no. 5, pp. 94–100, May 2017.
- [4] Z. Ding *et al.*, “A survey on non-orthogonal multiple access for 5G networks: Research challenges and future trends,” *IEEE J. Sel. Areas Commun.*, vol. 35, no. 10, pp. 2181–2195, Oct 2017.
- [5] R. Kassab, O. Simeone, P. Popovski, and T. Islam, “Non-orthogonal multiplexing of ultra-reliable and broadband services in fog-radio architectures,” *IEEE Access*, vol. 7, pp. 13 035–13 049, 2019.
- [6] R. Abreu *et al.*, “On the multiplexing of broadband traffic and grant-free ultra-reliable communication in uplink,” in *IEEE Vehicular Technology Conference (VTC)*, 2019, pp. 1–6.
- [7] A. A. Zaidi *et al.*, “Designing for the future: the 5G NR physical layer,” Ericsson, Tech. Rep., Jul. 2017. [Online]. Available: <https://www.ericsson.com/en/ericsson-technology-review/archive/2017/designing-for-the-future-the-5g-nr-physical-layer>
- [8] B. Bai, W. Chen, Z. Cao, and K. B. Letaief, “Max-matching diversity in OFDMA systems,” *IEEE Trans. Commun.*, vol. 58, no. 4, pp. 1161–1171, April 2010.
- [9] W. Yang, G. Durisi, T. Koch, and Y. Polyanskiy, “Quasi-static multiple-antenna fading channels at finite blocklength,” *IEEE Trans. Inf. Theory*, vol. 60, no. 7, pp. 4232–4265, July 2014.
- [10] F. Rosas, R. D. Souza, M. Verhelst, and S. Pollin, “Energy-efficient MIMO multihop communications using the antenna selection scheme,” in *IEEE Intern. Symp. on Wireless Commun. Systems (ISWCS)*, 2015, pp. 686–690.