

Block Error Performance of NOMA with HARQ-CC in Finite Blocklength

Dileepa Marasinghe, Nandana Rajatheva, Matti Latva-aho

Centre for Wireless Communications

University of Oulu

Oulu, Finland

{dileepa.marasinghe, nandana.rajatheva, matti.latva-aho}@oulu.fi

Abstract—This paper investigates the performance of a two-user downlink non-orthogonal multiple access (NOMA) system using hybrid automatic repeat request with chase combining (HARQ-CC) in finite blocklength. First, an analytical framework is developed by deriving closed-form approximations for the individual average block error rate (BLER) of the near and the far user. Based upon that, the performance of NOMA is discussed in comparison to orthogonal multiple access (OMA), which draws the conclusion that NOMA outperforms OMA in terms of user fairness. Further, an algorithm is devised to determine the required blocklength and power allocation coefficients for NOMA that satisfies reliability targets for the users. The required blocklength for NOMA is compared to OMA, which shows NOMA has a lower blocklength requirement in high transmit signal-to-noise ratio (SNR) conditions, leading to lower latency than OMA when reliability requirements in terms of BLER for the two users are in the order of 10^{-5} .

Index Terms—non-orthogonal multiple access, hybrid automatic repeat request, chase combining, short packet communications, block error rate, ultra-reliable communications.

I. INTRODUCTION

With the advent of new use-cases requiring high reliability and low-latency in 5G and beyond, transmission with finite blocklength becomes inevitable to reduce latency. In contrast to classical information-theoretic principles, the use of finite blocklength results in a non-negligible decoder error probability. Hybrid automatic repeat request (HARQ) procedures are used to improve the accuracy in decoding by exploiting time-diversity at the expense of increased latency. Thus, achieving high reliability and low-latency are Pareto-optimal, which calls for a trade-off between the two. Concurrently, non-orthogonal multiple access (NOMA) has gained widespread attention in research due to the ability to outperform its counterpart, orthogonal multiple access (OMA) in terms of spectral efficiency and user fairness.

Studies on NOMA with HARQ can be found in the literature [1]–[4]. Authors in [1] show that NOMA with successive interference cancellation (SIC) employing HARQ with incremental redundancy (HARQ-IR) can outperform OMA in outage probability. In [2], a power allocation strategy for HARQ-IR with NOMA is presented. The outage performance of NOMA with the HARQ-CC scheme has been studied in [3] and [4], by deriving closed-form approximations for outage probability. Analysis of HARQ based systems using finite blocklength have been presented in [5]–[7]. Authors in [5]

investigate the blocklength which maximizes throughput and minimizes average delay while authors in [6] investigate power allocation for systems using type-I ARQ in finite blocklength. In [7], a closed-form derivation of the outage probabilities on HARQ-IR in finite blocklength is provided. A power allocation method for HARQ-CC with finite blocklength, which targets reliability constraints is proposed in [8]. Analysis of NOMA in the finite blocklength regime is reported in [9], [10]. Authors in [9] and [10] investigate the finite blocklength performance of a two-user downlink NOMA system for single and multiple antenna base station (BS) and demonstrate that having a common blocklength for both NOMA users is optimal and NOMA outperforms OMA in latency.

Improving reliability and minimizing the latency are two targets that are conflicting with each other to be achieved simultaneously. The reason is improving reliability would be supported by re-transmissions including the use of longer packets, which will then increase the latency. A trade-off between latency and reliability that fit different use-cases is required when considering ultra-reliable low-latency communications (URLLC). Some URLLC use-cases are remote surgery and factory automation which have stricter targets like 1×10^{-9} with 1 ms latency and V2X communications, and tactile internet which have reliability around 1×10^{-5} and latency requirements ranging from 1 ms to 100 ms. The motivation behind this work is to analyze the performance of a system comprising of the three enablers; NOMA combined with HARQ in finite blocklength. While NOMA allows higher spectral efficiency by utilizing the same frequency-time resource, the use of HARQ improves the reliability, and the use of short packets allows reducing latency. The main goal is to investigate the ability of NOMA to deliver ultra-reliability using HARQ combined with the use of short packets to reduce latency. The reliability is investigated by characterizing the average block error rate (BLER). Also, determining the required number of channel uses or blocklength for NOMA, which satisfies given reliability targets for the two users and comparing the performance of NOMA with OMA are presented.

The next sections are organized as follows. Section II describes the system model. In Section III analytical approximations for the average BLERs are derived for the two users. Section IV presents asymptotic BLER approximations

considering high SNR conditions and an algorithm is devised to determine the required blocklength for the system to meet given reliability constraints for the two users. In Section V numerical results are provided to validate the derived approximations in Section III and comparison of the blocklength requirement between NOMA and OMA is provided. Section VI concludes the paper. The proofs of the main results are provided in Appendices.

II. SYSTEM MODEL

Consider a downlink power domain NOMA system that uses short-packets for communications. The system comprises of a single antenna base station and two users u_1, u_2 equipped with single antennas. Without loss of generality, assume that u_1 is located close to the BS, thus having a higher channel gain, while u_2 is located far from the BS with a lower channel gain. Further, assume the channel gain of the users are known at the BS. To enhance the reliability of the transmission, the system uses the HARQ-CC scheme. The BS serves the users following the NOMA principle. Let x_1 and x_2 be the unit energy messages to u_1 and u_2 , respectively. The BS encodes these messages using the superposition coding technique with power allocation coefficients α_1 and α_2 such that $\alpha_1 + \alpha_2 = 1$ with a total power of P . According to the NOMA principle, BS allocates more power to the far user by setting $\alpha_1 < \alpha_2$ ensuring user fairness. Therefore, the transmitted signal s can be expressed as

$$s = \sqrt{\alpha_1 P} x_1 + \sqrt{\alpha_2 P} x_2. \quad (1)$$

The received signal y_i at u_i , $i = 1, 2$ in the t^{th} transmission round can be expressed as

$$y_i = \tilde{h}_{i,t} (\sqrt{\alpha_1 P} x_1 + \sqrt{\alpha_2 P} x_2) + n_i, \quad (2)$$

where $\tilde{h}_{i,t} = \frac{h_{i,t}}{\sqrt{1+d_i^\eta}}$, $h_{i,t} \sim \mathcal{CN}(0, 1)$ is the independent and identically distributed (i.i.d) fading coefficient of u_i with equal blocklength M in the t^{th} transmission round, d_i is the distance between u_i and the BS, η is the path loss exponent and n_i is the additive white Gaussian noise (AWGN) with variance σ^2 .

The far user, u_2 attempts to decode the received signal treating u_1 's signal as interference. Then the received signal-to-noise-plus-interference ratio (SINR) at u_2 for decoding its message at the t^{th} transmission round is

$$\gamma_{22}^t = \frac{\rho \alpha_2 |\tilde{h}_{2,t}|^2}{\rho \alpha_1 |\tilde{h}_{2,t}|^2 + 1}, \quad (3)$$

where ρ is the transmit SNR such that $\rho = \frac{P}{\sigma^2}$.

The near user, u_1 applies SIC in decoding the messages, which means u_1 decodes u_2 's message first and then its own message without interference. The SINRs for decoding at u_1 are given by

$$\gamma_{12}^t = \frac{\rho \alpha_2 |\tilde{h}_{1,t}|^2}{\rho \alpha_1 |\tilde{h}_{1,t}|^2 + 1} \quad \text{and} \quad \gamma_{11}^t = \rho \alpha_1 |\tilde{h}_{1,t}|^2. \quad (4)$$

In the HARQ-CC procedure, in case of a failure to decode its message, the user retains the received signal and sends a negative acknowledgement (NACK) to the BS. If a NACK is received to the BS from any of the two users, BS retransmits the same encoded signal. Users employ maximum ratio combining (MRC) for decoding by combining the received signals stored during previous rounds and the new signal received. In case of successful decoding, the user will send a positive acknowledgement (ACK). BS transmits a new signal when it receives ACKs from both users. This work assumes the feedback channel, which ACKs/NACKs are sent, to be a one-bit error-free channel. The number of transmission rounds is limited to a maximum of T . The SINR for decoding u_j 's signal at u_i where $i, j = 1, 2$ after T rounds of transmissions [3] is

$$\gamma_{ij} = \sum_{t=1}^T \gamma_{ij}^t. \quad (5)$$

III. AVERAGE BLER OF NOMA WITH HARQ-CC IN FINITE BLOCKLENGTH

A. Preliminaries

Short-packets are used in the system for achieving low-latency in communications with a finite blocklength. Based on the recent work by Polyanskiy et al. [11], the decoder error probability or the BLER of u_i for decoding u_j 's information, ϵ_{ij} in finite blocklength is given by

$$\epsilon_{ij} \approx Q \left(\frac{\log_2(1 + \gamma_{ij}) - \frac{N_j}{M}}{\sqrt{\frac{v_{i,j}}{M}}} \right) \triangleq \Phi(\gamma_{ij}, N_j, M). \quad (6)$$

where N_j is the number of information bits transferred using a blocklength of M channel uses, $\gamma_{i,j}$ is the SINR, $v_{i,j}$ is the channel dispersion defined by $v_{i,j} = (\log_2 e)^2 \left(1 - \frac{1}{(1+\gamma_{i,j})^2}\right)$, $Q(\cdot)$ is the Q function. This approximation holds when M is sufficiently large [10], such as $M \geq 100$.

The user u_1 uses SIC in decoding, so the instantaneous BLER depends on the two stages in the SIC procedure. The success of the first stage affects the BLER in decoding at the second stage. Therefore, the instantaneous BLER for u_1 is given by

$$\epsilon_1 = \epsilon_{12} + (1 - \epsilon_{12})\epsilon_{11}. \quad (7)$$

Here ϵ_{12} is the BLER resulting from the first stage of the SIC decoding and $1 - \epsilon_{12}$ denotes the success in the first stage. The average BLER ϵ_{11} , results from the interference-free decoding in the second stage. These are respectively given by

$$\epsilon_{12} = \Phi(\gamma_{12}, N_2, M) \quad \text{and} \quad \epsilon_{11} = \Phi(\gamma_{11}, N_1, M). \quad (8)$$

The user u_2 directly decodes its message, so the instantaneous BLER ϵ_2 is

$$\epsilon_2 = \epsilon_{22} = \Phi(\gamma_{22}, N_2, M). \quad (9)$$

Then the average BLERs at the two users are obtained by

$$\bar{\epsilon}_1 = \mathbb{E}[\epsilon_1] \quad \text{and} \quad \bar{\epsilon}_2 = \mathbb{E}[\epsilon_2]. \quad (10)$$

By taking the expectation of the instantaneous BLER over the SINR distribution average BLER $\bar{\epsilon}_{ij}$ is given as

$$\bar{\epsilon}_{ij} = \int_0^\infty \Phi(\gamma_{ij}, N_j, M) f_{\gamma_{ij}}(x) dx \quad (11)$$

$$\approx \int_0^\infty Q\left(\frac{\log_2(1 + \gamma_{ij}) - \frac{N_j}{M}}{\sqrt{\frac{v_{ij}}{M}}}\right) f_{\gamma_{ij}}(x) dx, \quad (12)$$

where $f_{\gamma_{ij}}(x)$ is the probability density function (PDF) of the SINR γ_{ij} . Equation (12), does not have a closed form solution and based on work the by Makki et al. [7], $Q\left(\frac{\log_2(1 + \gamma_i) - \frac{N_i}{M}}{\sqrt{\frac{v_i}{M}}}\right) \approx \Xi_i(\gamma_i)$ can be approximated as

$$\Xi_i(\gamma_i) = \begin{cases} 1, & \gamma_i \leq v_i, \\ \frac{1}{2} - \lambda_i(\gamma_i - \theta_i), & v_i < \gamma_i < \tau_i, \\ 0, & \gamma_i \geq \tau_i, \end{cases} \quad (13)$$

where

$$\lambda_i = \sqrt{\frac{M}{2\pi\left(2^{\frac{2N_i}{M}} - 1\right)}}, \quad \theta_i = 2^{\frac{N_i}{M}} - 1, \quad (14)$$

$$v_i = \theta_i - \frac{1}{2\lambda_i} \quad \text{and} \quad \tau_i = \theta_i + \frac{1}{2\lambda_i}. \quad (15)$$

Using this approximation in (12), the average BLER $\bar{\epsilon}_i$ is given by

$$\bar{\epsilon}_i = \lambda_i \int_{v_i}^{\tau_i} F_{\gamma_i}(x) dx, \quad (16)$$

where $F_{\gamma_i}(x)$ is the cumulative distribution function (CDF) of the SINR γ_i .

B. Average BLER for Decoding Far User's Information

Based on the work by Cai et al. [3], the CDF of the SINR γ_{i2} for decoding of u_2 's message with HARQ-CC is derived as

$$F_{\gamma_{i2}}(r) \approx c_i^T \sum_{\{p_1, \dots, p_N\} \in \mathbb{P}} \Lambda \left[\prod_{n=1}^N \Psi^{p_n}(a_n) \right] \times \sum_{k=1}^L (\omega_k \ln 2) E_1\left(\frac{S_{k,N}}{r}\right). \quad (17)$$

The description of the variables and functions is given under (18). The proof is provided in Appendix A as an extension of the work in [3].

With the CDF of γ_{i2} in (17), an approximation for the average BLER, $\bar{\epsilon}_{i2}$ can be computed using (16) as

$$\bar{\epsilon}_{i2} \approx \lambda_2 c_i^T \sum_{\{p_1, \dots, p_N\} \in \mathbb{P}} \Lambda \left[\prod_{n=1}^N \Psi^{p_n}(a_n) \right] \times \sum_{k=1}^L (\omega_k \ln 2) [\Omega(v_2, S_{k,N}) - \Omega(\tau_2, S_{k,N})] \quad (18a)$$

where,

$$c_i = \frac{2\pi\kappa\alpha_2}{N\mu_i\rho} e^{\frac{1}{\mu_i\rho\alpha_1}}, \quad \kappa = \frac{\alpha_2}{\alpha_1}, \quad \mu_i = \frac{1}{1 + d_i\eta}, \quad (18b)$$

$$a_n = \cos\left(\frac{2n-1}{2N}\pi\right) \text{ for } n = 1, 2, \dots, N, \quad (18c)$$

$$\Lambda = \frac{T!}{\prod_{n=1}^N p_n!}, \quad \mathbb{P} = \left\{ p_1, \dots, p_N \mid T = \sum_{n=1}^N p_n \right\}, \quad (18d)$$

$$S_{k,N} = \frac{k\kappa \ln 2}{2} \sum_{n=1}^N p_n (a_n + 1), \quad (18e)$$

$$\Psi(a_n) = \frac{\sqrt{1 - a_n^2}}{(2\alpha_2 - \alpha_1\kappa(a_n + 1))^2} e^{-\frac{2\alpha_2}{\mu_i\rho\alpha_1(2\alpha_2 - \alpha_1\kappa(a_n + 1))}}, \quad (18f)$$

$$\omega_k = (-1)^{\frac{k}{2} + k} \sum_{[j=\frac{k+1}{2}]^{\min(k, \frac{L}{2})}} \frac{j^{\binom{k}{2} + 1}}{\left(\frac{L}{2}\right)!} \binom{L}{j} \binom{2j}{j} \binom{j}{k-j}, \quad (18g)$$

$$\Omega(x, y) = xe^{-\frac{y}{x}} - (x + y)E_1\left(\frac{y}{x}\right), \quad (18h)$$

$E_1(x) = \int_x^\infty \frac{e^{-t}}{t} dt$ is the exponential integral function and N, L are complexity-accuracy trade-off parameters. The proof is provided in Appendix B.

C. Average BLER for Interference-free Decoding of the Near User's Information

For u_1 decoding its information with HARQ-CC after T transmissions, the SNR is given by (4) which is

$$Z = \sum_{t=1}^T \gamma_{11}^t = \sum_{t=1}^T \rho\alpha_1 |\tilde{h}_{1,t}|^2 = \sum_{t=1}^T \rho\alpha_1 \mu_1 |h_{1,t}|^2,$$

where $\mu_1 = \frac{1}{\sqrt{1 + d_1\eta}}$. Since $h_{1,t} \sim \mathcal{CN}(0, 1)$, $|h_{1,t}|^2$ is an exponential variable, $\rho\alpha_1 \mu_1 |h_{1,t}|^2$ is exponentially distributed such that $|h_{1,t}|^2 \sim \text{Exp}\left(\frac{1}{\rho\alpha_1 \mu_1}\right)$. The sum of T exponential random variables is a Gamma distributed random variable with T degrees of freedom. Therefore Z can be described as

$$Z \sim \text{Gamma}\left(T, \frac{1}{\rho\alpha_1 \mu_1}\right). \quad (19)$$

Then the CDF of Z is,

$$F_Z(r) = \frac{1}{\Gamma(T)} \gamma\left(T, \frac{r}{\rho\alpha_1 \mu_1}\right) \quad (20)$$

where $\Gamma(k) = \int_0^\infty t^{k-1} e^{-t} dt$ is the Gamma function and $\gamma(k, x) = \int_0^x t^{k-1} e^{-t} dt$ is the lower incomplete Gamma function.

Therefore, $\bar{\epsilon}_{11}$ can be computed using (16) resulting in

$$\bar{\epsilon}_{11} = \lambda_1 (\Upsilon(\tau_1) - \Upsilon(v_1)), \quad (21a)$$

where,

$$\Upsilon(x) = \frac{1}{\Gamma(T)} \left[\gamma \left(T, \frac{x}{\rho\alpha_1\mu_1} \right) - \rho\alpha_1\mu_1\gamma \left(T+1, \frac{x}{\rho\alpha_1\mu_1} \right) \right], \quad (21b)$$

and the μ_1, λ_1, τ_1 and ν_1 as defined before. The proof is provided in Appendix C.

If the users are served using OMA, the blocklength or the number of channels uses available for transmission, M would be shared between the two users and their messages will be transmitted utilizing the full power for that particular number of channel uses without interference from the other user. Note that the average BLER for OMA will have the same form as in (21) with $\alpha_1 = 1$ and μ_1 will be replaced by μ_i for $i = 1, 2$.

IV. BLOCKLENGTH AND POWER ALLOCATION

A. Asymptotic BLER approximations

Due to the mathematical complexity of the derived expressions in Section III, asymptotic expressions are derived in high SNR conditions. In short packet communications, the rate $\frac{N_i}{M}$ is small [10], which leads to $\tau_{i,M} - \nu_{i,M}$ being smaller. Thus, the integration in (16) can be approximated using the Riemann integral approximation, $\int_a^b g(x)dx = (b-a)g(\frac{a+b}{2})$ such that

$$\bar{\epsilon}_i^\infty \approx \lambda_i(\tau_i - \nu_i)F_{\gamma_i}\left(\frac{\tau_i + \nu_i}{2}\right) = F_{\gamma_i}(\theta_i), \quad (22)$$

where the superscript ∞ denotes the asymptotic approximation.

The average BLER targets for ultra reliable communication are in the order of 10^{-5} or lower and can be achieved with high transmit SNR. Therefore, $1 - \epsilon_{12}^\infty \approx 1$ in (10), which results in $\epsilon_1^\infty \approx \epsilon_{12}^\infty + \epsilon_{11}^\infty$. Therefore, $\bar{\epsilon}_1^\infty$ approximates to $\mathbb{E}[\epsilon_{12}^\infty] + \mathbb{E}[\epsilon_{11}^\infty] = \bar{\epsilon}_{12}^\infty + \bar{\epsilon}_{11}^\infty$ and $\bar{\epsilon}_2^\infty$ can be obtained by $\mathbb{E}[\epsilon_{22}^\infty] = \bar{\epsilon}_{22}^\infty$.

B. Required blocklength and power allocation

The problem of finding the required blocklength M , which guarantees the target BLERs can be stated as

$$\text{find } M \quad (23a)$$

$$\text{s.t } \bar{\epsilon}_1 = \bar{\epsilon}_1^R \quad (23b)$$

$$\bar{\epsilon}_2 = \bar{\epsilon}_2^R \quad (23c)$$

$$\alpha_1 + \alpha_2 = 1 \quad (23d)$$

$$0 < \alpha_1 < 0.5 \quad (23e)$$

$$\text{for given } \rho, \mu_1, \mu_2, N_1, N_2, T, \quad (23f)$$

where the required BLERs for the two users are $\bar{\epsilon}_1^R$ and $\bar{\epsilon}_2^R$. The conditions in (23b) and (23c) ensure the reliability targets of the users while (23d) and (23e) arise from the NOMA principle. Since $\alpha_2 = 1 - \alpha_1$, α_2 can be omitted from the expressions.

According to Section IV-A, $\bar{\epsilon}_1^R$ and $\bar{\epsilon}_2^R$ can be expressed as,

$$\bar{\epsilon}_1^R = \bar{\epsilon}_{12}^\infty + \bar{\epsilon}_{11}^\infty \quad \text{and} \quad \bar{\epsilon}_2^R = \bar{\epsilon}_{22}^\infty. \quad (24)$$

Since u_1 is the stronger user with high channel gain and according to NOMA principle more power is allocated to u_2 ,

average BLER for decoding u_2 's information in the first stage of SIC at u_1 is smaller than the average BLER for the second stage of SIC when interference-free decoding of u_1 is done. Therefore, for simplicity, the average BLER for the first stage of SIC in u_1 is considered as $\bar{\epsilon}_{12}^\infty = \delta\bar{\epsilon}_{11}^\infty$, where $\delta < 1$. Then, from (24) $\bar{\epsilon}_{11}^\infty$ can be written as

$$\bar{\epsilon}_{11}^\infty = \frac{\bar{\epsilon}_1^R}{1 + \delta}. \quad (25)$$

Using the approximation with the Riemann integral as in (22), $\bar{\epsilon}_{11}^\infty$ and $\bar{\epsilon}_{22}^\infty$ can be obtained as

$$\bar{\epsilon}_{11}^\infty \approx F_{\gamma_{11}}(\theta_1) \quad (26)$$

$$\bar{\epsilon}_{22}^\infty \approx F_{\gamma_{22}}(\theta_2). \quad (27)$$

Therefore, from (26) the blocklength M , which satisfies the reliability targets can be found as

$$M = \frac{N_1}{\log_2 \left(1 + \mu_1 \rho \alpha_1 \Gamma(T) \gamma^{-1} \left(T, \frac{\bar{\epsilon}_1^R}{1 + \delta} \right) \right)}, \quad (28)$$

where $\gamma^{-1}(k, x)$ is the inverse of the lower incomplete Gamma function. With the use of (28) in (27) the required M can be found. Also, the required blocklength for OMA can be obtained by the addition of the blocklengths needed to achieve their reliability targets using a similar expression to (28), with $\alpha_1 = 1$ and μ_1 replaced by μ_i for $i = 1, 2$.

Let G be a function such that $G(\alpha_1) \triangleq \bar{\epsilon}_2 - \bar{\epsilon}_2^R$ according to the condition in (23c). Solving $G(\alpha_1) = 0$ will give the α_1 needed to achieve for the required blocklength, which can be used to find the required blocklength M_{req} using (28). Noting that $G(\alpha_1)$ is highly nonlinear, the solution can be computed using Algorithm 1.

Algorithm 1 Power Allocation and Required Blocklength for NOMA with HARQ CC

- 1: **Input** : $\bar{\epsilon}_1^R, \bar{\epsilon}_2^R, \rho, \mu_1, \mu_2, N_1, N_2, T, \delta$ and tolerance ν
 - 2: **Output** : M_{req} and α_1^*
 - 3: **Initialize** : $\alpha_1^- = 0$ and $\alpha_1^+ = 0.5$
 - 4: **while** $|G(\alpha_1^c)| > \nu$ **do**
 - 5: set $\alpha_1^c \leftarrow (\alpha_1^+ + \alpha_1^-)/2$
 - 6: compute $G(\alpha_1^c)$ based on (27) and (28)
 - 7: **if** $G(\alpha_1^c)G(\alpha_1^+) > 0$ **then** set $\alpha_1^+ \leftarrow \alpha_1^c$
 - 8: **else** : set $\alpha_1^- \leftarrow \alpha_1^c$
 - 9: **end while**
 - 10: set $\alpha_1^* \leftarrow \alpha_1^c$
 - 11: compute M_{req} using (28) with α_1^*
-

V. NUMERICAL RESULTS

Monte Carlo simulations are carried out based on the results for the decoding error probability in short blocklengths to verify the accuracy of the approximations derived in Section III. In all the simulations, the complexity-accuracy parameters $N = 30$ and $L = 18$ to ensure the numerical accuracy. The path loss exponent $\eta = 2$ while $d_1 = 3$ m and $d_2 = 7$ m.

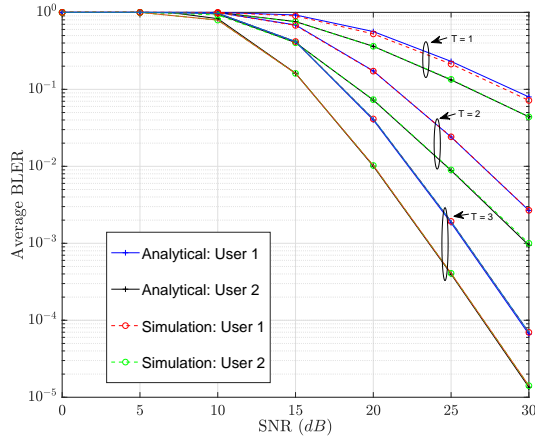


Fig. 1. Average BLER vs. transmit SNR (ρ) for different number of transmission rounds (T) with $\alpha_1 = 0.1$, $\alpha_2 = 0.9$, $N_1 = N_2 = 160$ and $M = 200$.

Figure 1 shows the average BLERs plotted against the transmit SNR (ρ) for different maximum transmission rounds. The approximations derived match with the Monte Carlo simulation results, which prove the accuracy of the expressions in (18) and (21). According to Figure 1, the far user always has a smaller average BLER than the near user, u_1 . The reason is that higher power is allocated for the far user for user fairness in the NOMA principle. Also, with the increasing number of maximum transmission rounds allowed, the average BLER decreases for a particular transmit SNR.

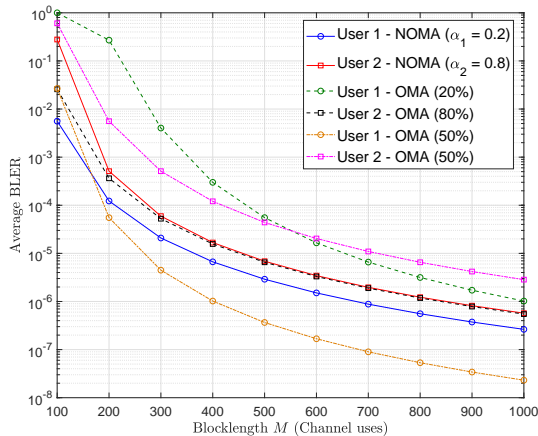


Fig. 2. Average BLER vs. blocklength for NOMA with $\alpha_1 = 0.2$ and OMA with 20% and 50% for u_1 , $\rho = 30$ dB, $T = 3$, $N_1 = N_2 = 300$.

In Figure 2, average BLER is plotted with the blocklength at $\rho = 30$ dB, $T = 3$ with power allocation α_1 set to 0.2 and N_1 and N_2 set to 300. The comparison with OMA is provided for two scenarios as 20%-80% and 50%-50% blocklength share for u_1, u_2 respectively. It is clear from Figure 2, when the blocklength increases the average BLERs in all scenarios

decrease monotonically which is desirable. One interesting result is that the performance of u_2 in NOMA and OMA with 80% share is almost similar with increasing blocklength. However, u_1 has a lower average BLER when NOMA is used compared to OMA with 20% share of blocklength. Nevertheless, as the blocklength increase, this difference in performance between NOMA and OMA decreases. For the second scenario, blocklength is shared equally between the two users. The performance of u_2 degrades significantly compared to the performance with 80% share. However, u_1 achieves a lower BLER than NOMA since a higher number of channel uses is available to u_1 . Although u_1 has a lower average BLER with an equal share in OMA than NOMA, u_2 's average BLER degrades significantly. Therefore, the NOMA scheme delivers fairness to both users, unlike OMA, since the difference in average BLER between two users is smaller than in OMA while achieving considerable average BLER performance for both users.

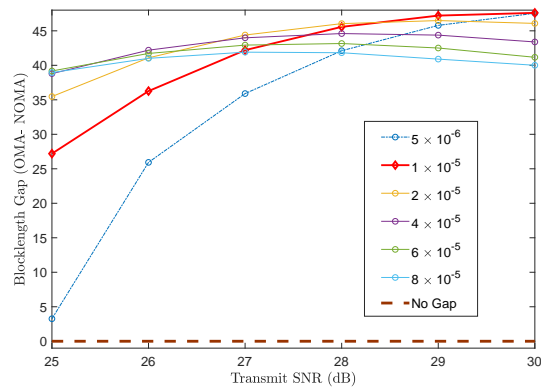


Fig. 3. Blocklength gap between OMA and NOMA vs. transmit SNR for $\bar{\epsilon}_1^R = 1 \times 10^{-5}$ and varying $\bar{\epsilon}_2^R$ with $N_1 = N_2 = 300$, $T = 3$, $\delta = 0.1$.

Figure 3 shows the gap between OMA and NOMA for the required blocklength to achieve a given reliability target of 1×10^{-5} for u_1 and varying $\bar{\epsilon}_2^R$ for u_2 using Algorithm 1. The value of δ is set to 0.1. Here the gap is taken by subtracting the NOMA blocklength from the OMA blocklength. It can be seen from Figure 3, NOMA has a smaller blocklength than OMA for the given reliability targets since the gap is positive. The bold red curve represents both users having the same reliability target of 1×10^{-5} and NOMA always has a lower blocklength requirement and this gap increases as the transmit SNR increases. For lower reliability target such as 5×10^{-6} for u_2 , the gap is smaller as seen from the dashed curve. Thus, NOMA has a lower blocklength requirement than OMA which leads to having lower latency when the reliability targets are in the order of 10^{-5} .

VI. CONCLUSION

This paper analyzed the performance of NOMA with HARQ-CC for finite blocklength by deriving tight closed-form approximations for the average BLER for two users. The comparison with OMA was done proving that NOMA could

meet lower average BLER requirements such as 1×10^{-5} while ensuring user fairness better than OMA. Further, an algorithm to determine the blocklength required to meet the reliability requirements of the two users was developed based upon the asymptotic expressions considering high SNR conditions. Simulations proved that NOMA has a lower blocklength requirement in high SNR leading to lower latency compared to OMA when the reliability requirements are in the order of 10^{-5} . Analysis with multiple antennas is intended to be done as future work.

APPENDIX A PROOF OF EQUATION (16)

The PDF of the SINR for decoding u_2 s signal at u_i where $i = 1, 2$ after T rounds of transmissions, $Z = \gamma_{i2} = \sum_{t=1}^T \gamma_{i2}^t$ is given by Equation (35) in [3]. The CDF is calculated by extending that work. By taking the integral over $f_Z(z)$ with respect to z as

$$F_Z(r) = \int_{-\infty}^r f_Z(z) dz = \int_0^r f_Z(z) dz.$$

$$\approx c_i^T \sum_{\{p_1, \dots, p_N\} \in \mathbb{P}} \Lambda \left[\prod_{n=1}^N \Psi^{p_n}(a_n) \right] \sum_{k=1}^L (\omega_k \ln 2) I_k,$$

$$\text{where } I_k = \int_0^r \frac{1}{z} e^{-\frac{k \ln 2}{2z} \sum_{n=1}^N p_n (a_n + 1)} dz.$$

By change of variables with $u = \frac{1}{z}$ integral in I_k converts to,

$$I_k = \int_{\frac{S_{k,N}}{r}}^{\infty} \frac{1}{u} e^{-u} du = E_1 \left(\frac{S_{k,N}}{r} \right)$$

where $S_{k,N}$ as defined in (18e) and $E_1(x)$ is the exponential integral function defined by $E_1(x) = \int_x^{\infty} \frac{e^{-t}}{t} dt$.

APPENDIX B PROOF OF EQUATION (18)

CDF of γ_{i2} is given by (17). Computation of the approximation for average BLER, $\bar{\epsilon}_{i2}$ is provided here using (16).

$$\bar{\epsilon}_{i2} \approx \lambda_2 \int_{v_2}^{\tau_2} c_i^T \sum_{\{p_1, \dots, p_N\} \in \mathbb{P}} \Lambda \left[\prod_{n=1}^N \Psi^{p_n}(a_n) \right]$$

$$\times \sum_{k=1}^L (\omega_k \ln 2) E_1 \left(\frac{S_{k,N}}{x} \right) dx$$

$$= \lambda_2 c_i^T \sum_{\{p_1, \dots, p_N\} \in \mathbb{P}} \Lambda \left[\prod_{n=1}^N \Psi^{p_n}(a_n) \right] \sum_{k=1}^L (\omega_k \ln 2) J_k,$$

$$\text{where } J_k = \int_{v_2}^{\tau_2} E_1 \left(\frac{S_{k,N}}{x} \right) dx = -S_{k,N} \int_{\frac{S_{k,N}}{v_2}}^{\frac{S_{k,N}}{\tau_2}} \frac{E_1(v)}{v^2} dv.$$

Using integration by parts twice and Leibniz integral rule

$$= -S_{k,N} \left(\left[-\frac{E_1(v)}{v} + \frac{e^{-v}}{v} \right]_{\frac{S_{k,N}}{v_2}}^{\frac{S_{k,N}}{\tau_2}} + \int_{\frac{S_{k,N}}{v_2}}^{\frac{S_{k,N}}{\tau_2}} \frac{e^{-v}}{v} dv \right)$$

$$= \Omega(v_2, S_{k,N}) - \Omega(\tau_2, S_{k,N}),$$

where $\Omega(x, y)$ as defined (18h) which completes the proof.

APPENDIX C PROOF OF EQUATION (21)

The average BLER, $\bar{\epsilon}_{11}$ can be computed using (16) with the CDF for γ_{11} given by the (20) as

$$\bar{\epsilon}_{11} = \lambda_1 \int_{v_1}^{\tau_1} F_{\gamma_{11}}(x) dx = \lambda_1 \int_{v_1}^{\tau_1} \frac{1}{\Gamma(T)} \gamma \left(T, \frac{x}{\rho \alpha_1 \mu_1} \right) dx.$$

Using integration by parts,

$$= \lambda_1 \frac{1}{\Gamma(T)} \left(\left[x \gamma \left(T, \frac{x}{\rho \alpha_1 \mu_1} \right) \right]_{v_1}^{\tau_1} - Q_1 \right), \quad (31)$$

where $Q_1 = \int_{v_1}^{\tau_1} x \frac{d}{dx} \gamma \left(T, \frac{x}{\rho \alpha_1 \mu_1} \right) dx$. By definition of the lower incomplete Gamma function, change of variables with $u = \delta_1 x$ where $\delta_1 = \frac{1}{\rho \alpha_1 \mu_1}$ and Leibniz rule

$$Q_1 = \frac{1}{\delta_1} \int_{\delta_1 v_1}^{\delta_1 \tau_1} u \frac{d}{du} \left(\int_0^u t^{T-1} e^{-t} dt \right) du.$$

$$= \rho \alpha_1 \mu_1 \left[\gamma \left(T + 1, \frac{x}{\rho \alpha_1 \mu_1} \right) \right]_{v_1}^{\tau_1}.$$

Using the result of Q_1 in (31) completes the proof.

REFERENCES

- [1] Choi J. (2008) H-ARQ Based Non-Orthogonal Multiple Access with Successive Interference Cancellation. In: IEEE GLOBECOM 2008 - 2008 IEEE Global Telecommunications Conference, pp. 1-5.
- [2] Choi J. (2016) On HARQ-IR for Downlink NOMA Systems. IEEE Transactions on Communications 64, pp. 3576-3584.
- [3] Cai D., Ding Z., Fan P. & Yang Z. (2018) On the Performance of NOMA With Hybrid ARQ. IEEE Transactions on Vehicular Technology 67, pp. 10033-10038.
- [4] Xu Y., Cai D., Fang F., Ding Z., Shen C. & Zhu G. (2018) Outage Analysis and Power Allocation for HARQ-CC Enabled NOMA Downlink Transmission. In: 2018 IEEE Global Communications Conference (GLOBECOM), pp. 1-6.
- [5] Devassy R., Durisi G., Popovski P. & Strm E.G. (2014) Finite-blocklength analysis of the ARQ-protocol throughput over the Gaussian collision channel. In: 2014 6th International Symposium on Communications, Control and Signal Processing (ISCCSP), pp. 173-177.
- [6] Makki B., Svensson T. & Zorzi M. (2014) Green communication via Type-I ARQ: Finite block-length analysis. In: 2014 IEEE Global Communications Conference, pp. 2673-2677.
- [7] Makki B., Svensson T. & Zorzi M. (2014) Finite Block-Length Analysis of the Incremental Redundancy HARQ. IEEE Wireless Communications Letters 3, pp. 529-532.
- [8] Dosti E., Shehab M., Alves H. & Latva-aho M. (2017) Ultra reliable communication via CC-HARQ in finite block-length. In: 2017 European Conference on Networks and Communications (EuCNC), pp. 1-5.
- [9] Yu Y., Chen H., Li Y., Ding Z. & Vucetic B. (2018) On the Performance of Non-Orthogonal Multiple Access in Short-Packet Communications. IEEE Communications Letters 22, pp. 590-593.
- [10] Huang X. & Yang N. (2019) On the Block Error Performance of Short-Packet Non-Orthogonal Multiple Access Systems. In: ICC 2019 - 2019 IEEE International Conference on Communications (ICC), pp. 1-7.
- [11] Polyanskiy Y., Poor H.V. & Verdú S. (2010) Channel Coding Rate in the Finite Blocklength Regime. IEEE Trans. Inf. Theor. 56, pp. 2307-2359. URL: <http://dx.doi.org/10.1109/TIT.2010.2043769>.