

Coexistence of Heterogeneous Services in the Uplink with Discrete Signaling and Treating Interference as Noise

Min Qiu[†], Yu-Chih Huang^{*}, and Jinhong Yuan[†]

[†]School of Electrical Engineering and Telecommunications, University of New South Wales, Sydney, Australia

^{*}Institute of Communications Engineering, National Yang Ming Chiao Tung University, Hsinchu City, Taiwan

E-mail: {min.qiu, j.yuan}@unsw.edu.au, jerryhuang@nycu.edu.tw

Abstract—The problem of enabling the coexistence of heterogeneous services, e.g., different ultra-reliable low-latency communications (URLLC) services and/or enhanced mobile broadband (eMBB) services, in the uplink is studied. Each service has its own error probability and blocklength constraints and the longer transmission block suffers from heterogeneous interference. Due to the latency concern, the decoding of URLLC messages cannot leverage successive interference cancellation (SIC) and should always be performed before the decoding of eMBB messages. This can significantly degrade the achievable rates of URLLC users when the interference from other users is strong. To overcome this issue, we propose a new transmission scheme based on discrete signaling and treating interference as noise decoding, i.e., without SIC. Guided by the deterministic model, we provide a systematic way to construct discrete signaling for handling heterogeneous interference effectively. We demonstrate theoretically and numerically that the proposed scheme can perform close to the benchmark scheme based on capacity-achieving Gaussian signaling with the assumption of perfect SIC.

Index Terms—Multiple access channels, finite blocklength, discrete modulations, treating interference as noise.

I. INTRODUCTION

The fifth generation (5G) wireless communication systems and beyond are expected to support a variety of service types such as enhanced mobile broadband (eMBB) and ultra-reliable low-latency communications (URLLC). To this end, the idea of enabling the coexistence of heterogeneous services in the same radio access network has attracted many interests [1]–[6]. The current proposal is to slice the network by allocating orthogonal resources, e.g., in time/frequency domain, for each service [1]. Since different services are isolated from each other, their own quality-of-service requirements are guaranteed. However, this approach could lead to very low spectral and energy efficiency when the number of devices is large.

The work of Min Qiu and Jinhong Yuan was supported in part by the Australian Research Council (ARC) Discovery Project under Grant DP220103596, and in part by the ARC Linkage Project under Grant LP200301482. The work of Yu-Chih Huang was supported by the National Science and Technology Council, Taiwan, under Grant 111-2221-E-A49-069-MY3. This work was also supported in part by the Higher Education Sprout Project of the National Yang Ming Chiao Tung University and Ministry of Education (MOE), Taiwan.

Various research has been carried out to investigate efficient coexistence mechanisms for heterogeneous services with improved efficiency and user fairness. Notably, [2], [3] introduced a non-orthogonal uplink transmission scheme for enabling simultaneous communication between different types of services and the base station. In the Gaussian multiple access channel (MAC) with homogeneous and infinite blocklength, the entire capacity region can be efficiently achieved by superposition coding and successive interference cancellation (SIC), together with time-sharing [7]. However, as pointed out in [2], [3], one salient difference between the conventional and heterogeneous cases is that the decoding of a URLLC transmission block cannot leverage SIC due to the latency constraints. That is, the base station should always decode URLLC services first before decoding eMBB of services regardless of channel conditions. Without SIC, the achievable rates of URLLC users could be severely degraded when the interference is strong. Note also that SIC can fail under finite blocklength constraint [8], [9]. In addition, to ensure ultra reliability for URLLC services, performing SIC by decoding eMBB first then URLLC would require the decoding of eMBB services to achieve at least the same reliability as for URLLC services. However, this is an overkill as eMBB messages are coded based on a lower reliability requirement than URLLC services. Further, SIC can introduce extra decoding delay, complexity, and error propagation, which can become pronounced when the number of users is large. In light of the above, the non-orthogonal schemes relying on SIC may not fully address the challenges from the coexistence of URLLC and other services; thereby, new proposals are called for.

The achievable rate region of the Gaussian MAC under finite blocklength in terms of the second-order term and the third-order term were rigorously derived in [10] and [11], respectively. In order to achieve the optimal first- and second-order rates, both works considered the input distribution to be shell codes, i.e., codewords drawn from a power shell. Moreover, the achievability bounds therein were derived under joint typically and joint maximum-likelihood decoding, respectively, with global error probability formalism and homogeneous blocklength constraints. Nevertheless, these results are not

applicable to the case with heterogeneous error probability and blocklength constraints. Moreover, joint decoding has a much higher complexity than SIC decoding, which is difficult to realize in practice.

Motivated by the benefits and challenges of heterogeneous services coexistence in the uplink, we aim to design a new and practical scheme based on discrete signaling and treating interference as noise (TIN) decoding. We emphasize that the proposed discrete signaling is formed by practical channel coding and discrete constellations, e.g., quadrature amplitude modulation (QAM), which are the prevailing setup in current communication systems [12]. Meanwhile, the TIN decoding only has single-user decoding complexity and latency, which is more suitable for decoding URLLC services. The designs of discrete signaling with TIN with infinite blocklength and homogeneous interference were investigated for the Gaussian broadcast channel (GBC) and the Gaussian interference channel with infinite blocklength in [13] and [14], respectively. Recently, we proposed a transmission scheme for the GBC under heterogeneous blocklength and error probability constraints [15]. However, the scheme therein cannot be applied here due to the subtle difference between the broadcast nature (i.e., all users' signals go through the same channel to one receiver) and the multiple access nature (i.e., each user's signal goes through a different channel to the receiver). Hence, it remains unclear on how to design the discrete signaling for each uplink user to effectively handle heterogeneous interference with TIN to achieve a rate close to that under Gaussian signaling with perfect SIC.

In this paper, we provide the design and analysis of a new multiple access scheme based on discrete signaling and TIN decoding to support heterogeneous services coexistence in the uplink. We first approximate the MAC channel as a linear deterministic model [16] and design a capacity-achieving coding scheme with TIN. We then translate the scheme from the deterministic model into the coded modulation scheme with TIN for the MAC channel model under heterogeneous constraints. The main feature of the proposed scheme is that the user with a longer transmission block can use different sets of constellations in different parts of its block depending on heterogeneous interference statistics. Whereas the encoding and decoding are single-user based. Our theoretical and numerical results show that the proposed scheme can achieve rate pairs close to the benchmark scheme that assumes Gaussian signaling and perfect SIC.

A. Notations

Random variables are written in upper case sans-serif fonts, e.g., X , while their realizations are in lower case form, e.g., x . $X_1 \stackrel{d}{=} X_2$ means that X_1 and X_2 are equal in distribution. All logarithms are base 2. $X \sim \text{unif QAM}(|\Lambda|, d_{\min}(\Lambda))$ represents that X is uniformly distributed over a zero mean regular QAM Λ with cardinality $|\Lambda|$, minimum distance $d_{\min}(\Lambda)$, and average energy $E_\Lambda = d_{\min}^2(\Lambda) \frac{|\Lambda|-1}{6}$. $\lceil x \rceil$ rounds x to the nearest integer greater than or equal to x . We define the operation $(x)^+ \triangleq \max\{0, x\}$. The binary field, the collections

of binary vectors of size n and binary matrices of size $m \times n$ are denoted by \mathbb{F}_2 , \mathbb{F}_2^n , and $\mathbb{F}_2^{m,n}$, respectively.

II. SYSTEM MODEL

We consider a two-user uplink MAC that consists of two senders and one receiver. We leave the extension of our work to the K -user case in our future work. Let $k \in \{1, 2\}$ be the user index. Let $\mathbf{x}_k \in \mathbb{C}^{N_k}$ represent the coded symbols of user k , where N_k is the symbol length. Due to heterogeneous blocklength constraints, we assume $N_1 \leq N_2$ without loss of generality. Each user's coded symbols satisfy the individual maximum power constraint per codeword given by

$$\frac{1}{N_k} \sum_{j=1}^{N_k} |x_k[j]|^2 \leq P_k. \quad (1)$$

Let $h_k \in \mathbb{C}$ represent the complex channel for user k . The j -th received symbol is

$$y[j] = \begin{cases} h_1 x_1[j] + h_2 x_2[j] + z[j], & j = 1, \dots, N_1 \\ h_2 x_2[j] + z[j], & j = N_1 + 1, \dots, N_2 \end{cases}, \quad (2)$$

where $z[j] \sim \mathcal{CN}(0, 1)$ is the i.i.d. Gaussian noise for $j \in \{1, \dots, N_2\}$. In addition, we assume the channel state information is available so that each transmitter can eliminate the channel phase by rotating its signal by $\frac{h_k^*}{|h_k|}$. In this regard, the channel model in (2) can be equivalently expressed as that with some $h_k \in \mathbb{R}$.

From (2), it can be seen that the first N_1 symbols of \mathbf{x}_2 and \mathbf{x}_1 are superimposed. This is because we consider that user 2, e.g., a URLLC user, needs to transmit its signals as soon as possible due to its urgency. In this case, the received symbol block of user 2 suffers from *heterogeneous interference* across its symbols. Note that user 2 can be a different URLLC user or a user of other service types that are less urgent, e.g., eMBB. We denote by $\text{SNR}_k = P_k |h_k|^2$, R_k , and ϵ_k the signal-to-noise ratio (SNR), achievable rate, and the requirement of the decoding error probability of user k , respectively.

III. PROPOSED DISCRETE SIGNALING SCHEME WITH TIN

In this section, we first present the finite blocklength achievable rate of the MAC with an arbitrary choice of discrete input distributions, blocklength, and error probability requirements. Using the insight obtained from the finite blocklength analysis, we approximate the considered channel model by using the linear deterministic model [16] and construct a capacity-achieving coding scheme with TIN. Next, we translate the proposed scheme of the deterministic channel to the coding scheme for the original MAC under heterogeneous constraints. We show analytically that the gap between the mutual information under TIN and the corresponding single-user capacity is upper bounded by a constant gap.

A. Finite Blocklength Achievable Rate

In this work, we consider that both users employ discrete and finite constellations. Specifically, user 1's symbol satisfies $X_1[j] \stackrel{\text{unif}}{\sim} \Lambda_1$ for $j = 1, \dots, N_1$, where Λ_1 is user 1's

constellation set. In contrast, user 2 uses two sets of constellations to handle heterogeneous interference. Thus, its symbol satisfies $X_2[j] \stackrel{\text{unif}}{\sim} \Lambda_{2,1}$ for $j = 1, \dots, N_1$ and $X_2[j] \stackrel{\text{unif}}{\sim} \Lambda_{2,2}$ for $j = N_1 + 1, \dots, N_2$. For ease of presentation, we let $X_{2,1} \stackrel{d}{=} X_2[j]$ and $Y_1 \stackrel{d}{=} Y[j]$ for $j = 1, \dots, N_1$, and let $X_{2,2} \stackrel{d}{=} X_2[j]$ and $Y_2 \stackrel{d}{=} Y[j]$ for $j = N_1 + 1, \dots, N_2$.

The finite blocklength achievable rate pairs given constellation tuples $(\Lambda_1, \Lambda_{2,1}, \Lambda_{2,2})$ are provided in Theorem 1.

Theorem 1. Let ϵ_k be the upper bound on the average TIN decoding error probability of user k . For the channel model in (2), the achievable rates of users 1 and 2 with TIN are

$$R_1 \leq I(X_1; Y_1) - \sqrt{\frac{V(X_1; Y_1)}{N_1}} Q^{-1}(\epsilon_1) + O\left(\frac{\log N_1}{N_1}\right), \quad (3)$$

$$R_2 \leq \frac{N_1 I(X_{2,1}; Y_1) + (N_2 - N_1) I(X_{2,2}; Y_2)}{N_2} - \frac{\sqrt{N_1 V(X_{2,1}; Y_1) + (N_2 - N_1) V(X_{2,2}; Y_2)}}{N_2} \times Q^{-1}(\epsilon_2) + O\left(\frac{\log N_2}{N_2}\right), \quad (4)$$

where we have user 2's mutual information $I(X_{2,t}; Y_t) = \mathbb{E}[i(X_{2,t}; Y_t)]$ and dispersion $V(X_{2,t}; Y_t) = \text{Var}[i(X_{2,t}; Y_t)]$ for $t \in \{1, 2\}$ with information densities $i(X_{2,1}; Y_1) = \log\left(\frac{\sum_{x_{2,1} \in \Lambda_{2,1}} e^{-|y-h_1 x_1 - h_2 x_{2,1}|}}{\frac{1}{|\Lambda_{2,1}|} \sum_{x_{2,1} \in \Lambda_{2,1}} \sum_{x_1 \in \Lambda_1} e^{-|y-h_1 x_1 - h_2 x_{2,1}|}}\right)$, and $i(X_{2,2}; Y_2) = \log\left(\frac{e^{-|y-h_2 x_{2,2}|}}{\frac{1}{|\Lambda_{2,2}|} \sum_{x_{2,2} \in \Lambda_{2,2}} e^{-|y-h_2 x_{2,2}|}}\right)$, and $Q^{-1}(x)$ is the inverse of Q function $Q(x) = \int_x^\infty \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt$. User 1's mutual information $I(X_1; Y_1)$ and dispersion $V(X_1; Y_1)$ can be easily obtained from $I(X_{2,1}; Y_1)$ and $V(X_{2,1}; Y_1)$ by swapping the argument between subscripts "1" and "2".

The proof of Theorem 1 follows from the steps in [15, Sec. V] and is omitted due to space limitation.

Remark 1. The impacts of the interfering symbol blocklength on the achievable rate is clearly shown in (4). Under heterogeneous interference, the achievable rate is no longer determined by a single signal-to-interference-plus-noise ratio (SINR) as in the homogeneous blocklength case. In addition, the relationship between heterogeneous blocklengths and achievable rate in (4) holds for other i.i.d. inputs, e.g., i.i.d. Gaussian codes. However, for the shell codes commonly used for deriving larger finite blocklength achievable rates in many works [8], [10], [11], this relation may no longer hold as the distribution of each symbol of shell codes is not independent. ■

From Theorem 1, we note that for given power, blocklength, and error probability constraints, larger mutual information and smaller dispersion are desirable for achieving a larger rate. In the subsequent sections, we introduce the proposed design and justify its effects on mutual information and dispersion.

B. Deterministic MAC under Heterogeneous Interference

Theorem 1 hints that to attain a larger achievable rate, different signaling designs are required for interfering and interference-free parts. To this end, we approximate the MAC channel under heterogeneous blocklength constraints as a concatenation of two deterministic channels and propose a capacity-achieving coding scheme with TIN. Let $n_k \triangleq \lceil \log \text{SNR}_k \rceil^+$ be the approximated single-user capacity of users k , $k \in \{1, 2\}$. By adopting the idea in [16], the channel model of (2) is approximated as

$$Y_1 = \mathbf{S}^{q-n_1} X_1 \oplus \mathbf{S}^{q-n_2} X_{2,1}, \quad (5)$$

$$Y_2 = X_{2,2}, \quad (6)$$

where the multiplication and summation \oplus are over \mathbb{F}_2 , $X_1, X_{2,1}, Y_1 \in \mathbb{F}_2^q$, $X_{2,2}, Y_2 \in \mathbb{F}_2^{n_2}$ are binary column vectors and the subscripts "1" and "2" follow those defined in the first paragraph of Sec. III-A, $q = \max\{n_1, n_2\}$, and \mathbf{S} is a $q \times q$ down shift matrix. The operational meaning of the binary column vector is that each of its entries represents a power level or a bit. $\mathbf{S}^{q-n_1} X_1$ models the channel effect, where the lowest $q - n_1$ bits of X_1 are shifted down below the noise level and truncated.

We first assume $|h_1| > |h_2|$ such that $q = n_1$. For the channel model in (5) and (6), the set of non-negative rate tuples $(m_1, m_{2,1}, m_{2,2})$ are achievable if

$$m_1 + m_{2,1} \leq n_1, \quad (7)$$

$$m_{2,1} \leq n_2, \quad (8)$$

$$m_{2,2} \leq n_2. \quad (9)$$

Observe that (7) and (8) are the deterministic MAC capacity region and (9) is the deterministic point-to-point channel capacity, which are within 1 bit of the capacity of their Gaussian counterparts [16], respectively. It is worth emphasizing that the considered deterministic MAC is different from the one in [16] where all user 2's bits experience the same level of interference strength. Different from the achievable scheme that relies on SIC [16], we show how to design input distributions $(X_1, X_{2,1}, X_{2,2})$ to achieve the above rate region with TIN.

Let $U_1 \in \mathbb{F}_2^{m_1}$ and $U_{2,1} \in \mathbb{F}_2^{m_{2,1}}$ be users 1 and 2's message vectors, respectively, with each entry drawn independently and uniformly distributed from \mathbb{F}_2 . Let $G_1 \in \mathbb{F}_2^{q, m_1}$ and $G_{2,1} \in \mathbb{F}_2^{q, m_{2,1}}$ be the generator matrices such that $X_1 = G_1 U_1$ and $X_{2,1} = G_{2,1} U_{2,1}$. The achievable rate of user 1 under TIN is

$$\begin{aligned} I(X_1; Y_1) &= H(Y_1) - H(Y_1|X_1) \\ &= H(\mathbf{S}^{q-n_1} G_1 U_1 \oplus \mathbf{S}^{q-n_2} G_{2,1} U_{2,1}) - H(\mathbf{S}^{q-n_2} G_{2,1} U_{2,1}) \\ &= \text{rank}([\mathbf{S}^{q-n_1} G_1, \mathbf{S}^{q-n_2} G_{2,1}]) - \text{rank}(\mathbf{S}^{q-n_2} G_{2,1}), \end{aligned} \quad (10)$$

where the multiplication and addition are over \mathbb{F}_2 . User 2's rate under TIN $I(X_{2,1}; Y_1)$ can be easily obtained from (10) by swapping the arguments between subscripts "1" and "2".

To achieve the rate region of (7)-(8), we propose

$$\mathbf{G}_1 = \begin{bmatrix} \mathbf{F}_{1,1} \\ \mathbf{0}^{n_1-n_2-m_{1,1}, m_1} \\ \mathbf{0}^{m_{2,1}, m_1} \\ \mathbf{0}^{n_2-m_{2,1}-m_{1,2}, m_1} \\ \mathbf{F}_{1,2} \end{bmatrix}, \mathbf{G}_{2,1} = \begin{bmatrix} \mathbf{F}_{2,1} \\ \mathbf{0}^{n_2-m_{2,1}, m_{2,1}} \\ \mathbf{0}^{n_1-n_2, m_{2,1}} \end{bmatrix}, \quad (11)$$

where $\mathbf{F}_{1,1} \in \mathbb{F}_2^{m_{1,1}, m_1}$, $\mathbf{F}_{1,2} \in \mathbb{F}_2^{m_{1,2}, m_1}$, $\mathbf{F}_{2,1} \in \mathbb{F}_2^{m_{2,1}, m_{2,1}}$, are submatrices with linearly independent rows, $m_{1,1} = \min\{(m_1 + m_{2,1} - n_2)^+, n_1 - n_2\}$, $m_{1,2} = \min\{m_1, n_2 - m_{2,1}\}$ such that $m_1 = m_{1,1} + m_{1,2}$. Substituting the proposed \mathbf{G}_1 and $\mathbf{G}_{2,1}$ into (10), we have that

$$\begin{aligned} & \text{rank}([\mathbf{S}^{q-n_1} \mathbf{G}_1, \mathbf{S}^{q-n_2} \mathbf{G}_{2,1}]) \\ &= \text{rank} \left(\begin{bmatrix} \mathbf{F}_{1,1} & \mathbf{0}^{n_1-n_2, m_{2,1}} \\ \mathbf{0}^{n_k-n_{2,1}-m_{1,1}, m_1} & \mathbf{0}^{n_1-n_2, m_{2,1}} \\ \mathbf{0}^{m_{2,1}, m_1} & \mathbf{F}_{2,1} \\ \mathbf{0}^{n_2-m_{2,1}-m_{1,2}, m_1} & \mathbf{0}^{n_2-m_{2,1}, m_{2,1}} \\ \mathbf{F}_{1,2} & \mathbf{0}^{n_2-m_{2,1}, m_{2,1}} \end{bmatrix} \right) \quad (12) \\ &= m_1 + m_{2,1}, \quad (13) \end{aligned}$$

$$\text{rank}(\mathbf{S}^{q-n_2} \mathbf{G}_{2,1}) = \text{rank}(\mathbf{F}_{2,1}) = m_{2,1}. \quad (14)$$

Substituting (13) and (14) into (10) results in $I(X_1; Y_1) = m_1$. Similarly, we can follow the above steps to obtain that $I(X_{2,1}; Y_1) = m_{2,1}$.

Let $\mathbf{U}_{2,2} \in \mathbb{F}_2^{m_{2,2}}$ be another part of user 2's message vector. To achieve the rate in (9), we design

$$\mathbf{G}_{2,2} = \begin{bmatrix} \mathbf{F}_{2,2} \\ \mathbf{0}^{n_2-m_{2,2}, m_{2,2}} \end{bmatrix}, \quad (15)$$

where $\mathbf{F}_{2,2} \in \mathbb{F}_2^{m_{2,2}, m_{2,2}}$. Then, we have that $X_{2,2} = \mathbf{G}_{2,2} \mathbf{U}_{2,2}$, which is interference-free according to (6). As a result, we have that $I(X_{2,2}; Y_2) = \text{rank}(\mathbf{F}_{2,2}) = m_{2,2}$. From here, we see that the proposed scheme can achieve the rate region of (7)-(9) with TIN.

In the subsequent sections, we focus on achieving the boundary points of this capacity region such that only the equalities of (7) and (9) are active. In this case, we set $m_{1,1} = n_1 - n_2$, $m_{1,2} = n_2 - m_{2,1} = m_1 + n_2 - n_1$, and $m_{2,2} = n_2$.

When $|h_1| < |h_2|$, the design of \mathbf{G}_1 and $\mathbf{G}_{2,1}$ can be easily obtained from (11) by swapping the argument between subscripts "1" and "2". The design of $\mathbf{G}_{2,2}$ remains unchanged.

The above deterministic approach allows us to obtain a systematic design on the input distributions for the original model, which we will show in the next section.

C. Proposed Coded Modulation Scheme With TIN

We translate the proposed scheme in Sec. III-B into the coded modulation scheme for the MAC under heterogeneous interference. We assume $|h_1| > |h_2|$ for illustrate purposes. The proposed scheme consists of the following steps.

1) *Encoding*: User k encodes its message into codeword \mathbf{c}_k by using a length- L_k rate- r_k binary code. Then, codeword \mathbf{c}_k is interleaved by employing bit-interleaved coded modulation (BICM) technique [17]. The interleaved codeword $\tilde{\mathbf{c}}_k$ is mapped to a length- N_k symbol sequence \mathbf{x}_k , where the modulation design and mapping steps will be described in Sec. III-C2 and Sec. III-C3, respectively. It is worth emphasizing that each user only uses a *single* channel code such that the encoding complexity is the same as for the single-user case.

2) *Modulation Design*: The proposed scheme for the deterministic model in Sec. III-B is systematically translated into the following QAM signaling

$$X_1[j] = \eta \sqrt{P_1} (\mathbf{F}_{1,2} + 2^{\frac{n_2}{2}} \mathbf{F}_{1,1}) \stackrel{\text{unif}}{\sim} \Lambda_1, \text{ for } j = 1, \dots, N_1, \quad (16)$$

$$X_2[j] = \begin{cases} \eta \sqrt{P_2} \cdot 2^{\frac{\log \text{SNR}_1 - m_{2,1} + n_2 - \log \text{SNR}_2}{2}} \mathbf{F}_{2,1} \stackrel{\text{unif}}{\sim} \Lambda_{2,1}, \\ \text{for } j = 1, \dots, N_1, \\ \eta' \sqrt{P_2} \mathbf{F}_{2,2} \stackrel{\text{unif}}{\sim} \Lambda_{2,2}, \text{ for } j = N_1 + 1, \dots, N_2, \end{cases} \quad (17)$$

where $\mathbf{F}_{k,t} \stackrel{\text{unif}}{\sim} \text{QAM}(2^{\text{rank}(\mathbf{F}_{k,t})}, 1)$ for $k, t \in \{1, 2\}$ and the rank number follows (11) and (15), Λ_1 is the superposition of two scaled regular QAM constellations with total cardinality 2^{m_1} whereas $\Lambda_{2,1}$ and $\Lambda_{2,2}$ are two scaled regular QAM constellations with cardinalities $2^{m_{2,1}}$ and $2^{m_{2,2}}$, respectively, η and η' are the normalization factors to ensure that $\mathbb{E}[|X_k[j]|^2] \leq P_k$ for $j = 1, \dots, N_k$. Specifically, the normalization factors are

$$\eta = \sqrt{\frac{1}{\max\{E_1, E_{2,1}\}}}, \quad (18)$$

$$\eta' = \frac{1}{\mathbb{E}[|\mathbf{F}_{2,2}|^2]} = \frac{6}{2^{n_2} - 1}, \quad (19)$$

where we define

$$E_1 \triangleq \mathbb{E} \left[\left| \mathbf{F}_{1,2} + 2^{\frac{n_2}{2}} \mathbf{F}_{1,1} \right|^2 \right] \quad (20)$$

$$= \frac{2^{n_2-m_{2,1}} - 1 + 2^{n_1} - 2^{n_2}}{6}, \quad (21)$$

$$E_{2,1} \triangleq \mathbb{E} \left[\left| 2^{\frac{\log \text{SNR}_1 - m_{2,1} + n_2 - \log \text{SNR}_2}{2}} \mathbf{F}_{2,1} \right|^2 \right] \quad (22)$$

$$= \frac{2^{n_2 + \log \text{SNR}_1 - \log \text{SNR}_2} (1 - 2^{-m_{2,1}})}{6}. \quad (23)$$

The power coefficient in front of each random variable \mathbf{F} is obtained by counting the number of rows below its corresponding submatrix \mathbf{F} in \mathbf{G} . For example, $2^{\frac{n_2}{2}}$ is due to that the number of rows below $\mathbf{F}_{1,1}$ in \mathbf{G}_1 is n_2 . For $\mathbf{F}_{2,1}$, the number of rows below $\mathbf{F}_{2,1}$ in $\mathbf{G}_{2,1}$ is $n_1 - m_{2,1} \approx \log \text{SNR}_1 - m_{2,1} + n_2 - \log \text{SNR}_2$. The reason for using these power coefficients will become clear in Sec. III-D.

For the case of $|h_1| < |h_2|$, the modulation design can be trivially obtained from (16) and (17) by swapping the argument between subscripts "1" and "2".

3) *Modulation Mapping*: From (2), we note that user 2's transmission block suffers from interference in the first N_1 symbols whereas the last $N_2 - N_1$ symbols are interference-free. Thus, user 2 uses two sets of constellations as shown in (17) to handle heterogeneous interference. Following the modulation design in (17), user 2 maps the first sub-block of the interleaved codeword $[\tilde{c}_2[1], \dots, \tilde{c}_2[N_1 m_{2,1}]]$ to $[x_2[1], \dots, x_2[N_1]]$ and the last sub-block $[\tilde{c}_2[N_1 m_{2,1} + 1], \dots, \tilde{c}_2[L_2]]$ to $[x_2[N_1 + 1], \dots, x_2[N_2]]$, respectively. Likewise, user 1 maps \mathbf{c}_1 to $[x_1[1], \dots, x_1[N_1]]$ following (16).

4) *TIN Decoding*: At the receiver, each user's message is decoded by treating other user's signals as noise. Specifically, the receiver starts to decode user 1's message upon receiving $[y[1], \dots, y[N_1]]$ while the decoding of user 2 starts upon receiving $[y[1], \dots, y[N_2]]$. Unlike SIC decoding, TIN decoding can be performed in parallel, which is more favorable if both users are URLLC users. The decoding process is the same as that in the point-to-point channel using BICM [17].

D. Performance Analysis

In this section, we analyze the performance of the proposed scheme in the Gaussian MAC with heterogeneous block-lengths. We first bound the mutual information of user 1 with TIN decoding as follows.

$$I(\mathbf{X}_1; \mathbf{Y}_1) = h(\mathbf{Y}_1) - h(\mathbf{Y}_1 | \mathbf{X}_1) \quad (24)$$

$$= h(\mathbf{Y}_1) - h(\mathbf{Z}) - (h(h_2 \mathbf{X}_{2,1} + \mathbf{Z}) - h(\mathbf{Z})) \quad (25)$$

$$= I(h_1 \mathbf{X}_1 + h_2 \mathbf{X}_{2,1}; \mathbf{Y}_1) - I(h_2 \mathbf{X}_{2,1}; h_2 \mathbf{X}_{2,1} + \mathbf{Z}) \quad (26)$$

$$\geq I(h_1 \mathbf{X}_1 + h_2 \mathbf{X}_{2,1}; \mathbf{Y}_1) - H(\mathbf{X}_{2,1}). \quad (27)$$

It remains to bound $I(h_1 \mathbf{X}_1 + h_2 \mathbf{X}_{2,1}; \mathbf{Y}_1)$. We note that $h_1 \mathbf{X}_1 + h_2 \mathbf{X}_{2,1}$ is a two-dimensional discrete constellation, whose achievable rate can be bounded by applying [14, Lemma 5] as

$$I(h_1 \mathbf{X}_1 + h_2 \mathbf{X}_{2,1}; \mathbf{Y}_1) \geq H(h_1 \mathbf{X}_1 + h_2 \mathbf{X}_{2,1}) - \log 2\pi e \left(\frac{1}{4} + \frac{4}{\pi d_{\min}^2(h_1 \mathbf{X}_1 + h_2 \mathbf{X}_{2,1})} \right). \quad (28)$$

We can also easily obtain user 2's mutual information under TIN $I(\mathbf{X}_{2,1}; \mathbf{Y}_1)$ from (27) whose lower bound is also determined by $d_{\min}(h_1 \mathbf{X}_1 + h_2 \mathbf{X}_{2,1})$ as in (28). Thus, a large $d_{\min}(h_1 \mathbf{X}_1 + h_2 \mathbf{X}_{2,1})$ is desirable for both users as it improves the mutual information lower bounds of both users under TIN decoding. With the proposed design in (16) and (17), we lower bound the minimum distance of $h_1 \mathbf{X}_1 + h_2 \mathbf{X}_{2,1}$ as follows.

$$d_{\min}(h_1 \mathbf{X}_1 + h_2 \mathbf{X}_{2,1}) \quad (29)$$

$$= \eta \cdot d_{\min} \left(2^{\frac{\log \text{SNR}_1}{2}} \mathbf{F}_{1,2} + 2^{\frac{\log \text{SNR}_1 - m_{2,1} + n_2}{2}} \mathbf{F}_{2,1} + 2^{\frac{\log \text{SNR}_1 + n_2}{2}} \mathbf{F}_{1,1} \right) \quad (30)$$

$$= \sqrt{\frac{1}{2^{-\log \text{SNR}_1} \max\{E_1, E_{2,1}\}}}$$

$$\times d_{\min} \left(\mathbf{F}_{1,2} + 2^{\frac{n_2 - m_{2,1}}{2}} \mathbf{F}_{2,1} + 2^{\frac{n_2}{2}} \mathbf{F}_{1,1} \right) \quad (31)$$

$$\geq \sqrt{\frac{6}{2^{n_1 - \log \text{SNR}_1}}} \quad (32)$$

$$\geq \sqrt{3}, \quad (33)$$

where in (32) we have used the fact that $\max\{E_1, E_{2,1}\} \leq \frac{2^{n_1}}{6}$ and $d_{\min}(\mathbf{F}_{1,2} + 2^{\frac{n_2 - m_{2,1}}{2}} \mathbf{F}_{2,1} + 2^{\frac{n_2}{2}} \mathbf{F}_{1,1}) = 1$ by applying Lemma 1 in Appendix A twice, and (33) is due to that $n_1 - \log \text{SNR}_1 \leq 1$. From (33), it can be seen that the proposed constellation and power coefficient designs in (16) and (17) guarantees that the superimposed constellation $h_1 \mathbf{X}_1 + h_2 \mathbf{X}_{2,1}$ has a constant minimum distance lower bound.

Note also that according to Lemma 1 in Appendix A, $\mathbf{F}_{1,2} + 2^{\frac{n_2 - m_{2,1}}{2}} \mathbf{F}_{2,1} + 2^{\frac{n_2}{2}} \mathbf{F}_{1,1}$ forms a regular QAM with zero mean, minimum distance 1, and cardinality 2^{n_1} . In this case, (28) can be refined to

$$I(h_1 \mathbf{X}_1 + h_2 \mathbf{X}_{2,1}; \mathbf{Y}_1) \geq H(h_1 \mathbf{X}_1 + h_2 \mathbf{X}_{2,1}) - \log \left(\frac{2\pi e}{12} \right) - \log \left(1 + \frac{12}{d_{\min}^2(h_1 \mathbf{X}_1 + h_2 \mathbf{X}_{2,1})} \right) \quad (34)$$

$$\stackrel{(33)}{=} m_1 + m_{2,1} - \log \left(\frac{5\pi e}{6} \right), \quad (35)$$

where (34) is obtained by following the proof of [14, Lemma 5] and setting the packing random variable therein $\mathbf{U} \stackrel{\text{unif}}{\sim} \mathbb{Z}^2$ with $\mathbf{G}_{\mathbb{Z}^2} = d_{\min}(h_1 \mathbf{X}_1 + h_2 \mathbf{X}_{2,1}) \mathbf{I}_2$, where \mathbf{I}_2 denotes a 2×2 identity matrix. Substituting (35) into (27) gives

$$I(\mathbf{X}_1; \mathbf{Y}_1) \geq m_1 - \log \left(\frac{5\pi e}{6} \right). \quad (36)$$

Following the steps from (24) to (36), we obtain the mutual information for user 2 under TIN decoding as

$$I(\mathbf{X}_{2,t}; \mathbf{Y}_t) \geq m_{2,t} - \log \left(\frac{5\pi e}{6} \right), \quad (37)$$

where $t \in \{1, 2\}$ is the sub-block index. Since the deterministic rate region $(m_1, m_{2,1}, m_{2,2})$ is within 1 bit to the corresponding Gaussian counterpart [16], the proposed scheme with TIN is capable of achieving the capacity region of the Gaussian MAC to within a *constant gap* for all channel parameters according to (36) and (37).

According to Theorem 1, the minimum distance affects both mutual information and dispersion. Unfortunately, deriving a closed-form bound for the dispersion term based on discrete constellations is very challenging. In the next section, we will show numerically that the proposed scheme leads to lower dispersion than that under Gaussian signaling.

IV. NUMERICAL RESULTS

In this section, we provide numerical results to show the finite blocklength achievable rates of the proposed scheme using QAM with TIN. We consider the channel model in Sec. II, where $(\text{SNR}_1, \text{SNR}_2) = (12, 24)$ dB, $(\epsilon_1, \epsilon_2) = (10^{-6}, 10^{-5})$,

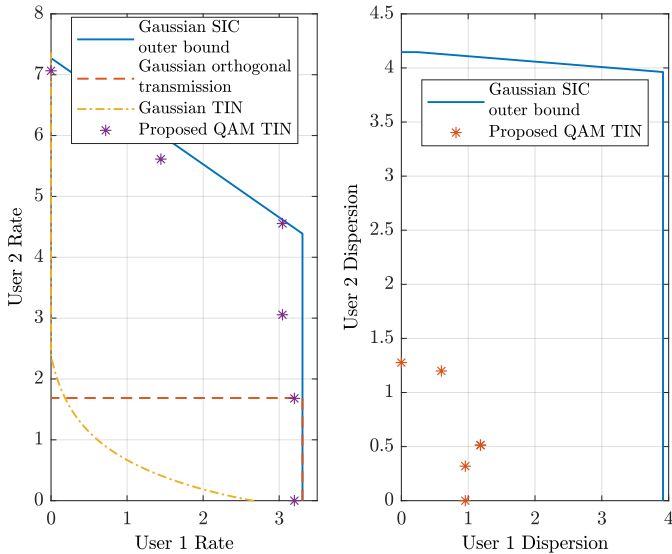


Fig. 1. (a) Achievable rate pairs (bits/s/Hz); (b) dispersion pairs (bits²/s/Hz).

and $(N_1, N_2) = (150, 200)$. The achievable rate pairs of the proposed scheme are computed by taking the first two terms of (3) and (4) and shown in Fig. 1(a). The proposed scheme uses the following modulation orders $(m_1, m_{2,1}, m_{2,2}) = (0, 8, 8), (2, 6, 8), (4, 4, 8), (4, 2, 8), (4, 0, 8), (4, 0, 0)$, which corresponds to the rate pairs orders from left to right in the figure.

In the same figure, we also include the achievable rates of three benchmark schemes based on Gaussian signaling and perfect SIC, orthogonal transmission, and TIN, respectively. For the first benchmark scheme based on perfect SIC, the rate region is generated by using the convex combination of all corner points. However, it should be noted that under finite blocklength, SIC can fail and time-sharing is not able to provide a convex region [18]. Thus, the rate region of the first benchmark scheme is an outer bound of the actual rate region under imperfect SIC and time-sharing, which is a hypothetical rate region used for comparison purposes only. It should be kept in mind that in the heterogeneous case, the messages of the more urgent user have to be decoded without waiting for the reception of the whole transmission block of the other user; therefore, SIC is infeasible in URLLC communication scenarios. The orthogonal transmission means that user 2 does not transmit any symbols for the first sub-block such that $x_2[j] = 0$ for $j = 1, \dots, N_1$. In this way, the transmission of both users do not interfere with each other.

Fig. 1(a) shows that the rate region of Gaussian TIN is the smallest. Intuitively, Gaussian interference lacks structures which may be difficult to be exploited when TIN is used. In contrast, thanks to the structural interference introduced by our carefully designed discrete input signaling, the proposed scheme with QAM and TIN can significantly outperform the second and third benchmark schemes based on Gaussian signaling. More importantly, the proposed scheme can perform very close to the benchmark scheme based on Gaussian

signaling and perfect SIC.

To further investigate the impacts of the proposed design on the second-order term of the achievable rate, we plot the dispersion pairs of the proposed scheme and those of the first benchmark scheme in Fig. 1(b). Since the second and third benchmark schemes perform poorly as shown in Fig. 1(a), there is no need to discuss their dispersions further. Observe that the proposed scheme has much smaller dispersion pairs than the benchmark scheme with perfect SIC. It is also worth emphasizing that the main features of URLLC communications are short blocklength and ultra-low error probability. This means that even a small increase in the dispersion can have a non-negligible impact on the achievable rate. Therefore, a small dispersion is desirable. To sum up, the superior performance of the proposed scheme with QAM and TIN is owing to close-to-capacity mutual information and smaller dispersion. As a result, we believe that the proposed scheme is promising for supporting heterogeneous services in the uplink.

V. CONCLUSION

In this paper, we have investigated the problem of the coexistence of heterogeneous services in the uplink. We have constructed a capacity-achieving coding scheme with TIN for the deterministic MAC with heterogeneous interference. Then, we have translated the aforementioned coding scheme into the coded modulation scheme for the Gaussian MAC with heterogeneous interference. We have proved that the mutual information of the proposed scheme with QAM and TIN is constant gap optimal. Numerical results have shown that the proposed scheme can achieve rate pairs close to the benchmark scheme based on Gaussian signaling and perfect SIC.

APPENDIX A

A USEFUL LEMMA FOR SUPERIMPOSED CONSTELLATIONS

Lemma 1. Consider a pair of regular QAM constellations (Λ_1, Λ_2) with zero mean and $d_{\min}(\Lambda_1) = d_{\min}(\Lambda_2) = \delta > 0$. The superimposed constellation $\Lambda_1 + 2^{\frac{\log |\Lambda_1|}{2}} \Lambda_2$ is a regular QAM with zero mean, $d_{\min}(\Lambda_1 + 2^{\frac{\log |\Lambda_1|}{2}} \Lambda_2) = \delta$, and cardinality $|\Lambda_1| \cdot |\Lambda_2|$.

Proof: We first look at the superposition of Λ_1 and Λ_2 in real parts by fixing the imaginary parts, i.e., $\Re(\Lambda_1) + 2^{\frac{\log |\Lambda_1|}{2}} \Re(\Lambda_2)$. For $k \in \{1, 2\}$, we have

$$\Re(\Lambda_k) = \left\{ -\frac{2^{\frac{\log |\Lambda_k|}{2}} - 1}{2} \delta, -\frac{2^{\frac{\log |\Lambda_k|}{2}} + 1}{2} \delta, \dots, \frac{2^{\frac{\log |\Lambda_k|}{2}} - 1}{2} \delta \right\}. \quad (38)$$

Clearly, for any two neighboring constellation points $\lambda_{k,1}, \lambda_{k,2} \in \Re(\Lambda_k)$ and $\lambda_{k,1} < \lambda_{k,2}$, we have

$$\lambda_{k,2} - \lambda_{k,1} = \delta. \quad (39)$$

Next, we treat $\mathfrak{R}(\Lambda_1)$ as a cluster and compute the inter-cluster distance of $\mathfrak{R}(\Lambda_1) + 2^{\frac{\log|\Lambda_1|}{2}}\mathfrak{R}(\Lambda_2)$ as

$$d(\mathfrak{R}(\Lambda_1) + \lambda_{2,1}, \mathfrak{R}(\Lambda_1) + \lambda_{2,2}) = -\min\{\mathfrak{R}(\Lambda_1)\} \\ = +2^{\frac{\log|\Lambda_1|}{2}}\lambda_{2,2} - \max\{\mathfrak{R}(\Lambda_1)\} - 2^{\frac{\log|\Lambda_1|}{2}}\lambda_{2,1} \quad (40)$$

$$= -\frac{2^{\frac{\log|\Lambda_1|}{2}}-1}{2}\delta \cdot 2 + 2^{\frac{\log|\Lambda_1|}{2}}\delta \quad (41)$$

$$=\delta, \quad (42)$$

where (41) follows by using (38) with $k = 1$ and (39) with $k = 2$. By (39) and (42), we know that for any pair of neighboring points $\lambda'_1, \lambda'_2 \in \mathfrak{R}(\Lambda_1) + 2^{\frac{\log|\Lambda_1|}{2}}\mathfrak{R}(\Lambda_2)$ and $\lambda'_1 < \lambda'_2$, we have

$$\lambda'_2 - \lambda'_1 = \delta. \quad (43)$$

As a result, the following holds.

$$\mathfrak{R}(\Lambda_1) + 2^{\frac{\log|\Lambda_1|}{2}}\mathfrak{R}(\Lambda_2) = \left\{ -\frac{2^{\frac{\log(|\Lambda_1| \cdot |\Lambda_2|) - 1}{2}}}{2}\delta, \right. \\ \left. -\frac{2^{\frac{\log(|\Lambda_1| \cdot |\Lambda_2|) + 1}{2}}}{2}\delta, \dots, \frac{2^{\frac{\log(|\Lambda_1| \cdot |\Lambda_2|) - 1}{2}}}{2}\delta \right\}. \quad (44)$$

Due to the symmetry property of the regular QAM, we also have that

$$\mathfrak{S}(\Lambda_1) + 2^{\frac{\log|\Lambda_1|}{2}}\mathfrak{S}(\Lambda_2) = \mathfrak{R}(\Lambda_1) + 2^{\frac{\log|\Lambda_1|}{2}}\mathfrak{R}(\Lambda_2). \quad (45)$$

Combining (44) and (45), we arrive at the conclusion stated in Lemma 1. ■

REFERENCES

- [1] H. Zhang, N. Liu, X. Chu, K. Long, A.-H. Aghvami, and V. C. M. Leung, "Network slicing based 5G and future mobile networks: Mobility, resource management, and challenges," *IEEE Commun. Mag.*, vol. 55, no. 8, pp. 138–145, 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] R. Kassab, O. Simeone, and P. Popovski, "Coexistence of URLLC and eMBB services in the C-RAN uplink: An information-theoretic study," in *Proc. IEEE Globecom*, 2018, pp. 1–6.
- [4] M. B. Shahab, R. Abbas, M. Shirvanimoghaddam, and S. J. Johnson, "Grant-free non-orthogonal multiple access for IoT: A survey," *IEEE Commun. Surveys Tut.*, vol. 22, no. 3, pp. 1805–1838, 2020.
- [5] O. Dizdar, Y. Mao, Y. Xu, P. Zhu, and B. Clerckx, "Rate-splitting multiple access for enhanced URLLC and eMBB in 6G," in *Proc. 17th Int. Symp. Wireless Commun. Syst. (ISWCS)*, Sep. 2021, pp. 1–6.
- [6] P.-H. Lin, S.-C. Lin, P.-W. Chen, M. Mross, and E. A. Jorswieck, "Rate region of Gaussian broadcast channels with heterogeneous blocklength constraints," in *Proc. IEEE Int. Conf. Commun. (ICC)*, May 2022, pp. 2144–2150.
- [7] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. New York, NY, USA: Wiley-Interscience, 2006.
- [8] Y. Polyanskiy, H. V. Poor, and S. Verdú, "Channel coding rate in the finite blocklength regime," *IEEE Trans. Inf. Theory*, vol. 56, no. 5, pp. 2307–2359, May 2010.
- [9] X. Sun, S. Yan, N. Yang, Z. Ding, C. Shen, and Z. Zhong, "Short-packet downlink transmission with non-orthogonal multiple access," *IEEE Trans. Wireless Commun.*, vol. 17, no. 7, pp. 4550–4564, Jul. 2018.
- [10] E. MolavianJazi and J. N. Laneman, "A second-order achievable rate region for Gaussian multi-access channels via a central limit theorem for functions," *IEEE Trans. Inf. Theory*, vol. 61, no. 12, pp. 6719–6733, 2015.
- [11] R. C. Yavas, V. Kostina, and M. Effros, "Gaussian multiple and random access channels: Finite-blocklength analysis," *IEEE Trans. Inf. Theory*, vol. 67, no. 11, pp. 6983–7009, 2021.
- [12] 3GPP, "5G;NR; Multiplexing and channel coding," 3rd Generation Partnership Project (3GPP), TR 38.212, Jan. 2022.
- [13] M. Qiu, Y.-C. Huang, S.-L. Shieh, and J. Yuan, "A lattice-partition framework of downlink non-orthogonal multiple access without SIC," *IEEE Trans. Commun.*, vol. 66, no. 6, pp. 2532 – 2546, Jun. 2018.
- [14] M. Qiu, Y.-C. Huang, and J. Yuan, "Discrete signaling and treating interference as noise for the Gaussian interference channel," *IEEE Trans. Inf. Theory*, vol. 67, no. 11, pp. 7253–7284, Nov. 2021.
- [15] M. Qiu, Y.-C. Huang, and J. Yuan, "Downlink transmission with heterogeneous URLLC services: Discrete signaling with single-user decoding," *IEEE J. Sel. Areas Commun.*, vol. 41, no. 7, pp. 2261–2277, 2023.
- [16] A. S. Avestimehr, S. N. Diggavi, and D. N. C. Tse, "Wireless network information flow: A deterministic approach," *IEEE Trans. Inf. Theory*, vol. 57, no. 4, pp. 1872–1905, 2011.
- [17] A. Guillén i Fàbregas, A. Martínez, and G. Caire, "Bit-interleaved coded modulation," *Found. Trends Commun. Inf. Theory*, vol. 5, no. 1–2, pp. 1–153, 2008.
- [18] V. Y. F. Tan and O. Kosut, "On the dispersions of three network information theory problems," *IEEE Trans. Inf. Theory*, vol. 60, no. 2, pp. 881–903, 2014.