

Impact of Interference Subtraction on Grant-Free Multiple Access with Massive MIMO

Lorenzo Valentini, Alberto Faedi, Marco Chiani, Enrico Paolini
CNIT/WiLab, DEI, University of Bologna, Italy

Email: {lorenzo.valentini13, alberto.faedi6, marco.chiani, e.paolini}@unibo.it

Abstract—The design of highly scalable multiple access schemes is a main challenge in the evolution towards future massive machine-type communications, where reliability and latency constraints must be ensured to a large number of uncoordinated devices. In this scenario, coded random access (CRA) schemes, where successive interference cancellation algorithms allow large improvements with respect to classical random access protocols, have recently attracted an increasing interest. Impressive performance can be potentially obtained by combining CRA with massive multiple input multiple output (MIMO). In this paper we provide an analysis of such schemes focusing on the effects of imperfect channel estimation on successive interference cancellation. Based on the analysis we then propose an innovative signal processing algorithm for CRA in massive MIMO systems.

I. INTRODUCTION

Next generation massive multiple access (MMA) protocols should be designed to achieve very high scalability (number of simultaneously active users the system can support) in presence of reliability and latency constraints [1]–[6]. In this respect, grant-free multiple access schemes have gained an increasing interest, owing to their capability to substantially reduce control signalling for connection establishment, which is beneficial in terms of both scalability and latency. Examples of grant-free schemes are the ones recently proposed in [7]–[12]. Uncoordinated protocols based on the CRA paradigm [13]–[20], a particular class of grant-free access schemes, ensure high reliability and are currently regarded as candidates for 6G [21] due to their capability of bridging random access with iterative decoding of codes on sparse graphs.

The performance of CRA schemes does not depend only on the medium access control (MAC) protocol; it also heavily relies on the effectiveness of the physical (PHY) layer processing algorithm. Although part of the literature on CRA tends to model the PHY layer signal processing (including packet detection, channel estimation, and interference cancellation) as ideal, signal processing in a realistic setting may introduce considerable losses with respect to the performance under idealized conditions. Moreover, especially in terrestrial scenarios, the often employed collision or multi-packet reception (MPR) channel models tend to be inaccurate, jeopardizing effective system design and optimization [14], [15], [22], [23].

In this paper, we investigate successive interference cancellation (SIC) algorithms for CRA in MMA applications. In particular, we review and in-depth analyze a low-complexity SIC processing proposed in [8] discussing its vulnerabilities.

Motivated by this analysis, we propose an innovative massive MIMO SIC technique able to improve scalability performance. The algorithm relies on the fact that, in CRA, it is possible to retrieve all packets sent by a user when a subset of them has been decoded. Exploiting this knowledge, it is possible to accurately estimate the channel coefficients, which are needed to subtract the interference. Due to imperfections of the SIC procedure in real scenarios, from now on we adopt the nomenclature successive interference subtraction (SIS), to emphasize the non-ideality of this step.

This paper is organized as follows. Section II introduces preliminary concepts, the system model, and some background material. Section III describes the proposed SIS technique along with an analysis which show the improvement of the proposed protocol. Numerical results are shown in Section IV. Finally, conclusions are drawn in Section V.

II. PRELIMINARIES AND BACKGROUND

In this section we define the reference scenario, including the MAC protocols and the channel model, also reviewing some physical layer signal processing techniques useful in the next sections. Throughout the paper, capital and lowercase bold letters denote matrices and vectors, respectively, $(\cdot)^H$ stands for conjugate transposition, $\|\cdot\|$ indicates the Euclidean norm, $\mathbb{E}\{\cdot\}$ denotes expectation, and $\mathbb{V}\{\cdot\}$ is used for variance.

A. Scenario Definition

We consider an MMA scenario with a very large number K of single-antenna transmitters (also referred to as users or devices), and one receiving base station (BS) equipped with multiple antennas. The K users are not all simultaneously active, since they are assumed to wake up unpredictably to transmit one data packet. We assume K_a out of the K users are active and the receiver has no prior knowledge of K_a .

We focus on grant-free MAC protocols to send uplink data from the active users to the BS. The schemes of interest belong to the class of coded slotted ALOHA (CSA) [15], which includes schemes using repetition codes such as contention resolution diversity slotted ALOHA (CRDSA) [13] and irregular repetition slotted ALOHA [14]. Variations on the packets transmission schedule have been proposed in [20]. In this paper we consider CSA with repetition codes of a given rate $1/r$ for all users. The time is organized in frames, each frame is divided into N slots, and users are frame- and slot-

synchronous. Hence, each active user generates r replicas of its packet and transmits them in r slots of the frame.

The availability of a BS with a massive number of antenna elements is a key feature to enable MPR at the receiver. In this respect, the use of orthogonal pilot sequences (or preambles) represents a simple approach to obtain MPR capabilities. Since in MMA K is typically much larger than the number of available pilots N_P , each active user picks one pilot randomly from the set of N_P preambles, without any coordination with the other active users. The use of CSA-based access and random pilot selection was proposed in [8]. Synchronization is achieved exploiting, for example, a beacon transmitted by the BS at the beginning of each frame.

Regarding the channel model, we consider a block Rayleigh fading channel with additive white Gaussian noise (AWGN). The channel coherence time is assumed equal to the slot duration T_s , which implies statistical independence of the channel coefficients of the same user across different slots. We do not consider shadowing effects owing to the assumption of perfect power control. Coherently with the above-mentioned access protocol and use of orthogonal pilots, each user active in a slot transmits a packet replica composed of one of the N_P orthogonal pilot sequences, of length N_P symbols, concatenated with a payload of length N_D symbols. Denoting the number of BS antennas by M , the signal received in a slot may be expressed as $[\mathbf{P}, \mathbf{Y}] \in \mathbb{C}^{M \times (N_P + N_D)}$ where

$$\mathbf{P} = \sum_{k \in \mathcal{A}} \mathbf{h}_k \mathbf{s}(k) + \mathbf{Z}_p, \quad \mathbf{Y} = \sum_{k \in \mathcal{A}} \mathbf{h}_k \mathbf{x}(k) + \mathbf{Z}. \quad (1)$$

In (1), \mathcal{A} is the set of users transmitting a replica in the considered slot, while $\mathbf{h}_k = (h_{k,1}, \dots, h_{k,M})^T \in \mathbb{C}^{M \times 1}$ is the vector of channel coefficients of the k -th user. The elements of \mathbf{h}_k are modeled as zero-mean, circularly symmetric, complex Gaussian independent and identically distributed (i.i.d.) random variables, i.e., $h_{k,i} \sim \mathcal{CN}(0, \sigma_h^2)$ for all $k \in \mathcal{A}$ and $i \in \{1, \dots, M\}$. Moreover, $\mathbf{s}(k) \in \mathbb{C}^{1 \times N_P}$ and $\mathbf{x}(k) \in \mathbb{C}^{1 \times N_D}$ are the orthogonal pilot sequence picked by user k in the current slot and the user payload, respectively, both with a unitary average energy per symbol. Finally, $\mathbf{Z}_p \in \mathbb{C}^{M \times N_P}$ and $\mathbf{Z} \in \mathbb{C}^{M \times N_D}$ are matrices whose elements are Gaussian noise samples. The elements of both \mathbf{Z}_p and \mathbf{Z} are i.i.d. random variables with distribution $\mathcal{CN}(0, \sigma_n^2)$. Due to power control, through the paper we adopt the normalization $\sigma_h^2 = 1$ for all users' channel coefficients.

B. Channel and Payload Estimation

As mentioned above, the BS receives a signal in the form $[\mathbf{P}, \mathbf{Y}]$ in each slot of the frame. The processing can be split into two phases [8], [20]. In the first one, the BS attempts channel estimation for all possible pilots by computing $\phi_j \in \mathbb{C}^{M \times 1}$, for all $j \in \{1, \dots, N_P\}$, as

$$\phi_j = \frac{\mathbf{P} \mathbf{s}_j^H}{\|\mathbf{s}_j\|^2} = \sum_{k \in \mathcal{A}^j} \mathbf{h}_k + \mathbf{z}_j \quad (2)$$

where \mathcal{A}^j is the set of active devices employing pilot j in the current slot, $\mathbf{s}_j \in \mathbb{C}^{1 \times N_P}$ is the j -th pilot sequence, and

$\mathbf{z}_j \in \mathbb{C}^{M \times 1}$ is a noise vector with i.i.d. $\mathcal{CN}(0, \sigma_n^2/N_P)$ entries. Note that in absence of noise, when pilot j is picked by a single user in the current slot, ϕ_j equals the vector of channel coefficients for that user.

In the second phase, the BS computes the quantities $\mathbf{f}_j \in \mathbb{C}^{1 \times N_D}$ and $g_j \in \mathbb{R}$ as

$$\begin{aligned} \mathbf{f}_j &= \phi_j^H \mathbf{Y} \\ &= \sum_{k \in \mathcal{A}^j} \left(\|\mathbf{h}_k\|^2 + \sum_{m \in \mathcal{A}^j \setminus \{k\}} \mathbf{h}_k^H \mathbf{h}_m \right) \mathbf{x}(k) \\ &+ \sum_{m \in \mathcal{A} \setminus \mathcal{A}^j} \left(\sum_{k \in \mathcal{A}^j} \mathbf{h}_k^H \mathbf{h}_m \right) \mathbf{x}(m) + \tilde{\mathbf{z}}_j \end{aligned} \quad (3)$$

and

$$\begin{aligned} g_j &= \|\phi_j\|^2 \\ &= \sum_{k \in \mathcal{A}^j} \left(\|\mathbf{h}_k\|^2 + \sum_{m \in \mathcal{A}^j \setminus \{k\}} \mathbf{h}_k^H \mathbf{h}_m \right) + \tilde{n}_j \end{aligned} \quad (4)$$

where $\tilde{\mathbf{z}}_j \in \mathbb{C}^{1 \times N_D}$ and \tilde{n}_j are noise terms. Then, the BS attempts estimation of the payload using conventional maximal ratio combining (MRC) as

$$\hat{\mathbf{x}} = \frac{\mathbf{f}_j}{g_j} = \frac{\phi_j^H \mathbf{Y}}{\|\phi_j\|^2}. \quad (5)$$

In the case where a generic user ℓ is the only one transmitting with pilot j in a given slot, hereafter referred to as singleton user ($\mathcal{A}^j = \{\ell\}$), we have $\hat{\mathbf{x}} \approx \mathbf{x}_\ell$. Upon successful channel decoding, the packet symbols are stored in a buffer waiting for the successive interference subtraction phase. The aim of this iterative processing, that will be explained in detail in the next section, is to subtract the interference of a packet in a slot using the information retrieved in another slot from one of its replicas. In fact, whenever a packet is successfully decoded, the BS acquires information about the positions of its replicas along with the employed preambles. This can be implemented in several ways, e.g., letting this information be a function of the information bits. This information can be used to subtract interference from a slot and attempt the decoding procedure again. Here, we separately computed \mathbf{f}_j and g_j for reasons that will be clear in Section III-A.

III. ANALYSIS OF SUCCESSIVE INTERFERENCE SUBTRACTION TECHNIQUES

In this section we present our main contributions. We first review in detail a state-of-the-art SIS technique for CSA with massive MIMO [8], discussing some critical points. Then, we present a theoretical analysis to evaluate the role of interference. Motivated by this analysis, we propose a SIS algorithm to improve the overall CSA scheme performance.

A. Squared-Norm-Based Interference Subtraction

Consider the low-complexity SIS algorithm, here indicated as squared norm based (SNB), proposed in [8] (also recently exploited in [20]). It relies on the assumption, whose validity is

analyzed and discussed later, that in a massive MIMO setting (3) and (4) can be approximated as

$$\mathbf{f}_j \approx \sum_{k \in \mathcal{A}^j} \|\mathbf{h}_k\|^2 \mathbf{x}(k) + \tilde{\mathbf{z}} \quad (6)$$

$$g_j \approx \sum_{k \in \mathcal{A}^j} \|\mathbf{h}_k\|^2 + \tilde{n} \quad (7)$$

respectively. Assume that we initially compute \mathbf{f}_j and g_j , $j = 1, \dots, N_P$, in all slots and that user ℓ is successfully decoded in a slot. Then, the above approximations lead naturally to the SIS procedure where we update $\mathbf{f}_j \leftarrow \mathbf{f}_j - \|\mathbf{h}_\ell\|^2 \mathbf{x}(\ell)$ and $g_j \leftarrow g_j - \|\mathbf{h}_\ell\|^2$ in all slots with replicas of the ℓ -th user. In other words, this algorithm subtracts only the main interfering term from (3) and (4). The update requires to know $\|\mathbf{h}_\ell\|^2$ in the replica slots where, due to the block fading assumption, the channel coefficients are different. For this issue, we can use the property that $\|\mathbf{h}_\ell\|^2/M$ tends to 1 for large M .

Importantly, the approximations (6) and (7) are not accurate when the cardinality of \mathcal{A} is large. In fact, since for $m \neq k$

$$\mathbb{E}\{\mathbf{h}_k^H \mathbf{h}_m\} = 0, \quad \mathbb{V}\{\mathbf{h}_k^H \mathbf{h}_m\} = M \quad (8)$$

the corresponding interfering terms in (3) and (4) may prevent from decoding a user packet even if it is the only one with a specific pilot. In the following we analyze this phenomenon by evaluating the probability that a user, being the only one with a specific pilot in a slot, is nevertheless not decoded.

B. Theoretical Analysis of the Interference Effects

We use the terminology ‘‘logical’’ to refer to an idealized setting in which: (i) whenever a user is the only one using a pilot in a given slot it is successfully decoded with probability one; (ii) channel estimation is perfect so that interference subtraction is ideal. Hereafter we provide a theoretical analysis of the effects of interference by removing hypotheses (i) and (ii), to understand their impact in a realistic setting.

Let us consider a situation where $|\mathcal{A}|$ users transmit simultaneously in a slot, $|\mathcal{A}^j|$ of them using pilot j . Assume $|\mathcal{A}^j| - 1$ users from the set \mathcal{A}^j have been successfully decoded in other slots. Then, in the current slot, we can apply SNB interference subtraction which, as mentioned above, mitigate but does not eliminate completely the interference. At this point, there is only one undecoded user adopting the j -th pilot (singleton). To analyze the probability that this user is successfully decoded, we focus on the interfering terms in (3). To highlight the effects of the interference we here neglect the noise contribution. Then, from (3) we can write

$$\mathbf{f}_j = \sum_{k \in \mathcal{A}^j} \|\mathbf{h}_k\|^2 \mathbf{x}(k) + \mathbf{I}_j \quad (9)$$

with $\mathbf{I}_j = \sum_{i=1}^{|\mathcal{A}^j| \cdot (|\mathcal{A}|-1)} \boldsymbol{\xi}_i$. Each term $\boldsymbol{\xi}_i$ is expressible as $\mathbf{h}_k^H \mathbf{h}_m \mathbf{x}$, where \mathbf{h}_k and \mathbf{h}_m are length- M vectors whose entries are modeled as i.i.d. $\mathcal{CN}(0, 1)$ random variables and \mathbf{x} is a length- N_D payload vector with i.i.d. entries. It follows

that each $\boldsymbol{\xi}_i$ is a vector whose generic entry fulfills

$$\mathbb{E}\{\boldsymbol{\xi}_i\} = 0, \quad \mathbb{V}\{\boldsymbol{\xi}_i\} = M. \quad (10)$$

We can therefore make the approximation

$$\mathbf{I}_j \approx \sum_{i=1}^{|\mathcal{A}^j| \cdot (|\mathcal{A}|-1)} \boldsymbol{\psi}_i \quad (11)$$

where $\boldsymbol{\psi}_i$ are independent random vectors with i.i.d. $\mathcal{CN}(0, M)$ entries. Due to subtraction of the interference generated by the $|\mathcal{A}^j| - 1$ users decoded in other slots, only one user remains using pilot j . Residues of imperfect interference cancellation are incorporated in \mathbf{I}_j , yielding a resulting interference term in the form

$$\tilde{\mathbf{I}}_j \approx \sum_{i=1}^{N_{\text{it}}} \boldsymbol{\psi}_i \quad (12)$$

where $N_{\text{it}} = |\mathcal{A}^j| \cdot |\mathcal{A}| - 1$ is the total number of interfering terms. Performing the estimation as in (5), we can write

$$\hat{\mathbf{x}}(\ell) = \mathbf{x}(\ell) + \frac{1}{M} \tilde{\mathbf{I}}_j \quad (13)$$

where the subscript ℓ denotes the only remaining user employing pilot j .

For a realistic analysis we also consider modulation and channel coding. For example, with a quadrature phase-shift keying (QPSK) constellation and hard-decision decoding, the symbol error probability is given by

$$P_e = \text{erfc}\left(\sqrt{\frac{M}{2N_{\text{it}}}}\right) - \frac{1}{4} \text{erfc}^2\left(\sqrt{\frac{M}{2N_{\text{it}}}}\right). \quad (14)$$

Finally, assume an error correcting code with bounded-distance decoding, able to correct up to t errors, and Gray QPSK constellation mapping. We can express the probability that decoding of a user packet is unsuccessful in a slot where its $|\mathcal{A}^j| - 1$ pilot-interferers are subtracted and a total of $|\mathcal{A}|$ users were initially allocated in the slot as

$$P_{\text{fail}} \approx 1 - \sum_{d=0}^t \binom{N_D}{d} P_e^d (1 - P_e)^{N_D-d} \quad (15)$$

where N_D is the number of payload symbols.

We report in Fig. 1 the analytical approximation (15) with P_e given by (14) when $N_D = 256$, $t = 10$, and $M = 256$. Moreover, we plot the corresponding curves obtained by numerical simulation to validate the derived result. Despite the approximations, the analytical results provide a good estimate, in terms of location along the horizontal axis, of the simulated curves. To improve the system performance in terms of average number of supported users for a given packet error probability of a singleton user, we can increase either the number of BS antennas M or the code error correction capability t for fixed N_D (which decreases the error correcting code rate).

In the particular case $|\mathcal{A}^j| = 1$, no interference subtraction is necessary and the user experiences the most favorable

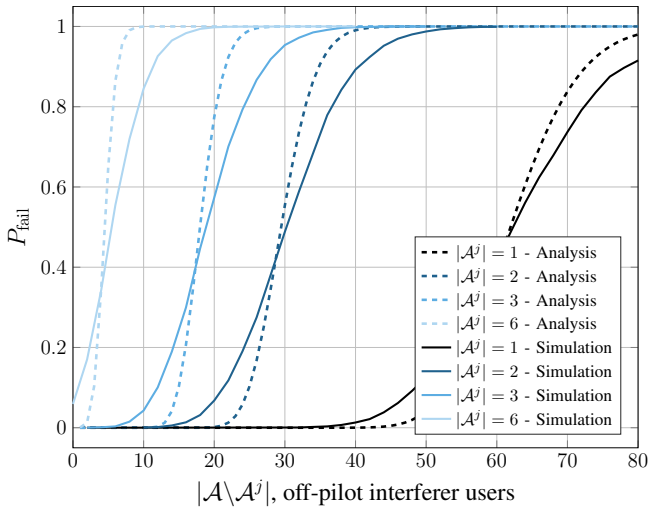


Fig. 1. Probability to unsuccessfully decode a singleton user after $|\mathcal{A}^j| - 1$ SNB iterations. Comparison between the analytical approximation and the simulation for $N_D = 256$, $t = 10$, and $M = 256$.

interference conditions. The $|\mathcal{A}^j| = 1$ curve in Fig. 1 reveals the actual performance of MRC payload estimation in (5) when interferers, using different orthogonal preambles, are captured in the model. Indeed, this is a major non-ideality, degrading the general performance of MAC protocols when a realistic channel model is accounted. On the other hand, when $|\mathcal{A}^j| > 1$, the estimation deteriorates even more, revealing the non-ideality of the SIS procedure. Moreover, we point out that, whenever a device using pilot j in the current slot is successfully decoded and SNB is performed, the interference on pilots different from j is not mitigated. This is the critical point of this SIS procedure and in Section III-C we will propose a technique able to overcome this problem.

C. Payload Aided Subtractions

Motivated by the analysis carried out in the previous subsection, we aim at changing the SIS algorithm to improve the overall performance. In repetition-based CSA, users send multiple copies of the same payload over the frame. Hereafter, we refer to the slots used to successfully decode a packet as “generator” slots.

Whenever a user ℓ is successfully decoded, the BS available information consists of the user’s payload of all packets, its channel coefficients in the generator slots (with an accuracy depending on the AWGN), the transmission slots, and the chosen pilots in each slot. Owing to this information, we can perform the update

$$\mathbf{P}^{(i+1)} = \mathbf{P}^{(i)} - \mathbf{h}_\ell \mathbf{s}(\ell), \quad \mathbf{Y}^{(i+1)} = \mathbf{Y}^{(i)} - \mathbf{h}_\ell \mathbf{x}(\ell) \quad (16)$$

in the generator slot, where we let $\mathbf{P}^{(0)} = \mathbf{P}$ and $\mathbf{Y}^{(0)} = \mathbf{Y}$. Regarding the replica slots, we exploit knowledge of the payload to estimate the channel coefficients as

$$\hat{\mathbf{h}}_\ell = \frac{\mathbf{Y} \mathbf{x}(\ell)^H}{\|\mathbf{x}(\ell)\|^2} = \mathbf{h}_\ell + \tilde{\mathbf{h}}_\ell$$

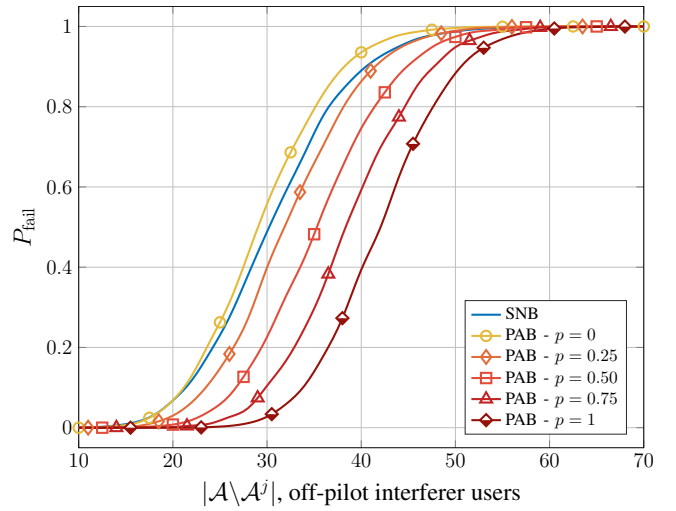


Fig. 2. Probability to unsuccessfully decode a singleton user after one SIS ($|\mathcal{A}^j| = 2$). Comparison between SNB and PAB for $N_D = 256$, $t = 10$, and $M = 256$.

$$= \mathbf{h}_\ell + \sum_{k \in \mathcal{A} \setminus \{\ell\}} \mathbf{h}_k \frac{\mathbf{x}(k) \mathbf{x}(\ell)^H}{\|\mathbf{x}(\ell)\|^2} + \mathbf{z}_h \quad (17)$$

where \mathbf{z}_h is the residual noise and \mathbf{Y} has not been modified yet by other interference subtractions. We can derive the statistical properties of the estimation error $\tilde{\mathbf{h}}_\ell$ given that the payload symbols are independent among users, as

$$\mathbb{E}\{\tilde{\mathbf{h}}_{\ell,n}\} = 0, \quad \mathbb{V}\{\tilde{\mathbf{h}}_{\ell,n}\} = \frac{|\mathcal{A}| - 1}{N_D} \quad (18)$$

where $n = 1, \dots, M$. As expected, increasing the number of payload symbols the accuracy of the channel coefficients estimation improves. We remark that using also knowledge of the preamble to perform channel estimation in slots where we wish to subtract interference may heavily deteriorate the estimation quality due to preamble collisions. Similarly to (16), we can now perform

$$\mathbf{P}^{(i+1)} = \mathbf{P}^{(i)} - \hat{\mathbf{h}}_\ell \mathbf{s}(\ell), \quad \mathbf{Y}^{(i+1)} = \mathbf{Y}^{(i)} - \hat{\mathbf{h}}_\ell \mathbf{x}(\ell) \quad (19)$$

in the replica slots. In this SIS algorithm, hereafter referred to as payload aided based (PAB), each time an update of the matrices \mathbf{P} and \mathbf{Y} has been carried out we re-compute (2) and (5) for each pilot in the current slot, to check if any other user can be successfully decoded after interference subtraction. In general, at the step $i = n_{\text{up}} + n_{\text{pa}}$ of the SIS algorithm, n_{up} subtractions are based on uncanceled pilots as from (16), and n_{pa} subtractions are based on payload aided channel coefficients estimation as from (19).

Fig. 2 illustrates the results of a variation of the experiment described in Section III-B for the two SIS techniques discussed in this paper and for $|\mathcal{A}^j| = 2$. More specifically, we assume that a fraction $0 \leq p \leq 1$ of users in the set $\mathcal{A} \setminus \mathcal{A}^j$ have been successfully decoded and subtracted. For them, we consider the worst case scenario ($n_{\text{up}} = 0$) where the SIS is performed using (19). As expected, the PAB performance improves as

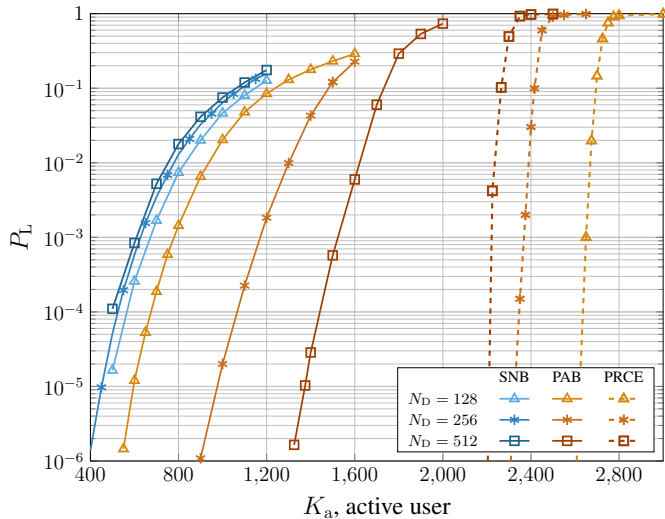


Fig. 3. Packet loss rate values of schemes characterized by different SIS techniques and payload sizes N_D . Baseline MAC with $N_P = 64$, $N = 78$. Comparison between the SNB, the proposed PAB and the ideal SIC case denoted as PRCE.

p increases. On the other hand, SNB is not influenced by p . As in a real scenario we have $n_{\text{up}} > 0$, the PAB technique is expected to outperform the SNB one; this is confirmed by the numerical results presented in the next section.

IV. PERFORMANCE EVALUATION

A. Simulation Setup

We present numerical results about the SIS techniques discussed in the previous sections using different MAC protocols. We consider a system where users transmit payloads encoded with an (n, k, t) narrow-sense binary Bose–Chaudhuri–Hocquenghem (BCH) code. A cyclic redundancy check (CRC) code is also used to validate decoded packets and avoid that the SIS procedure adds interference instead of subtracting it. Zero padding the BCH codeword with a final bit, we can map encoded bits onto a QPSK constellation with Gray mapping, obtaining N_D symbols per codeword. The QPSK symbol energy is normalized to 1. Simulations have been carried out with symbol rate $B_s = 1$ Msps, $M = 256$ BS antennas, and $\sigma_n^2 = 0.1$. We also impose a maximum latency constraint $\Omega = 50$ ms, leading to a number of slots per frame N equal to [20]

$$N = \left\lceil \frac{\Omega B_s}{2(N_P + N_D)} \right\rceil \quad (20)$$

where the number of orthogonal pilot symbols, N_P , equals the total number of available pilot sequences. These sequences are constructed using Hadamard matrices.

B. Numerical Results

We compare SIS techniques in terms of packet loss rate P_L for a given number K_a of active users in the frame. Regarding the MAC protocol, we adopt a standard repetition-based CSA protocol with a constant number r of replicas per packet

[13], referred to in the following as the “baseline MAC”. As a variation of the baseline protocol, we also show results for a second MAC protocol, namely, the recently proposed repetition-based CSA with intra-frame spatial coupling (SC) scheduling and acknowledgement (ACK) messages [20]. As a reference upper bound, we report the performance of a logical simulation. In this idealized setting a user transmitting alone in a “resource” (slot-pilot pair) uncollided is successfully decoded with probability one. Finally, as a second upper bound for the proposed scheme, we consider also a simulation which performs the PAB processing under the assumption that the subtractions are perfect (ideal SIC). In this scheme, denoted as perfect replica channel estimation (PRCE), the performance is limited by the payload estimation (5).

In Fig. 3 we report the packet loss rate (PLR) varying the symbol payload size N_D while keeping the rate of the BCH code constant, for the SNB, PAB, and PRCE techniques. To be precise, for $N_D \in \{128, 256, 512\}$ the corresponding BCH codes are $(255, 207, 6)$, $(511, 421, 10)$, and $(1023, 843, 18)$. In this particular example, we adopt the baseline MAC fixing $N = 78$ and $N_P = 64$ in order to show only the influence of N_D in the SIS processing. Looking at the curves of the SNB processing, we observe that the performance slightly degrades when N_D increases. The same behavior can be observed for PRCE. On the other hand, the PAB improves when N_D increases, as expected from (18). In addition, from Fig. 3 we can see that the gap between the PRCE and the PAB reduces, highlighting the effectiveness of the proposed technique in a complete scenario which accounts for both the PHY and MAC layer.

In Fig. 4 we plot the comparison between the SNB and the PAB imposing a maximum latency $\Omega = 50$ ms, $N_P = 64$ available pilots, repetition rate $r = 3$, a $(511, 421, 10)$ BCH code, using both the baseline MAC protocol and the SC with ACKs [20]. From these curves we observe that PAB remarkably improves the performance in comparison to SNB. For example, at $P_L = 10^{-4}$ SNB can support $K_a \approx 550$ users, while PAB more than twice $K_a \approx 1100$.

In Fig. 5 we extend Fig. 4 aiming to point out the gap between a real system represented by PAB and two idealized schemes, the PRCE and the logic one. As anticipated in Fig. 3, the distance between the PAB curve and the PRCE is mainly due to the channel estimation imperfections. On the other hand, the gap between the PRCE and the logical simulation is a consequence of the payload estimation non-ideality addressed in Section III-B. This plot reveals how neglecting the PHY layer processing in real scenario may lead to wrong conclusions and optimizations.

V. CONCLUSIONS

Interference poses a serious challenge for next generation grant-free massive multiple access. In this paper we provided an in-depth analysis of this problem for CRA schemes and proposed a massive MIMO interference subtraction processing able to ameliorate the scalability of state-of-the-art schemes, in the presence of reliability and latency constraints. For example,

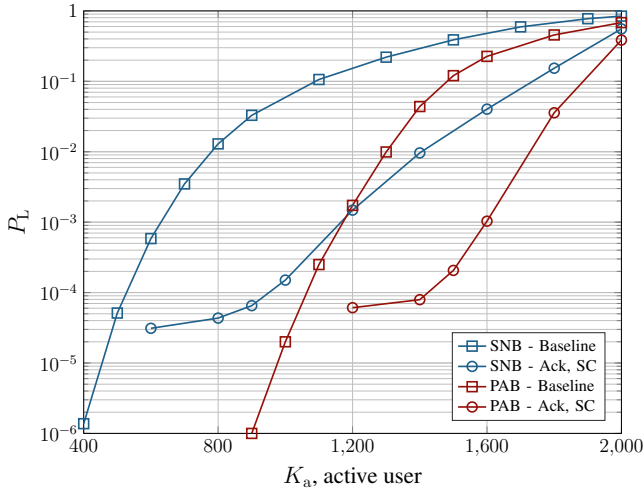


Fig. 4. Packet loss rate comparison between SNB and the proposed PAB, maximum latency $\Omega = 50$ ms, $N_P = 64$, $N = 78$, $N_D = 256$, for MAC protocols with or without spatial coupling (SC).

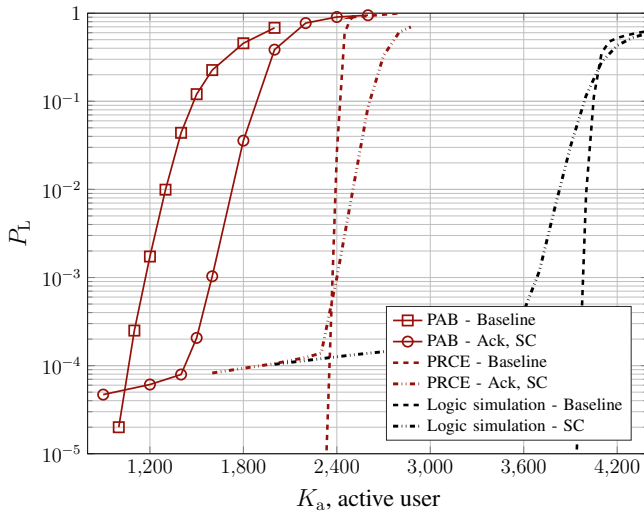


Fig. 5. Packet loss rate comparison between the proposed PAB and its bounds given by PRCE and logic simulation. Maximum latency $\Omega = 50$ ms, $N_P = 64$, $N = 78$, $N_D = 256$, for MAC protocols with or without spatial coupling (SC).

with a maximum latency of 50 ms and a target packet loss rate $P_L = 10^{-4}$, the proposed scheme is able to double the number of served users compared to the baseline. We also emphasized a large gap between results obtained under idealized and realistic condition, revealing how system design and analysis relying on collision-like channels may turn inaccurate.

ACKNOWLEDGEMENTS

The Authors would like to thank the anonymous Reviewers for comments and suggestions. This work has been carried out in the framework of the CNIT National Laboratory WiLab and the WiLab-Huawei Joint Innovation Center.

REFERENCES

- [1] M. Hasan, E. Hossain, and D. Niyato, "Random access for machine-to-machine communication in LTE-advanced networks: Issues and approaches," *IEEE Commun. Mag.*, vol. 51, no. 6, pp. 86–93, Jun. 2013.
- [2] L. Liu, E. G. Larsson, W. Yu, P. Popovski, C. Stefanovic, and E. De Carvalho, "Sparse signal processing for grant-free massive connectivity: A future paradigm for random access protocols in the internet of things," *IEEE Signal Process. Mag.*, vol. 35, no. 5, pp. 88–99, Sep. 2018.
- [3] X. Chen, D. W. K. Ng, W. Yu, E. G. Larsson, N. Al-Dhahir, and R. Schober, "Massive access for 5G and beyond," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 3, pp. 615–637, Mar. 2021.
- [4] G. Gui, M. Liu, F. Tang, N. Kato, and F. Adachi, "6G: Opening new horizons for integration of comfort, security, and intelligence," *IEEE Wireless Commun.*, vol. 27, no. 5, pp. 126–132, Oct. 2020.
- [5] C. Kalalal and J. Alonso-Zarate, "Massive connectivity in 5G and beyond: Technical enablers for the energy and automotive verticals," in *Proc. 2020 2nd 6G Wireless Summit*, Levi, Finland, Mar. 2020.
- [6] S. R. Pokhrel, J. Ding, J. Park, O.-S. Park, and J. Choi, "Towards enabling critical mMTC: A review of URLLC within mMTC," *IEEE Access*, vol. 8, pp. 131 796–131 813, Jul. 2020.
- [7] L. Liu and W. Yu, "Massive connectivity with massive MIMO—part I: Device activity detection and channel estimation," *IEEE Trans. Signal Process.*, vol. 66, no. 11, pp. 2933–2946, Mar. 2018.
- [8] J. H. Sørensen, E. De Carvalho, Č. Stefanovic, and P. Popovski, "Coded pilot random access for massive MIMO systems," *IEEE Trans. Wireless Commun.*, vol. 17, no. 12, pp. 8035–8046, Dec. 2018.
- [9] A. Fengler, S. Haghghatshoar, P. Jung, and G. Caire, "Grant-free massive random access with a massive MIMO receiver," in *2019 53rd Asilomar Conf. Signals, Systems, Computers*, Pacific Grove, CA, USA, Nov. 2019, pp. 23–30.
- [10] H. Han, Y. Li, W. Zhai, and L. Qian, "A grant-free random access scheme for M2M communication in massive MIMO systems," *IEEE Internet Things J.*, vol. 7, no. 4, pp. 3602–3613, Apr. 2020.
- [11] J. Choi, J. Ding, N. P. Le, and Z. Ding, "Grant-free random access in machine-type communication: Approaches and challenges," arXiv:2012.10550 [cs.IT], Dec. 2020.
- [12] A. T. Abebe and C. G. Kang, "MIMO-based reliable grant-free massive access with QoS differentiation for 5G and beyond," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 3, pp. 773–787, Mar. 2021.
- [13] E. Casini, R. De Gaudenzi, and O. del Rio Herrero, "Contention resolution diversity slotted ALOHA (CRDSA): An enhanced random access scheme for satellite access packet networks," *IEEE Trans. Wireless Commun.*, vol. 6, no. 4, pp. 1408–1419, Apr. 2007.
- [14] G. Liva, "Graph-based analysis and optimization of contention resolution diversity slotted ALOHA," *IEEE Trans. Commun.*, vol. 59, no. 2, pp. 477–487, Feb. 2011.
- [15] E. Paolini, G. Liva, and M. Chiani, "Coded slotted ALOHA: A graph-based method for uncoordinated multiple access," *IEEE Trans. Inf. Theory*, vol. 61, no. 12, pp. 6815–6832, Dec. 2015.
- [16] E. Paolini, Č. Stefanović, G. Liva, and P. Popovski, "Coded random access: Applying codes on graphs to design random access protocols," *IEEE Commun. Mag.*, vol. 53, no. 6, pp. 144–150, Jun. 2015.
- [17] F. Clazzer, C. Kissling, and M. Marchese, "Enhancing contention resolution ALOHA using combining techniques," *IEEE Trans. Commun.*, vol. 66, no. 6, pp. 2576–2587, Jun. 2018.
- [18] M. Berioli, G. Cocco, G. Liva, and A. Munari, "Modern random access protocols," *Foundations and Trends in Networking*, vol. 10, no. 4, pp. 317–446, 2016.
- [19] A. Munari, "Modern random access: An age of information perspective on irregular repetition slotted ALOHA," *IEEE Trans. Commun.*, vol. 69, no. 6, pp. 3572–3585, Jun. 2021.
- [20] L. Valentini, A. Faedi, M. Chiani, and E. Paolini, "Coded random access for 6G: Intra-frame spatial coupling with ACKs," in *Proc. 2021 IEEE Global Commun. Conf. Workshops*, Madrid, Spain, Dec. 2021.
- [21] N. H. Mahmood, H. Alves, O. A. López, M. Shehab, D. P. M. Osorio, and M. Latva-Aho, "Six key features of machine type communication in 6G," in *Proc. 2020 2nd 6G Wireless Summit*, Levi, Finland, Mar. 2020.
- [22] M. Ghanbarinejad and C. Schlegel, "Irregular repetition slotted ALOHA with multiuser detection," in *Proc. 2013 10th Annual Conf. Wireless On-demand Netw. Systems Services*, Banff, AB, Canada, Mar. 2013.
- [23] Č. Stefanović, E. Paolini, and G. Liva, "Asymptotic performance of coded slotted ALOHA with multipacket reception," *IEEE Commun. Lett.*, vol. 22, no. 1, pp. 105–108, Jan. 2018.