

Analysis and Optimization of Successful Symbol Transmission Rate for Grant-free Massive Access with Massive MIMO

Gang Chen, Ying Cui, Hei Victor Cheng, Feng Yang, and Lianghai Ding

Abstract—Grant-free massive access is an important technique for supporting massive machine-type communications (mMTC) for Internet-of-Things (IoT). Two important features in grant-free massive access are low-complexity devices and short-packet data transmission, making the traditional performance metric, achievable rate, unsuitable in this case. In this letter, we investigate grant-free massive access in a massive multiple-input multiple-output (MIMO) system. We consider random access control, and adopt approximate message passing (AMP) for user activity detection and channel estimation in the pilot transmission phase and small phase-shift-keying (PSK) modulation in the data transmission phase. We propose a more reasonable performance metric, namely successful symbol transmission rate (SSTR), for grant-free massive access. We obtain closed-form approximate expressions for the asymptotic SSTR in the cases of maximal ratio combining (MRC) and zero forcing (ZF) beamforming at the base station (BS), respectively. We also maximize the asymptotic SSTR with respect to the access parameter and pilot length.

I. INTRODUCTION

Grant-free massive access is an important technique for supporting massive machine-type communications (mMTC) for Internet-of-Things (IoT), which is one of the three main use cases for 5G. In grant-free massive access, there are two phases, i.e., the pilot transmission phase and the data transmission phase. A main technical challenge in massive access is to detect active users and estimate their channels in the pilot transmission phase in the presence of an excessive number of potential users. As only a small subset of users is active at any given time, the user activity detection and channel estimation problem can be modeled as a compressed sensing problem. Among the existing algorithms for compressed sensing, approximate message passing (AMP) algorithm is widely adopted, as it provides a good tradeoff between performance and complexity. In [1], [2], the authors adopt AMP for user activity detection and channel estimation in massive multiple-input multiple-output (MIMO) systems. The asymptotic performance of user activity detection and channel estimation is analyzed in [1], and the asymptotic achievable rate is analyzed in [2] (assuming perfect user activity detection). In [3], the authors propose channel-based access control and modified AMP for user activity detection, and analyze the performance of user activity detection. Note that in [1] and [3], performance analysis of the data transmission phase is not considered.

Two main features of data transmission in mMTC distinct it from data transmission in traditional human-type communications. Firstly, most data packets are *short*, i.e., usually

contain a few bytes. Secondly, low-complexity devices are used, and thus *small modulation* and *simple channel coding* are preferable. Thus, the achievable rate adopted in [2], which is an information-theoretic limit in the infinite blocklength regime, may not be a suitable performance metric for data transmission in mMTC. To the best of our knowledge, existing analytical results for data transmission cannot reflect the aforementioned features of mMTC. In addition, the authors in [2] optimize the pilot length to maximize the achievable rate for only one user activity realization, without considering the activity statistics, making the obtained pilot length less suitable for the case where the total number of active users has a large variance. Finally, the authors in [3] optimize the access control parameter to maximize the user identification performance, without considering the channel estimation and data transmission, making the obtained access control applicable only for limited scenarios.

In this letter, we would like to address the above issues. We study grant-free massive access in a massive MIMO system. We consider random access control, and adopt AMP for user activity detection and channel estimation. Considering low-complexity devices, we adopt small phase-shift-keying (PSK) modulation, e.g., BPSK and QPSK, for data transmission according to the standards [4]. In addition, considering transmission of short data packets, we propose a new performance metric, namely successful symbol transmission rate (SSTR), which reflects the performance of user activity detection and channel estimation in the pilot transmission phase and the performance of detection in the data transmission phase. The proposed SSTR is a more suitable performance metric for mMTC than the achievable rate [2], and its analysis is also more challenging. We first obtain closed-form approximate expressions for the asymptotic SSTR in the cases of maximal ratio combining (MRC) and zero forcing (ZF) beamforming at the base station (BS), respectively. The analytical results significantly facilitate the evaluation and optimization of the SSTR. Then, we maximize the asymptotic SSTR by optimizing the access parameter and pilot length. The optimization results provide practical guidelines for the design of mMTC systems. Finally, numerical results demonstrate the accuracy of the analysis and the importance of the optimization.

II. SYSTEM MODEL

Consider a massive access scenario arising from mMTC in a single cell with N users (devices) [1], [2], [5]. Let \mathcal{N} denote the set of all users. The BS is equipped with M antennas while each user is equipped with one antenna. We adopt a block-fading channel model where the channels within one coherence interval (CI) of length T symbols remain constant. We consider transmission in one CI, and denote the complex uplink channel vector from user n to the BS by $\mathbf{h}_n \in \mathbb{C}^{M \times 1}$.

G. Chen, Y. Cui, F. Yang, and L. Ding are with the Department of Electronic Engineering, Shanghai Jiao Tong University, Shanghai 200240, China (e-mail: yangfeng@sjtu.edu.cn).

H.V. Cheng is with The Edward S. Rogers Sr. Department of Electrical and Computer Engineering, University of Toronto, Toronto, ON M5S 3G4, Canada.

Assume $\mathbf{h}_n \sim \mathcal{CN}(\mathbf{0}, \gamma_n \mathbf{I}_M)$, where γ_n represents the path loss and shadowing component [1]. Assume that $\gamma_n, n \in \mathcal{N}$ are perfectly known at the BS, and that all users are perfectly synchronized. We consider random access control with access parameter ϵ . Within each CI, the users generate data with probability p_a , and access the channel with probability ϵ once they have data to send, both in i.i.d. manners. Thus, within each CI, the users send data via the channel (i.e., become active) with probability $p_a \epsilon$ in an i.i.d. manner. Note that p_a is a given system parameter, and ϵ is a design parameter for access control (controlling transmitting user sparsity) which will be optimized later. Denote by $\alpha_n \in \{1, 0\}$ the random activity state of user n with $\Pr[\alpha_n=1]=p_a \epsilon$.

We adopt a grant-free multiple-access scheme, where each user $n \in \mathcal{N}$ is assigned a unique pilot sequence with L symbols, denoted by $\mathbf{a}_n \triangleq (a_{n,1}, \dots, a_{n,L}) \in \mathbb{C}^{L \times 1}$. The pilot sequences and their correspondence to the user identities are known at the BS. In a massive access scenario, the pilot length is typically much smaller than the total number of users, i.e., $L \ll N$. Thus, it is not possible to assign mutually orthogonal pilot sequences to all N users. Note that L is a design parameter which will be optimized later. As in [1]–[3], [5], assume that for all $n \in \mathcal{N}$, the entries of \mathbf{a}_n are independently generated according to $\mathcal{CN}(0, 1/L)$. Each CI has two phases which will be illustrated below.

A. Pilot Transmission Phase

In the first phase, i.e., the pilot transmission phase, the active users synchronously send their pilot sequences to the BS. Therefore, the matrix of received signals at M antennas $\mathbf{Y}^{\text{pilot}} \in \mathbb{C}^{L \times M}$ is given by:

$$\mathbf{Y}^{\text{pilot}} = \sum_{n \in \mathcal{N}} \sqrt{L \rho_n^{\text{pilot}}} \alpha_n \mathbf{a}_n \mathbf{h}_n^T + \mathbf{Z}, \quad (1)$$

where $L \rho_n^{\text{pilot}}$ represents the transmit energy for the pilot sequence of user n , and $\mathbf{Z} \in \mathbb{C}^{L \times M}$ is the additive noise at the BS with each element following $\mathcal{CN}(0, \sigma^2)$. Denote $\mathbf{x}_n \triangleq \alpha_n \mathbf{h}_n \in \mathbb{C}^{M \times 1}, n \in \mathcal{N}$. The goal of the BS in the pilot transmission phase is to detect user activities and estimate the channels of active users by recovering $\mathbf{x}_n, n \in \mathcal{N}$ from the noisy observations $\mathbf{Y}^{\text{pilot}}$. As $p_a \epsilon \ll 1$, a lot of $\mathbf{x}_n, n \in \mathcal{N}$ are zero vectors. Thus, such a reconstruction problem is a compressed sensing problem. Following [1], this paper adopts a low-complexity AMP algorithm to recover $\mathbf{x}_n, n \in \mathcal{N}$, as it provides a good tradeoff between performance and computational complexity. For all $n \in \mathcal{N}$, based on the estimate $\hat{\mathbf{x}}_n$ of \mathbf{x}_n , the detected user activity $\hat{\alpha}_n \in \{0, 1\}$ can be obtained by hard-decision detection, and if $\hat{\alpha}_n = 1$, the estimated channel vector $\hat{\mathbf{h}}_n$ for \mathbf{h}_n is $\hat{\mathbf{x}}_n$. Denote $\Delta \mathbf{h}_n$ as the corresponding channel estimation error for each user n , i.e., $\mathbf{h}_n = \hat{\mathbf{h}}_n + \Delta \mathbf{h}_n$. Moreover, the convergence results of AMP provide the distributions of the estimates $\hat{\mathbf{x}}_n, n \in \mathcal{N}$ and estimation errors $\Delta \mathbf{x}_n \triangleq \mathbf{x}_n - \hat{\mathbf{x}}_n, n \in \mathcal{N}$.

B. Data Transmission Phase

In the second phase, i.e., the data transmission phase, the active users directly send their data to the BS using the remaining $T - L$ symbols. We adopt PSK modulation for

data transmission, e.g., BPSK and QPSK, as suggested in the standards [4]. Let s_n^W denote a W -array PSK symbol of user n with unit power, i.e., $\|s_n^W\|^2 = 1$, where $W \in \{2, 4, \dots\}$. Therefore, the received signal at the BS is expressed as:

$$\mathbf{y}^{\text{data}} = \sum_{n \in \mathcal{N}: \alpha_n=1} \sqrt{\rho_n^{\text{data}}} \mathbf{h}_n s_n^W + \mathbf{z}^{\text{data}}, \quad (2)$$

where ρ_n^{data} represents the transmit power for a data symbol of user n , and $\mathbf{z}^{\text{data}} \in \mathbb{C}^{M \times 1}$ is the additive noise at the BS with each element following $\mathcal{CN}(0, \sigma^2)$.

Based on the detected user activities and estimated channels, the BS tries to decode the data symbols of the users that are detected to be active using two linear receive beamforming strategies, namely MRC and ZF. Denote:

$$\hat{\mathbf{U}}^i \triangleq \begin{cases} \hat{\mathbf{G}}, & i = \text{MRC} \\ \hat{\mathbf{G}} (\hat{\mathbf{G}}^H \hat{\mathbf{G}})^{-1}, & i = \text{ZF} \end{cases}, \quad (3)$$

where $\hat{\mathbf{G}} \triangleq (\hat{\mathbf{h}}_n)_{n \in \mathcal{N}: \hat{\alpha}_n=1} \in \mathbb{C}^{M \times \hat{K}}$ with $\hat{K} \triangleq \sum_{n \in \mathcal{N}} \hat{\alpha}_n$ denoting the number of the users that are detected to be active. Let $\hat{\mathbf{u}}_n^i$ denote the column of $\hat{\mathbf{U}}^i$ that corresponds to user n with $\hat{\alpha}_n = 1$. Employing beamforming vector $\hat{\mathbf{u}}_n^i$, by (2) and $\mathbf{h}_n = \hat{\mathbf{h}}_n + \Delta \mathbf{h}_n$, we have:

$$\begin{aligned} \hat{r}_n^{i,W} &= \hat{\mathbf{u}}_n^{iH} \mathbf{y}^{\text{data}} \\ &= \hat{\mathbf{u}}_n^{iH} \left(\sum_{n \in \mathcal{N}: \alpha_n=1} \sqrt{\rho_n^{\text{data}}} (\hat{\mathbf{h}}_n + \Delta \mathbf{h}_n) s_n^W + \mathbf{z}^{\text{data}} \right) \\ &= \sqrt{\rho_n^{\text{data}}} \hat{\mathbf{u}}_n^{iH} \hat{\mathbf{h}}_n s_n^W + \hat{\mathbf{u}}_n^{iH} \sum_{n' \in \mathcal{N}: \alpha_{n'}=1, n' \neq n} \sqrt{\rho_{n'}^{\text{data}}} \hat{\mathbf{h}}_{n'} s_{n'}^W \\ &\quad + \hat{\mathbf{u}}_n^{iH} \sum_{n' \in \mathcal{N}: \alpha_{n'}=1} \sqrt{\rho_{n'}^{\text{data}}} \Delta \mathbf{h}_{n'} s_{n'}^W + \hat{\mathbf{u}}_n^{iH} \mathbf{z}^{\text{data}}. \end{aligned} \quad (4)$$

Then, the BS performs the minimum-distance detection on $\hat{r}_n^{i,W}$ by treating the term induced by channel estimation errors and interference from other users as additional noise, and obtains the estimated symbol $\hat{s}_n^{i,W}$ for user n with $\hat{\alpha}_n = 1$.

III. PERFORMANCE METRIC

In this letter, we use the SSTR, which represents the total number of symbols that can be correctly detected at the BS within a CI, as the performance metric for data transmission in grant-free massive access.

Definition 1: For given pilot length L and access parameter ϵ , the SSTR under the receive beamforming strategy i and the PSK modulation of size W is defined as:

$$\Phi^{(i,W)}(L, \epsilon) \triangleq \frac{T-L}{T} \mathbb{E} \left[\sum_{n \in \mathcal{N}} \mathbb{I}[\alpha_n=1, \hat{\alpha}_n=1, \hat{s}_n^W = s_n^W] \right], \quad (5)$$

where $\mathbb{I}[\cdot]$ represents the indicator function, and the expectation is taken over all sources of randomness.

Note that the SSTR captures user activity detection errors, channel estimation errors and data detection errors. The SSTR is a more suitable performance metric for grant-free massive access. However, in the general case, the analytical form of $\Phi^{(i,W)}(L, \epsilon)$ is not tractable, due to the complicated signal processing in grant-free massive access. Thus, as in [2], we focus

on the asymptotic case. Specifically, in Section III and Section IV, we consider the asymptotic analysis and optimization of the SSSTR at large M, N and L and high SNR under a simple power control policy, i.e., statistical channel inversion, which can reduce the channel gain differences between users, and is especially beneficial to users with relatively weaker channel gains [5].

With statistical channel inversion, $\rho_n^{\text{pilot}}, n \in \mathcal{N}$ and $\rho_n^{\text{data}}, n \in \mathcal{N}$ are chosen such that $\rho_n^{\text{pilot}} \gamma_n = \rho_n^{\text{data}} \gamma_n = \gamma, n \in \mathcal{N}$, where γ denotes the receive power for both pilot symbols and data symbols at each user. That is, the transmission powers of users scale inversely proportionally to their path-loss and shadowing components. With the same receive power, all users are statistically the same. Therefore, we can drop the user index n , and some dependence on $(\alpha_n)_{n \in \mathcal{N}}$ reduces to the dependence on the number of active users $K \triangleq \sum_{n \in \mathcal{N}} \alpha_n$. Note that K follows binomial distribution $\mathcal{B}(N, p_a \epsilon)$, i.e.,

$$\Pr[K=k] = C_N^k (p_a \epsilon)^k (1-p_a \epsilon)^{N-k} \triangleq q(N, k), \quad (6)$$

where $k=0 \cdots N$. When there are k active users and the pilot length is L , all k active users have the same average probability of missed detection, denoted by $p(k, L) \triangleq \mathbb{E}_{\mathbf{H}}[\Pr[\hat{\alpha}_n = 0 | \mathbf{H}, K=k, \alpha_n=1]]$, and the same average symbol error rate (SER) under receive beamforming strategy i and PSK modulation of size W , denoted by $\psi^{(i, W)}(k, L) \triangleq \mathbb{E}_{\mathbf{H}}[\Pr[\hat{s}_n^W \neq s_n^W | \mathbf{H}, K=k, \hat{\alpha}_n=\alpha_n=1]]$, where n represents the index of a typical active user, and $\mathbf{H} \triangleq (\mathbf{h}_n)_{n \in \mathcal{N}}$.

IV. ANALYSIS OF SSSTR

In this section, we derive an approximate expression of the asymptotic $\Phi^{(i, W)}(L, \epsilon)$ at large M, L, N and high SNR. In the regime of $L \leq k$ where AMP does not work, we assume that activity detection and data detection fail, i.e., $p(k, L) = 1$ and $\psi^{(i, W)}(k, L) = 1$. In the following, we focus on the asymptotic analysis of $p(k, L)$ and $\psi^{(i, W)}(k, L)$ in the regime of $k < L$. First, we use the asymptotic expression of $p(k, L)$ at large L, N, k and high SNR obtained in [1, Theorem 4] as an approximation for $p(k, L)$ at large L, N and high SNR and $k < L$.

Lemma 1: [1, Theorem 4] At large L, N and high SNR, for all $k < L$,

$$p(k, L) \approx \frac{\exp(-M(b(k, L) - 1 - \log(b(k, L))))}{2\sqrt{2\pi M}} \left((1-b(k, L))^{-1} + (\sqrt{2(b(k, L) - 1 - \log(b(k, L)))})^{-1} \right), \quad (7)$$

where $b(k, L) \triangleq \frac{\sigma^2}{\gamma(L-k)} \log \left(1 + \frac{\gamma(L-k)}{\sigma^2} \right)$.

Next, we derive an asymptotic approximation of $\psi^{(i, W)}(k, L)$.

Lemma 2: At large M, L, N and high SNR, for all $k < L$,

$$\psi^{(i, W)}(k, L) \approx \begin{cases} Q(\sqrt{2\Gamma^i(k, L)}), & W=2 \\ 2Q(\sqrt{\Gamma^i(k, L)}) - (Q(\sqrt{\Gamma^i(k, L)})^2), & W=4 \end{cases} \quad (8)$$

where $Q(x) = \frac{1}{\sqrt{2\pi}} \int_x^\infty \exp\left(-\frac{t^2}{2}\right) dt$, and

$$\Gamma^i(k, L) = \begin{cases} \frac{M\gamma^2}{(\gamma + \frac{\sigma^2}{L-k})(k\gamma + \sigma^2)}, & i = \text{MRC} \\ \frac{(M-k)(L-k)\gamma^2}{\sigma^2(\gamma L + \sigma^2)}, & M > k, i = \text{ZF} \end{cases}. \quad (9)$$

Proof: For notation simplicity, let C_{-n} denote the event that $\alpha_{n'} = \hat{\alpha}_{n'}, n' \in \mathcal{N}, n' \neq n$. At large M , we have:

$$\begin{aligned} & \Pr[\hat{s}_n^W \neq s_n^W | \mathbf{H}, K=k, \hat{\alpha}_n=\alpha_n=1] \\ &= \Pr[\hat{s}_n^W \neq s_n^W | \mathbf{H}, K=k, \hat{\alpha}_n=\alpha_n=1, C_{-n}] \Pr[C_{-n} | \mathbf{H}, K=k, \hat{\alpha}_n=\alpha_n=1] \\ & \quad + \Pr[\hat{s}_n^W \neq s_n^W | \mathbf{H}, K=k, \hat{\alpha}_n=\alpha_n=1, \overline{C_{-n}}] \\ & \quad \times \Pr[\overline{C_{-n}} | \mathbf{H}, K=k, \hat{\alpha}_n=\alpha_n=1] \\ & \stackrel{(a)}{\approx} \Pr[\hat{s}_n^W \neq s_n^W | \mathbf{H}, K=k, \hat{\alpha}_n=\alpha_n=1, C_{-n}], \end{aligned}$$

where (a) is due to $\Pr[C_{-n} | \mathbf{H}, K=k, \hat{\alpha}_n=\alpha_n=1] \rightarrow 1$ and $\Pr[\overline{C_{-n}} | \mathbf{H}, K=k, \hat{\alpha}_n=\alpha_n=1] \rightarrow 0$ as $M \rightarrow \infty$. Accordingly, we assume at large M , $\alpha_n = \hat{\alpha}_n, n \in \mathcal{N}$ [1], [2]. In the following, we analyze $\Pr[\hat{s}_n^W \neq s_n^W | \mathbf{H}, K=k, \hat{\alpha}_n=\alpha_n=1, C_{-n}]$ as an approximation of $\Pr[\hat{s}_n^W \neq s_n^W | \mathbf{H}, K=k, \hat{\alpha}_n=\alpha_n=1]$ at large M .

The SINR at a particular channel realization is $\tilde{\Gamma}_n^i(k, L) = \frac{\rho_n^{\text{data}} |\hat{\mathbf{u}}_n^{iH} \hat{\mathbf{h}}_n|^2}{F}$, where $F = \sum_{n' \in \mathcal{N}: \alpha_{n'}=1, n' \neq n} \rho_{n'}^{\text{data}} \mathbb{E}[|\hat{\mathbf{u}}_n^{iH} \hat{\mathbf{h}}_{n'}|^2] + \sum_{n' \in \mathcal{N}: \alpha_{n'}=1} \rho_{n'}^{\text{data}} \mathbb{E}[|\hat{\mathbf{u}}_n^{iH} \Delta \mathbf{h}_{n'}|^2] + \mathbb{E}[|\hat{\mathbf{u}}_n^{iH} \mathbf{z}^{\text{data}}|^2]$. At large M , in the regime of $k < M$,

$$\begin{aligned} & \tilde{\Gamma}_n^{\text{ZF}}(k, L) \\ & \stackrel{(b)}{=} \frac{\rho_n^{\text{data}}}{\mathbb{E} \left[\left[(\hat{\mathbf{G}}^H \hat{\mathbf{G}})^{-1} \right]_{nn} \right] \left(\sum_{n' \in \mathcal{N}: \alpha_{n'}=1} \rho_{n'}^{\text{data}} \mathbb{E}[\|\Delta \mathbf{h}_{n'}\|^2] + \mathbb{E}[\|\mathbf{z}^{\text{data}}\|^2] \right)} \\ & \stackrel{(c)}{\approx} \frac{(M-k)(L-k)\gamma^2}{\sigma^2(\gamma L + \sigma^2)} \triangleq \Gamma^{\text{ZF}}(k, L), \end{aligned}$$

where (b) is due to (3), and (c) is due to that $\hat{\mathbf{h}}_n \sim \mathcal{CN}(\mathbf{0}, \frac{\rho_n^{\text{pilot}}(L-k)\gamma_n^2}{\rho_n^{\text{pilot}}(L-k)\gamma_n + \sigma^2} \mathbf{I}_M)$ and $\Delta \mathbf{h}_n \sim \mathcal{CN}(\mathbf{0}, \frac{\gamma_n \sigma^2}{\rho_n^{\text{pilot}}(L-k)\gamma_n + \sigma^2} \mathbf{I}_M)$, as $M \rightarrow \infty$ [1], and $\hat{\mathbf{G}}^H \hat{\mathbf{G}} \sim \mathcal{W}_k \left(M, \frac{\rho_n^{\text{pilot}}(L-k)\gamma_n^2}{\rho_n^{\text{pilot}}(L-k)\gamma_n + \sigma^2} \mathbf{I}_M \right)$ [6]. At large M ,

$$\begin{aligned} & \tilde{\Gamma}_n^{\text{MRC}}(k, L) \stackrel{(d)}{\approx} \frac{\rho_n^{\text{data}} (\mathbb{E}[|\hat{\mathbf{h}}_n^H \hat{\mathbf{h}}_n|])^2}{D} \\ & \stackrel{(e)}{\approx} \frac{M\gamma^2}{(\gamma + \frac{\sigma^2}{L-k})(k\gamma + \sigma^2)} \triangleq \Gamma^{\text{MRC}}(k, L), \end{aligned}$$

where $D = \sum_{n' \in \mathcal{N}: \alpha_{n'}=1, n' \neq n} \rho_{n'}^{\text{data}} \mathbb{E}[|\hat{\mathbf{h}}_n^H \hat{\mathbf{h}}_{n'}|^2] + \sum_{n' \in \mathcal{N}: \alpha_{n'}=1} \rho_{n'}^{\text{data}} \mathbb{E}[|\hat{\mathbf{h}}_n^H \Delta \mathbf{h}_{n'}|^2] + \mathbb{E}[|\hat{\mathbf{h}}_n^H \mathbf{z}^{\text{data}}|^2] + \rho_n^{\text{data}} \mathbb{E}[|\hat{\mathbf{h}}_n^H \hat{\mathbf{h}}_n|^2] - \rho_n^{\text{data}} (\mathbb{E}[|\hat{\mathbf{h}}_n^H \hat{\mathbf{h}}_n|])^2$, (d) is from [7], and (e) is due to the distributions of $\hat{\mathbf{h}}_n$ and $\Delta \mathbf{h}_n$, as $M \rightarrow \infty$ [1]. Then, by [8], $\Pr[\hat{s}_n^W \neq s_n^W | \mathbf{H}, K=k, \hat{\alpha}_n=\alpha_n=1, C_{-n}] = \begin{cases} Q(\sqrt{2\Gamma^i(k, L)}), & W=2 \\ 2Q(\sqrt{\Gamma^i(k, L)}) - (Q(\sqrt{\Gamma^i(k, L)})^2), & W=4 \end{cases} \triangleq \tilde{\psi}^{(i, W)}(k, L)$. As $\psi^{(i, W)}(k, L) \approx \mathbb{E}_{\mathbf{H}}[\Pr[\hat{s}_n^W \neq s_n^W | \mathbf{H}, K=k, \hat{\alpha}_n=\alpha_n=1, C_{-n}]] = \tilde{\psi}^{(i, W)}(k, L)$, we complete the proof. ■

Based on Lemma 1 and Lemma 2, we obtain an approximate expression of $\Phi^{(i, W)}(L, \epsilon)$ at large M, N, L and high SNR.

Theorem 1: At large M, N, L and high SNR,

$$\begin{aligned} \Phi^{(i, W)}(L, \epsilon) & \approx \frac{T-L}{T} \sum_{k=1}^N k q(k) (1-p(k, L)) \left(1 - \tilde{\psi}^{(i, W)}(k, L) \right) \\ & \triangleq \tilde{\Phi}^{(i, W)}(L, \epsilon), \end{aligned}$$

where $p(k, L)$ is given by Lemma 1 and $\psi^{(i, W)}(k, L)$ is given

by Lemma 2.

Proof: We have:

$$\begin{aligned}
\Phi^{(i,W)}(L, \epsilon) &\stackrel{(a)}{=} \frac{T-L}{T} N \mathbb{E}_{\mathbf{H}} [\Pr[\alpha_n = 1, \hat{\alpha}_n = 1, \hat{s}_n^W = s_n^W | \mathbf{H}]] \\
&= \frac{T-L}{T} N \sum_{k=1}^N \mathbb{E}_{\mathbf{H}} [\Pr[\alpha_n = 1 | \mathbf{H}] \Pr[K = k | \mathbf{H}, \alpha_n = 1] \\
&\quad \times \Pr[\hat{\alpha}_n = 1 | \mathbf{H}, K = k, \alpha_n = 1] \Pr[\hat{s}_n^W = s_n^W | \mathbf{H}, K = k, \alpha_n = 1, \hat{\alpha}_n = 1]] \\
&\stackrel{(b)}{=} \frac{T-L}{T} N \Pr[\alpha_n = 1] \sum_{k=1}^N \Pr[K = k | \alpha_n = 1] \mathbb{E}_{\mathbf{H}} [\Pr[\hat{\alpha}_n = 1 | \mathbf{H}, \\
&\quad K = k, \alpha_n = 1] \Pr[\hat{s}_n^W = s_n^W | \mathbf{H}, K = k, \alpha_n = 1, \hat{\alpha}_n = 1]] \\
&\stackrel{(c)}{\approx} \frac{T-L}{T} N \Pr[\alpha_n = 1] \sum_{k=1}^N \Pr[K = k | \alpha_n = 1] \mathbb{E}_{\mathbf{H}} [\Pr[\hat{\alpha}_n = 1 | \mathbf{H}, \\
&\quad K = k, \alpha_n = 1]] \mathbb{E}_{\mathbf{H}} [\Pr[\hat{s}_n^W = s_n^W | \mathbf{H}, K = k, \alpha_n = 1, \hat{\alpha}_n = 1]] \quad (10) \\
&\stackrel{(d)}{=} \frac{T-L}{T} \sum_{k=1}^N k q(k) (1 - p(k, L)) (1 - \psi^{(i,W)}(k, L)),
\end{aligned}$$

where (a) is due to (5) and the statistical channel inversion, (b) is due to the independence between α , and \mathbf{H} , (c) is due that $\Pr[\hat{\alpha}_n = 1 | \mathbf{H}, K = k, \alpha_n = 1]$ and $\Pr[\hat{s}_n^W = s_n^W | \mathbf{H}, K = k, \alpha_n = 1, \hat{\alpha}_n = 1]$ become approximately independent at large M [1], and (d) is due to $\Pr[\alpha_n = 1] = p_a \epsilon$, $\Pr[K = k | \alpha_n = 1] = q(N-1, k-1)$, $\mathbb{E}_{\mathbf{H}} [\Pr[\hat{\alpha}_n = 1 | \mathbf{H}, K = k, \alpha_n = 1]] = 1 - p(k, L)$ and $\mathbb{E}_{\mathbf{H}} [\Pr[\hat{s}_n^W = s_n^W | \mathbf{H}, K = k, \alpha_n = 1, \hat{\alpha}_n = 1]] = 1 - \psi^{(i,W)}(k, L)$. ■

In Fig. 1, each analytical curve and the corresponding Monte-Carlo points indicate $\frac{T-L}{T} (1 - p(k, L)) (1 - \psi^{(i,W)}(k, L))$ and $\frac{T-L}{T} \mathbb{E}_{\mathbf{H}} [\Pr[\hat{\alpha}_n = 1 | \mathbf{H}, K = k, \alpha_n = 1] \Pr[\hat{s}_n^W = s_n^W | \mathbf{H}, K = k, \alpha_n = 1, \hat{\alpha}_n = 1]]$, respectively. From Fig. 1, we can see that each analytical curve and the corresponding Monte-Carlo points closely match. This demonstrates the accuracy of the approximations in (10), Lemma 1 and Lemma 2, and hence demonstrates the accuracy of Theorem 1. In Fig. 2, each analytical curve and the corresponding Monte-Carlo points indicate $\Phi^{(i,W)}(L, \epsilon)$ and $\tilde{\Phi}^{(i,W)}(L, \epsilon)$, respectively. The fact that each analytical curve and the corresponding Monte-Carlo points closely match further demonstrates the accuracy of Theorem 1. The computational complexity for evaluating $\tilde{\Phi}^{(i,W)}(L, \epsilon)$ is $O(N^3)$. The closed-form expression $\tilde{\Phi}^{(i,W)}(L, \epsilon)$ in Theorem 1 can be used for efficiently evaluating and optimizing the SSTR in practical systems.

From Lemma 1 and Lemma 2, we know that as M or SNR increases, $p(k, L)$ and $\psi^{(i,W)}(k, L)$ decrease, which results in the increment of $\tilde{\Phi}^{(i,W)}(L, \epsilon)$. Other system parameters influence $\tilde{\Phi}^{(i,W)}(L, \epsilon)$ in very complex manners, and their impacts have to be obtained using numerical evaluation. For example, from Fig. 2, we can see that when N, p_a, ϵ or L is small, $\tilde{\Phi}^{(i,W)}(L, \epsilon)$ increases with it and when N, p_a, ϵ or L is large, $\tilde{\Phi}^{(i,W)}(L, \epsilon)$ decreases with it. The reasons are as follows. As N, p_a or ϵ increases, on average, the number of users sending data (i.e., the number of transmitted data symbols) increases. When N, p_a or ϵ is small, the accuracy of user activity detection and channel estimation decreases slowly with N, p_a or ϵ , and hence $\tilde{\Phi}^{(i,W)}(L, \epsilon)$ increases with N, p_a or ϵ . When N, p_a or ϵ is large, the accuracy of user activity detection and channel estimation decreases fast with

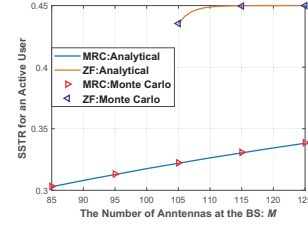
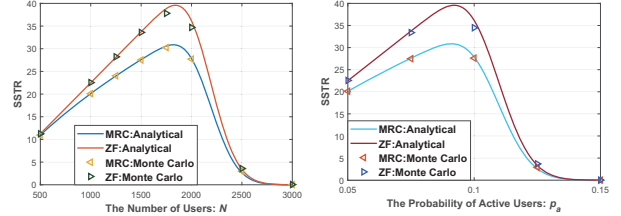
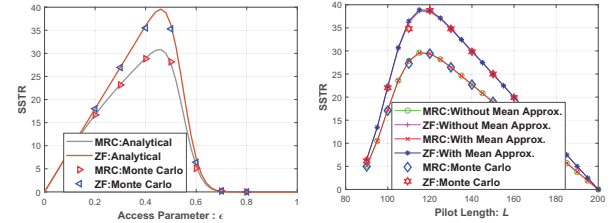


Fig. 1: SSTR for an active user at $N = 2000, k = 100, L = 110, T = 200$, SNR = 10dB and $W = 4$.



(a) SSTR versus N at $\epsilon = 0.5$, (b) SSTR versus p_a at $\epsilon = 0.5$, $L = 110, N = 2000$.



(c) SSTR versus ϵ at $L = 110$, (d) SSTR versus L at $\epsilon = 0.5$, $N = 2000, p_a = 0.1$.

Fig. 2: SSTR versus N, p_a, ϵ and L at $M = 128, T = 200$, SNR = 10dB and $W = 4$.

N, p_a or ϵ , and hence $\tilde{\Phi}^{(i,W)}(L, \epsilon)$ decreases with N, p_a or ϵ . In addition, a longer pilot length L leads to better user activity detection and channel estimation but fewer transmitted data symbols. When L is small, the accuracy of activity detection and channel estimation increases fast with L , and hence $\tilde{\Phi}^{(i,W)}(L, \epsilon)$ increases with L . When L is large, the accuracy of activity detection and channel estimation increases slowly with L , and hence $\tilde{\Phi}^{(i,W)}(L, \epsilon)$ decreases with L .

V. OPTIMIZATION OF SSTR

Fig. 2(c) and Fig. 2(d) indicate that it is important to carefully select the system design parameters ϵ and L so as to improve the SSTR. In this section, we consider the SSTR maximization with respect to ϵ and L .

A. Optimization of Access Parameter

In this part, we maximize the SSTR $\tilde{\Phi}^{(i,W)}(L, \epsilon)$ with respect to ϵ for given L :¹

$$g(L) \triangleq \max_{0 \leq \epsilon \leq 1} \tilde{\Phi}^{(i,W)}(L, \epsilon). \quad (11)$$

The problem in (11) is not in a convex form. By exploiting its structural properties, we have the following result.

Lemma 3: The optimization in (11) is equivalent to:

$$\begin{aligned}
g(L) &= \max_{\epsilon, t} \sum_{k=1}^N f(k, L) \epsilon^k t^{N-k} \\
\text{s.t.} \quad & 0 \leq p_a \epsilon + t \leq 1, \quad 0 \leq \epsilon \leq 1,
\end{aligned} \quad (12)$$

¹This problem is important for adjusting ϵ under abnormal conditions (e.g., p_a is far from its typical value).

where $f(k, L) = \frac{T-L}{T} C_N^k p_a^k k (1-p(k, L)) (1-\psi^{(i,W)}(k, L))$.

Proof: By Theorem 1, we have:

$$\tilde{\Phi}^{(i,W)}(L, \epsilon) = \sum_{k=1}^N f(k, L) \epsilon^k (1-p_a \epsilon)^{N-k}.$$

By introducing an auxiliary variable $t = 1 - p_a \epsilon$, the optimization in (11) can be equivalently transformed to:

$$\begin{aligned} \max_{\epsilon, t} \quad & \sum_{k=1}^N f(k, L) \epsilon^k t^{N-k} \\ \text{s.t.} \quad & t = 1 - p_a \epsilon, \quad 0 \leq \epsilon \leq 1. \end{aligned}$$

As $\sum_{k=1}^N f(k, L) \epsilon^k t^{N-k}$ is increasing in t , replacing the equality constraint $t = 1 - p_a \epsilon$ with the inequality constraint $t \leq 1 - p_a \epsilon$, i.e., $p_a \epsilon + t \leq 1$, in the optimization will not change the optimal solution (the inequality constraint is active at the optimal solution). In addition, as $t = 1 - p_a \epsilon$, we can add $t + p_a \epsilon \geq 0$ in the optimization without loss of optimality. Therefore, we complete the proof. ■

The optimization problem in (12) is a signomial geometric programming (SGP). A stationary point of it can be obtained using complementary geometric programming (CGP) [9]. We can run CGP multiple times, each with a random feasible initial point, and choose the stationary point with the largest objective value as a suboptimal solution of the optimization problem in (12). We omit the details due to page limitation. Fig. 2(c) demonstrates that the optimization with respect to ϵ for given L is of critical importance for SSTR improvement.

B. Optimization of Pilot Length

In this part, we maximize the SSTR $\tilde{\Phi}^{(i,W)}(L, \epsilon)$ with respect to L for given ϵ :²

$$\max_{L \in \{1, 2, \dots, T-1\}} \tilde{\Phi}^{(i,W)}(L, \epsilon). \quad (13)$$

This is a discrete optimization problem. Solving it requires computing $\tilde{\Phi}^{(i,W)}(L, \epsilon)$ (which is a sum of N terms) for all $L \in \{1, 2, \dots, T-1\}$. To reduce computational complexity, we adopt the mean approximation (i.e., approximating the expectation of a function of a random variable by the function of the expectation of the random variable) for $\tilde{\Phi}^{(i,W)}(L, \epsilon)$:

$$\begin{aligned} & \tilde{\Phi}^{(i,W)}(L, \epsilon) \\ &= \left(\sum_{k=1}^{L-1} q(k) \right) \frac{T-L}{T} \sum_{k=1}^{L-1} \frac{q(k)}{\sum_{k=1}^{L-1} q(k)} k (1-p(k, L)) (1-\psi^{(i,W)}(k, L)) \\ &\approx \frac{T-L}{T} \bar{K}_{<L} (1-p(\bar{K}_{<L}, L)) (1-\psi^{(i,W)}(\bar{K}_{<L}, L)) \sum_{k=1}^{L-1} q(k) \\ &= \frac{T-L}{T} (1-p(\bar{K}_{<L}, L)) (1-\psi^{(i,W)}(\bar{K}_{<L}, L)) \sum_{k=1}^{L-1} k q(k), \quad (14) \end{aligned}$$

where $\bar{K}_{<L} \triangleq \frac{\sum_{k=1}^{L-1} k q(k)}{\sum_{k=1}^{L-1} q(k)}$. Given the approximation of $\tilde{\Phi}^{(i,W)}(L, \epsilon)$ in (14), we only need to compute $p(\bar{K}_{<L}, L)$ and $\psi^{(i,W)}(\bar{K}_{<L}, L)$, and find the optimal L for given ϵ using exhaustive search (i.e., calculate $\tilde{\Phi}^{(i,W)}(L, \epsilon)$ for all $L \in \{1, 2, \dots, T-1\}$, and select L that achieves the maximum

among them). Fig. 2(d) shows that the error due to mean approximation is negligible. Fig. 2(d) also demonstrates that the optimization with respect to L for given ϵ is of great importance for SSTR improvement.

C. Joint Optimization of Pilot Length and Access Parameter

In this part, we jointly optimize L and ϵ to maximize the SSTR $\tilde{\Phi}^{(i,W)}(L, \epsilon)$:

$$\max_{0 \leq \epsilon \leq 1, L \in \{1, 2, \dots, T-1\}} \tilde{\Phi}^{(i,W)}(L, \epsilon), \quad (15)$$

which is equivalent to:

$$\max_{L \in \{1, 2, \dots, T-1\}} \max_{0 \leq \epsilon \leq 1} \tilde{\Phi}^{(i,W)}(L, \epsilon) = \max_{L \in \{1, 2, \dots, T-1\}} g(L),$$

where $g(L)$ is given by (11). Thus, we can solve the joint optimization problem in (15) based on the optimal solution of the problem in (11), and exhaustive search over $L \in \{1, 2, \dots, T-1\}$.

VI. CONCLUSION

In this letter, we investigated grant-free massive access in a massive MIMO system. We considered random access control, and adopted AMP for user activity detection and channel estimation in the pilot transmission phase and PSK modulation in the data transmission phase. We proposed a more reasonable performance metric, i.e., SSTR. We focused on the analysis and optimization of the asymptotic SSTR. Both analysis and optimization results offer important design insights for practical mMTC systems.

REFERENCES

- [1] 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, June 2018.
- [2] —, "Massive connectivity with massive MIMO Part II: Achievable rate characterization," *IEEE Trans. Signal Process.*, vol. 66, no. 11, pp. 2947–2959, June 2018.
- [3] Z. Sun, Z. Wei *et al.*, "Exploiting transmission control for joint user identification and channel estimation in massive connectivity," *IEEE Trans. Commun.*, pp. 1–1, 2019.
- [4] Y. E. Wang, X. Lin, A. Adhikary *et al.*, "A primer on 3gpp narrowband internet of things," *IEEE Commun. Mag.*, vol. 55, no. 3, pp. 117–123, March 2017.
- [5] K. Senel and E. G. Larsson, "Grant-Free massive MTC-Enabled massive MIMO: A compressive sensing approach," *IEEE Trans. Commun.*, vol. 66, no. 12, pp. 6164–6175, Dec 2018.
- [6] H. Q. Ngo, E. G. Larsson, and T. L. Marzetta, "Energy and spectral efficiency of very large multiuser MIMO systems," *IEEE Trans. Commun.*, vol. 61, no. 4, pp. 1436–1449, April 2013.
- [7] T. L. Marzetta and H. Yang, *Fundamentals of massive MIMO*. Cambridge University Press, 2016.
- [8] G. Proakis, John *et al.*, *Digital communications*. Mc-Graw-Hill, 2001.
- [9] M. Chiang *et al.*, "Geometric programming for communication systems," *Foundations and Trends® in Communications and Information Theory*, vol. 2, no. 1–2, pp. 1–154, 2005.

²This problem is important for the optimization of L without access control.