

# Variational Bayesian Inference Clustering Based Joint User Activity and Data Detection for Grant-Free Random Access in mMTC

Zhaoji Zhang, *Member, IEEE*, Qinghua Guo, *Senior Member, IEEE*, Ying Li, *Member, IEEE*,  
Ming Jin, *Member, IEEE*, and Chongwen Huang, *Member, IEEE*

**Abstract**—Tailor-made for massive connectivity and sporadic access, grant-free random access has become a promising candidate access protocol for massive machine-type communications (mMTC). Compared with conventional grant-based protocols, grant-free random access skips the exchange of scheduling information to reduce the signaling overhead, and facilitates sharing of access resources to enhance access efficiency. However, some challenges remain to be addressed in the receiver design, such as unknown identity of active users and multi-user interference (MUI) on shared access resources. In this work, we deal with the problem of joint user activity and data detection for grant-free random access. Specifically, the approximate message passing (AMP) algorithm is first employed to mitigate MUI and decouple the signals of different users. Then, we extend the data symbol alphabet to incorporate the null symbols from inactive users. In this way, the joint user activity and data detection problem is formulated as a clustering problem under the Gaussian mixture model. Furthermore, in conjunction with the AMP algorithm, a variational Bayesian inference based clustering (VBIC) algorithm is developed to solve this clustering problem. Simulation results show that, compared with state-of-art solutions, the proposed AMP-combined VBIC (AMP-VBIC) algorithm achieves a significant performance gain in detection accuracy.

**Index Terms**—Massive machine-type communications, grant-free, joint user activity and data detection, variational Bayesian inference, clustering, approximate message passing.

## I. INTRODUCTION

Internet of Things (IoT) facilitates information exchange among objects in the physical world, and motivates the development for a diversity of novel applications, such as the smart city, smart grid, factory automation, etc. As an important constituent scenario in 5G, massive machine-type communications (mMTC) has been proposed to accommodate

diversified IoT services [1]. Compared with conventional scenarios, mMTC is characterized by (i) the massiveness and low activation probability of user equipments (UEs), (ii) short data packets from activated UEs, and (iii) demand for low power consumption by low-cost UEs. Furthermore, these features will become more prominent with the evolution of 5G and 6G.

In medium access control (MAC) protocols, the random access mechanism configures connection setup for uplink transmission, i.e., the random access procedure allocates transmission resources to randomly activated UEs. However, the massiveness of UEs and shortage of uplink resources in mMTC have made random access a bottleneck problem for MAC designs [2], [3]. Existing random access schemes can be roughly divided into two categories, i.e. *grant-based* and *grant-free* schemes. In grant-based random access schemes, a handshaking procedure is needed to exchange the control signaling between the base station (BS) and active UEs. However, this handshaking procedure may incur prohibitively high signaling overhead for mMTC, which undermines the transmission efficiency of the small-sized data packets.

As an alternative to grant-based schemes, *grant-free* random access has emerged in recent years. In grant-free schemes, the handshaking procedure is skipped, while active UEs can share the uplink access resources, and directly transmit their data packets without the grant from the BS. To ensure successful data recovery under grant-free random access, several critical problems need to be addressed at the BS. For example, the BS needs to solve the user-activity detection (UAD) problem to identify the active UEs, as well as the channel estimation (CE) problem to obtain the channel state information (CSI) for these active UEs. After that, the BS needs to solve the multi-user detection (MUD) problem to detect the data from active UEs. Considering different enabling techniques for grant-free random access, the state-of-art solutions to above-mentioned problems are reviewed as follows.

### A. Grant-Free Random Access Enabled by MIMO and OFDM

As important enabling techniques for mMTC, the multiple input multiple output (MIMO) technique and orthogonal frequency division multiplexing (OFDM) technique can exploit spatial diversity and frequency diversity respectively to support the massive connectivity. On the other hand, the mobile traffic report [4] shows that only a small fraction of UEs will be activated in typical IoT applications. To exploit

Zhaoji Zhang and Ying Li are with the School of Telecommunications Engineering, Xidian University, Xi'an 710071, China (email: zhaojizhang@xidian.edu.cn; yli@mail.xidian.edu.cn).

Qinghua Guo is with the School of Electrical, Computer and Telecommunications Engineering, University of Wollongong, Wollongong, NSW 2522, Australia (e-mail: qguo@uow.edu.au)

Ming Jin is with the Faculty of Electrical Engineering and Computer Science, Ningbo University, Ningbo 315211, China (e-mail: jinming@nbu.edu.cn).

Chongwen Huang is with the College of Information Science and Electronic Engineering, Zhejiang University, Hangzhou 310027, China, also with the International Joint Innovation Center, Zhejiang University, Haining 314400, China, and also with the Zhejiang Provincial Key Laboratory of Information Processing, Communication and Networking (IPCAN), Hangzhou 310027, China (e-mail: chongwenhuang@zju.edu.cn).

(Corresponding Author: Ying Li)

this sparseness of active UEs, the framework of compressed sensing (CS) [5], [6] has received extensive studies. Under this CS framework, each UE is allocated with a unique pilot sequence, which will be transmitted with its data packet if this UE is activated. In this way, MIMO-enabled and OFDM-enabled grant-free random access share similar formulation of the detection problem, and the entire detection procedure is typically divided into two steps. Firstly, the *joint UAD and CE* problem is formulated as a *sparse-signal recovery problem*. Different CS algorithms have been proposed for this problem, such as the modified Bayesian CS algorithm [7], the block orthogonal matching pursuit (BOMP) algorithm [8], the approximate message passing (AMP) algorithm [9]–[11], the deep neural network-aided sparse Bayesian learning algorithm [12]. In the second step, the MUD problem can be readily addressed according to the UAD and CE results.

### B. Grant-Free Random Access Enabled by Spreading

The spreading technique serves as another enabling technique for mMTC with intriguing implementation feasibility. In spreading-enabled grant-free access mechanisms [13]–[17], each data symbol is spread with a UE-specific spreading sequence, while all the spread symbols of each UE experience the same *scalar* channel gain during transmission. In this way, the CE problem is significantly simplified, and spreading-enabled grant-free random access enjoys a much simpler problem formulation for receiver design. Then, different solutions have been proposed for the joint UAD and MUD problem. For example, an iterative order recursive least square (IORLS) algorithm [13] was proposed to exploit the joint sparsity of the data matrix to improve the detection accuracy. A joint expectation maximization and AMP (EM-AMP) algorithm was proposed in [14], where the data matrix is detected from the received signal by the AMP algorithm [18], while the activity detection is addressed by the EM algorithm [19]. In addition, a structured iterative support detection (SISD) algorithm is proposed in [15]. In [16], a block sparsity adaptive subspace pursuit (BSASP) algorithm is proposed for the joint UAD and MUD problem, while the CE problem is addressed with a reference symbol. Recently, a joint UAD, CE, and signal detection (JUICESD) algorithm was proposed in [17], where the AMP algorithm is employed for signal detection and the detected signals are also used to refine the CE result.

These above-mentioned solutions [13]–[17] involve some infeasible assumptions or deficiencies. For example, the sparsity level, i.e. the *exact* number of active UEs is assumed known to the BS in [13], while the schemes in [14], [15] require perfect knowledge on CSI at receiver (CSIR) even for inactive UEs. Such information is commonly unavailable in mMTC scenarios due to the massiveness and random activity of UEs. In addition, the subspace pursuit principle in [16] fails to address the inherent modulation constraint of data symbols, which undermines the data-detection accuracy. The UAD in [17] relies on a non-deterministic detection threshold, while fine-tuning this threshold may incur tedious work under complicated mMTC scenarios. Recently, some advances on MUD techniques have inspired new ideas to tackle these deficiencies, and the details are explained in the next subsection.

### C. Clustering and Variational Bayesian Inference for MUD

It is noted that modulated data symbols are discrete, while the received signals corrupted by fading and noise approximately follow the Gaussian distribution. Inspired by this fact, an unsupervised clustering approach is proposed in [20] for the joint CE and MUD problem. Specifically, the Gaussian-mixture model (GMM) is used to model the noise-corrupted received signals, where each cluster in the GMM is associated with one data symbol. Then, the EM algorithm is adopted for this clustering problem. However, the successive interference cancellation (SIC) principle is adopted for MUD in [20], which requires sufficiently large power difference among different users. For mMTC scenarios with densely deployed UEs, the received power of different UEs can be strongly correlated, which undermines the detection accuracy of SIC-based MUD. In addition, the variational Bayesian inference (VBI) method was employed for CE and MUD in one-bit quantized MIMO system [21]. With its powerful inference capability for intractable distributions, the VBI could effectively infer the distributions of the CSI and the data symbols from the received signals, which are heavily distorted after one-bit quantization.

### D. Motivations and Contributions

Intrigued by the implementation feasibility, we consider the spreading technique to enable grant-free random access for mMTC in this paper. In order to address the deficiencies of existing solutions and improve the detection accuracy, an AMP-combined variational Bayesian inference-based clustering (AMP-VBIC) algorithm is proposed for joint user activity and data detection. Specifically, the decoupling operations in the AMP framework are adopted to mitigate multi-user interference (MUI) and decouple the signals of different UEs. Given the decoupled signals, we first *extend* the data symbol alphabet to incorporate the null symbols from inactive UEs, and then formulate the joint user activity and data detection as a novel clustering problem under the GMM. Then, we develop a variational Bayesian inference based clustering (VBIC) algorithm for this clustering problem, where the CE result is also refined during the clustering procedure. The major contributions of this paper are summarized as follows.

- (i) With the extended symbol alphabet, the joint user activity and data detection is formulated as a clustering problem under GMM. Then, we derive the VBIC algorithm for this clustering problem, which iteratively works in conjunction with the AMP decoupling module to refine the detection accuracy.
- (ii) In the VBIC algorithm, the CE result is iteratively updated with the clustering result of all the data symbols, which in return improves the UAD and MUD accuracy.
- (iii) Analyses are provided to demonstrate the favorable linear complexity of the proposed AMP-VBIC algorithm, while simulation results show its superior detection accuracy over the state-of-art solutions.

The remainder of this paper is organized as follows. Section II describes the system model, and the AMP-VBIC algorithm is proposed in Section III for the joint user activity and data detection problem. Simulation results are provided in Section IV, and Section V concludes this paper.

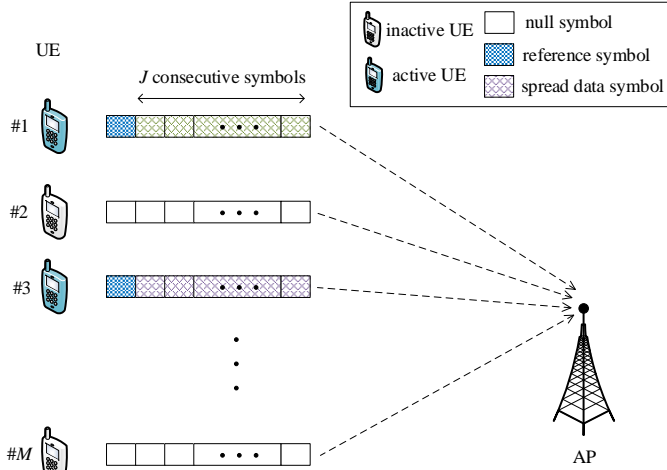


Fig. 1. Grant-free random access system model.

*Notations:* Scalar variables are written in italic letters. Vectors (or a set of variables) are written in boldface lower-case letters, and matrices in boldface upper-case letters. Unless stated otherwise, all the vectors are column vectors.  $(\cdot)^T$  and  $(\cdot)^H$  are the transpose and conjugate-transpose operations, respectively.  $E[\cdot]$  and  $\text{Var}[\cdot]$  take the expectation and variance of a random variable, respectively.  $X \sim \mathcal{CN}(\mu, v)$  means that a random variable  $X$  follows a complex Gaussian distribution with mean  $\mu$  and variance  $v$ , and  $\mathcal{CN}(x|\mu, v)$  is the probability density function (pdf) of this complex Gaussian distribution.

## II. SYSTEM MODEL

As shown in Fig. 1, we consider a spreading-based uplink grant-free random access system with an access point (AP) serving  $M$  user-equipments (UEs). Each UE is randomly activated with a probability  $p_a$ , while each active UE transmits  $J$  consecutive symbols in one transmission block. The  $j$ -th modulated symbol of the  $m$ -th UE is denoted by  $d_{[m]}^j$ , which will be spread over a time-spreading sequence  $\mathbf{a}_m$  of length  $N$  before transmission. Then, the  $j$ -th received-signal vector  $\mathbf{y}^j$  at the AP can be represented as

$$\mathbf{y}^j = \mathbf{A}\mathbf{U}\mathbf{d}^j + \mathbf{w}^j, \quad (1)$$

where  $\mathbf{y}^j$  is the received signal vector with length  $N$ , and  $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_M]$  is a  $N \times M$  spreading matrix for all the  $M$  UEs,  $\mathbf{U} = \text{diag}\{\mu_1, \mu_2, \dots, \mu_M\}$  is a diagonal matrix, and the  $m$ -th diagonal element  $\mu_m \sim \mathcal{CN}(0, 1)$  represents the Rayleigh channel coefficient of the  $m$ -th UE. The term  $\mathbf{d}^j = [d_{[1]}^j, d_{[2]}^j, \dots, d_{[M]}^j]^T$  is the  $j$ -th transmitted symbol vector of  $M$  UEs, and  $\mathbf{w}^j$  represents the additive white Gaussian noise (AWGN) vector with length  $N$ .

We assume an overloaded system with a large number of UEs, i.e.  $N < M$ . However, due to the sporadic activation of UEs in grant-free random access, there are only a small number of active UEs in each transmission frame. To facilitate the joint activity and data detection, we introduce an *extended symbol alphabet*  $\Delta = \{d_1, \Delta_a\}$  for the transmitted symbols

$d_{[m]}^j$ . Here,  $d_1 = 0$  represents the equivalent *null symbol* from inactive UEs,  $\Delta_a$  is the modulation symbol alphabet of active UEs. For example, if Quadrature Phase Shift Keying (QPSK) modulation is adopted for transmission, we have  $\Delta_a = \{\frac{1+1j}{\sqrt{2}}, \frac{1-1j}{\sqrt{2}}, \frac{-1+1j}{\sqrt{2}}, \frac{-1-1j}{\sqrt{2}}\}$ , where  $j = \sqrt{-1}$ . Furthermore, we denote  $K$  as the size of  $\Delta$ , i.e.  $\Delta = \{d_1, d_2, \dots, d_K\}$ .

Then, we consider the block transmission of  $J$  consecutive symbols, and obtain a matrix version of (1) as

$$\mathbf{Y} = \mathbf{A}\mathbf{U}\mathbf{D} + \mathbf{W} = \mathbf{A}\mathbf{X} + \mathbf{W}, \quad (2)$$

where  $\mathbf{Y} = [\mathbf{y}^1, \mathbf{y}^2, \dots, \mathbf{y}^J]$  is the received signal matrix of size  $N \times J$ ,  $\mathbf{D} = [\mathbf{d}^1, \mathbf{d}^2, \dots, \mathbf{d}^J]$  is the transmitted signal matrix of size  $M \times J$ , and  $\mathbf{W}$  is an AWGN matrix of size  $N \times J$ . The spreading matrix  $\mathbf{A}$  is known to the AP, and we assume a quasi-static block fading channel, i.e. the channel matrix  $\mathbf{U}$  remains unchanged over the entire block of  $J$  symbols.

It is noted that  $\mathbf{U}$  is unknown to the AP, and we define the *intermediate detection target*  $\mathbf{X}$  as  $\mathbf{X} = \mathbf{U}\mathbf{D}$  in (2), from which the decision on  $\mathbf{D}$  should be obtained. Since the data symbols in  $\Delta_a$  are usually *symmetric* for active UEs, we need to correct the *phase ambiguity* when recovering  $\mathbf{D}$  from  $\mathbf{X}$ . As shown in Fig. 1, we adopt a common solution to this phase ambiguity problem [16], i.e. inserting a reference symbol (RS) before the data symbols. [More details on correcting this phase ambiguity problem will be later explained in Remark 1 of Section III.](#) In addition, each inactive UE equivalently transmits  $J$  null symbols, i.e.  $d_{[m]}^j = d_1 = 0$  for  $j \in \{1, 2, \dots, J\}$ . In this way, both  $\mathbf{X}$  and  $\mathbf{D}$  exhibit the *row sparsity*. That is, the rows of  $\mathbf{X}$  and  $\mathbf{D}$  corresponding to inactive UEs only have zero elements, while the nonzero elements only reside in the rows corresponding to active UEs. The above-mentioned constraint is dubbed the *joint sparsity* for the elements in  $\mathbf{X}$  and  $\mathbf{D}$ , which will be used for activity detection. More details are explained as follows.

## III. VARIATIONAL BAYESIAN INFERENCE CLUSTERING FOR JOINT USER ACTIVITY AND DATA DETECTION

To address the joint UAD and MUD problem, we derive the following AMP-VBIC algorithm. Typically, the operations in the AMP algorithm are divided into two modules, i.e. the *decoupling module* which solves a linear mixing problem and a *denoiser module* which usually functions as a demodulator for the data-detection target. However, it is shown in (2) that both the data matrix  $\mathbf{D}$  and the unknown channel matrix  $\mathbf{U}$  are included in the intermediate detection target  $\mathbf{X}$ . As a result, the demodulator in the typical AMP framework is not applicable to the detection of  $\mathbf{X}$  under our model. As an alternative, we design the AMP-VBIC algorithm, where the denoiser module is now replaced with our proposed VBI clustering module. In this way, the VBI clustering module works in conjunction with the AMP decoupling module for the joint detection problem. The information exchange diagram between these two modules is illustrated in Fig. 2. More details are explained as follows.

### A. Pseudo Observation From AMP Decoupling Module

For the linear mixing problem in (2) with known spreading matrix  $\mathbf{A}$ , the decoupling operations of the AMP algorithm

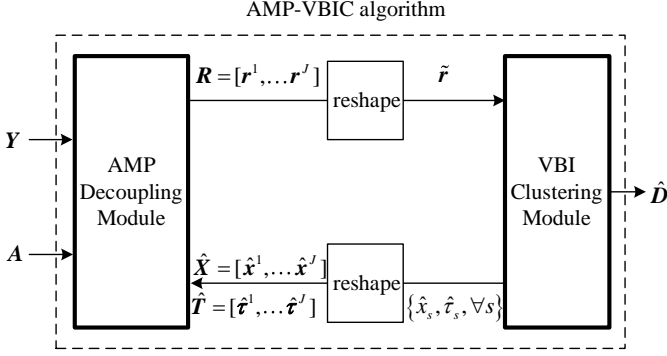


Fig. 2. Information exchange between the AMP decoupling module and the VBI clustering module.

(i.e. the column-by-column operations in Algorithm 1) can be readily adopted to decouple the *intermediate detection target*  $\mathbf{X}$  from the received signal matrix  $\mathbf{Y}$ . That is, at the output of the AMP decoupling module, we can obtain a pseudo-observation matrix  $\mathbf{R}$  for  $\mathbf{X}$ . Specifically, denote  $r_m^j$  and  $x_m^j$  as the element in the  $m$ -th row and  $j$ -th column of  $\mathbf{R}$  and  $\mathbf{X}$ , respectively. The pseudo observation  $r_m^j$  of the target element  $x_m^j$  is written as

$$r_m^j = x_m^j + n_m^j, \quad (3)$$

where the observation noise  $n_m^j$  follows the distribution  $n_m^j \sim \mathcal{CN}(0, \tau_m^j)$ . Both  $r_m^j$  and  $\tau_m^j$  are provided by the AMP decoupling module.

For each target element  $x_m^j$ , we have  $x_m^j = \mu_m d_{[m]}^j$ . Since  $\mu_m$  is unknown to the AP, the typical denoiser in the AMP algorithm (i.e. the demodulator) fails to demodulate  $d_{[m]}^j$  from the observation  $r_m^j$  of  $x_m^j$ . However, it is noted that each observation  $r_m^j$  is associated with one specific data symbol  $d_k \in \Delta$ . Therefore, we can use the Gaussian mixture model, and cluster these observations by the following VBI-based clustering (VBIC) algorithm. After that, we update the mean  $\hat{x}_m^j$  and variance  $\hat{\tau}_m^j$  of  $x_m^j$  in the VBI clustering module. All the mean  $\hat{x}_m^j$  and variance  $\hat{\tau}_m^j$  for  $\forall m, j$  will compose a mean-value matrix  $\hat{\mathbf{X}}$  and a variance matrix  $\hat{\mathbf{T}}$ , which will be fed into the AMP decoupling module for further refinement.

### B. Data Detection in VBI Clustering Module

For notational convenience, we first re-organize all the pseudo observations in  $\mathbf{R}$  into a column vector  $\tilde{\mathbf{r}} = [r_1^T, r_2^T, \dots, r_M^T]^T$ , where  $r_m^T$  is the  $m$ -th row vector of  $\mathbf{R}$ . Then,  $\tilde{\mathbf{r}}$  is further denoted as  $\tilde{\mathbf{r}} = [\tilde{r}_1, \tilde{r}_2, \dots, \tilde{r}_S]^T$ , where  $S = MJ$  is the total number of observations. With a little abuse of notations, the observation  $r_m^j$  in (3) is now re-written as  $\tilde{r}_s$  in  $\tilde{\mathbf{r}}$ , where  $s = j + (m-1)J$ . In the following context, this relation among the observation index  $s$ , the symbol index  $j$  and the UE index  $m$  will always hold, unless stated otherwise. Then we assume that these pseudo observations  $\tilde{r}_s$  are mutually independent with a Gaussian mixture model, i.e.

$$p(\tilde{\mathbf{r}}|\boldsymbol{\pi}, \boldsymbol{\mu}, \tau) = \prod_{s=1}^S \left( \sum_{k=1}^K \pi_{sk} \mathcal{CN}(\tilde{r}_s | \mu_m d_k, \tau^{-1}) \right) \quad (4)$$

where  $d_k \in \Delta$ ,  $\tau$  is a precision parameter, and  $\pi_{sk}$  is the mixing coefficient. Here, the mixing coefficient  $\pi_{sk}$  can be interpreted as the probability that the observation  $\tilde{r}_s$  is associated with the data symbol  $d_k \in \Delta$ . We further denote  $\boldsymbol{\mu} = \{\mu_1, \mu_2, \dots, \mu_M\}$  as the collection of channel gains, and denote  $\boldsymbol{\pi}_s = \{\pi_{s1}, \pi_{s2}, \dots, \pi_{sK}\}$  as the collection of mixing coefficients for each observation  $\tilde{r}_s$ , while  $\boldsymbol{\pi} = \{\boldsymbol{\pi}_1, \boldsymbol{\pi}_2, \dots, \boldsymbol{\pi}_S\}$  denotes the collection of  $\boldsymbol{\pi}_s$  for  $\forall s$ . For each observation  $\tilde{r}_s$ , we define a latent variable  $\mathbf{z}_s$ , which is a one-hot binary vector with length  $K$ . That is,  $\mathbf{z}_s = [z_{s1}, z_{s2}, \dots, z_{sK}]^T$ , and only one element in  $\mathbf{z}_s$  is 1. Here,  $z_{sk} = 1$  indicates the event that the observation  $\tilde{r}_s$  is *actually* associated with the symbol  $d_k \in \Delta$ . All the latent variables are collectively denoted as  $\mathbf{Z} = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_S]$ , and we have the following conditional distributions

$$p(\tilde{\mathbf{r}}|\mathbf{Z}, \boldsymbol{\mu}, \tau) = \prod_{s=1}^S \prod_{k=1}^K \mathcal{CN}(\tilde{r}_s | \mu_m d_k, \tau^{-1})^{z_{sk}}, \quad (5)$$

and

$$p(\mathbf{Z}|\boldsymbol{\pi}) = \prod_{s=1}^S \prod_{k=1}^K (\pi_{sk})^{z_{sk}}, \quad (6)$$

To facilitate the following variational Bayesian inference, we first introduce the conjugate priors for the parameters  $\boldsymbol{\mu}$ ,  $\tau$ , and  $\boldsymbol{\pi}$ . Specifically, we assume that the mixing coefficients  $\boldsymbol{\pi}_s$  are mutually independent for different observation  $\tilde{r}_s$ , i.e.

$$p(\boldsymbol{\pi}) = \prod_{s=1}^S p(\boldsymbol{\pi}_s). \quad (7)$$

Then, we choose a Dirichlet prior distribution for  $\boldsymbol{\pi}_s$ , i.e.,

$$p(\boldsymbol{\pi}_s) = \text{Dir}(\boldsymbol{\pi}_s | \boldsymbol{\alpha}_0^s) = C(\boldsymbol{\alpha}_0^s) \prod_{k=1}^K (\pi_{sk})^{\alpha_0^s[k]-1}, \quad (8)$$

where  $\boldsymbol{\alpha}_0^s = [\alpha_0^s[1], \alpha_0^s[2], \dots, \alpha_0^s[K]]^T$  is the parameter vector, and  $C(\boldsymbol{\alpha}_0^s)$  is a normalization constant for this Dirichlet distribution. By symmetry, we initialize the elements in  $\boldsymbol{\alpha}_0^s$  by the same constant  $\alpha_0$  for  $\forall s$  and  $\forall k$ . Then, the conjugate priors for  $\mu_m$  and  $\tau$  are given by

$$p(\boldsymbol{\mu}|\tau) = \prod_{m=1}^M \mathcal{CN}(\mu_m; \mu_0^m, (\lambda_0^m \tau)^{-1}), \quad (9)$$

$$p(\tau) = \text{Gam}(\tau; a_0, b_0) = \frac{1}{\Gamma(a_0)} b_0^{a_0} \tau^{a_0-1} e^{-b_0 \tau}, \quad (10)$$

where  $\mu_0^m$  and  $(\lambda_0^m \tau)^{-1}$  are the mean and variance for the complex Gaussian distribution of  $\mu_m$ , respectively. In addition,  $a_0$  and  $b_0$  are the parameters for the Gamma distribution of  $\tau$ ,  $\Gamma(\cdot)$  is the Gamma function. In this way, the joint distribution of all the variables are expressed as

$$p(\tilde{\mathbf{r}}, \mathbf{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \tau) = p(\tilde{\mathbf{r}}|\mathbf{Z}, \boldsymbol{\mu}, \tau) p(\mathbf{Z}|\boldsymbol{\pi}) p(\boldsymbol{\pi}) p(\boldsymbol{\mu}|\tau) p(\tau), \quad (11)$$

Then, we consider the variational distribution  $q$  of the latent variables and parameters with the following factorization

$$q(\mathbf{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \tau) = q(\mathbf{Z}) q(\boldsymbol{\pi}) q(\boldsymbol{\mu}, \tau). \quad (12)$$

Following the variational Bayesian inference procedure, we can infer different factor distributions in (12) as follows.

Firstly, for the latent variables  $\mathbf{Z}$ , we have

$$\begin{aligned} q(\mathbf{Z}) &= \exp\left(\mathbf{E}_{\boldsymbol{\mu}, \tau}[\ln p(\tilde{\mathbf{r}}, \mathbf{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \tau)] + \text{const}\right), \\ &= \exp\left(\mathbf{E}_{\boldsymbol{\mu}, \tau}[\ln p(\tilde{\mathbf{r}}|\mathbf{Z}, \boldsymbol{\mu}, \tau)] + \mathbf{E}_{\boldsymbol{\pi}}[\ln p(\mathbf{Z}|\boldsymbol{\pi})] + \text{const}\right) \\ &= \prod_{s=1}^S \prod_{k=1}^K e_{sk}^{z_{sk}}, \end{aligned} \quad (13)$$

where ‘‘const’’ in (13) refers to some constant numbers and they can be eliminated by the following normalization operations,  $\mathbf{E}_x[f(x)]$  represents the expectation of  $f(x)$  with respect to the random variable  $x$ , and the mixing coefficient  $e_{sk}$  is normalized as follows for each observation index  $s$ ,

$$e_{sk} = \frac{\rho_{sk}}{\sum_{k'=1}^K \rho_{sk'}}, \quad (14)$$

where

$$\begin{aligned} \rho_{sk} &= \exp\left(\mathbf{E}_{\boldsymbol{\mu}, \tau}\left[\mathcal{CN}(\tilde{r}_s|\mu_m d_k, \tau^{-1})\right] + \mathbf{E}_{\boldsymbol{\pi}}[\ln \pi_{sk}]\right), \\ &\stackrel{(a)}{=} \exp\left(\mathbf{E}_{\tau}[\ln \tau] - \ln \pi + \mathbf{E}_{\boldsymbol{\pi}}[\ln \pi_{sk}] - \mathbf{E}_{\boldsymbol{\mu}, \tau}[\tau|\tilde{r}_s - \mu_m d_k|^2]\right) \end{aligned} \quad (15)$$

where  $\pi$  in equation (a) of (15) is the circular constant, and the calculation of different terms in (15) will be explained later as in (23)-(25).

After the derivation of  $q(\mathbf{Z})$ , we consider the distribution  $q(\boldsymbol{\pi})$  in (12), and we have

$$\begin{aligned} q(\boldsymbol{\pi}) &= \exp\left(\mathbf{E}_{\mathbf{Z}, \boldsymbol{\mu}, \tau}[\ln p(\mathbf{Z}|\boldsymbol{\pi}) + \ln p(\boldsymbol{\pi})] + \text{const}\right), \\ &\stackrel{(b)}{=} \exp\left(\sum_{s=1}^S \sum_{k=1}^K (\mathbf{E}_{\mathbf{Z}}[z_{sk}] + \alpha_0^s[k] - 1) \ln \pi_{sk} + \text{const}\right), \\ &\stackrel{(c)}{=} \exp\left(\sum_{s=1}^S \sum_{k=1}^K (e_{sk} + \alpha_0^s[k] - 1) \ln \pi_{sk} + \text{const}\right), \\ &\stackrel{(d)}{=} \text{Dir}(\boldsymbol{\pi}|\bar{\boldsymbol{\alpha}}), \end{aligned} \quad (16)$$

where we have  $\mathbf{E}_{\mathbf{Z}}[z_{sk}] = e_{sk}$  in equation (b) of (16), and  $e_{sk}$  is calculated as in (15). In addition, we can conclude from equation (c) of (16) that  $q(\boldsymbol{\pi})$  still takes the form of a Dirichlet distribution. With some manipulations on the constant terms, the updated Dirichlet distribution is given in equation (d), where the updated parameter vector  $\bar{\boldsymbol{\alpha}}$  has components  $\bar{\alpha}_0^s[k]$ ,

$$\bar{\alpha}_0^s[k] = \alpha_0^s[k] + e_{sk}. \quad (17)$$

For the joint distribution  $q(\boldsymbol{\mu}, \tau)$  of  $\boldsymbol{\mu}$  and  $\tau$ , we have

$$\begin{aligned} q(\boldsymbol{\mu}, \tau) &\stackrel{(e)}{=} \exp\left(\mathbf{E}_{\mathbf{Z}, \boldsymbol{\pi}}[\ln p(\tilde{\mathbf{r}}|\mathbf{Z}, \boldsymbol{\mu}, \tau) + \ln p(\boldsymbol{\mu}, \tau)] + \text{const}\right), \\ &\stackrel{(f)}{=} \prod_{m=1}^M \mathcal{CN}\left(\mu_m; \bar{\mu}_0^m, (\bar{\lambda}_0^m \tau)^{-1}\right) \text{Gam}(\tau; \bar{a}_0, \bar{b}_0). \end{aligned} \quad (18)$$

After some mathematical manipulations on equation (e) of (18), it is observed that  $q(\boldsymbol{\mu}, \tau)$  still takes the distribution form as in (9) and (10). The updated distribution parameters in equation (f) of (18) are calculated as

$$\bar{\lambda}_0^m = \lambda_0^m + \sum_{s=(m-1)J+1}^{mJ} \sum_{k=1}^K e_{sk} |d_k|^2, \quad (19)$$

$$\bar{\mu}_0^m = \left( \lambda_0^m \mu_0^m + \sum_{s=(m-1)J+1}^{mJ} \sum_{k=1}^K e_{sk} d_k^* \tilde{r}_s \right) / \bar{\lambda}_0^m, \quad (20)$$

$$\bar{a}_0 = a_0 + S, \quad (21)$$

$$\bar{b}_0 = b_0 + \sum_{m=1}^M \lambda_0^m |\mu_0^m|^2 + \sum_{s=1}^S \sum_{k=1}^K e_{sk} |\tilde{r}_s|^2 - \sum_{m=1}^M \bar{\lambda}_0^m |\bar{\mu}_0^m|^2, \quad (22)$$

where  $d_k^*$  is the conjugate of  $d_k$ , and the cumulative summation over  $(m-1)J+1 \leq s \leq mJ$  indicates that only the observations of UE  $m$  are taken to update  $\bar{\lambda}_0^m$  and  $\bar{\mu}_0^m$ . In addition, it is noted that  $\bar{\mu}_0^m$  in (20) represents the updated channel estimate for UE  $m$ . In other words, the data-detection result  $e_{sk}$  is employed to refine the CE result  $\bar{\mu}_0^m$ .

Now we can calculate different terms in equation (a) of (15)

$$\mathbf{E}_{\boldsymbol{\pi}}[\ln \pi_{sk}] = \psi\left(\bar{\alpha}_0^s[k]\right) - \psi\left(\sum_{k'=1}^K \bar{\alpha}_0^s[k']\right), \quad (23)$$

$$\mathbf{E}_{\tau}[\ln \tau] = \psi(\bar{a}_0) - \ln \bar{b}_0, \quad (24)$$

$$\begin{aligned} &\mathbf{E}_{\boldsymbol{\mu}, \tau}[\tau|\tilde{r}_s - \mu_m d_k|^2] \\ &= \mathbf{E}_{\tau} \left[ \tau \mathbf{E}_{\boldsymbol{\mu}|\tau} \left[ |\tilde{r}_s|^2 + |\mu_m d_k|^2 - 2\text{Re}(\tilde{r}_s^* \mu_m d_k) \right] \right], \\ &= \mathbf{E}_{\tau} \left[ \tau |\tilde{r}_s|^2 + \tau |d_k|^2 \left( |\bar{\mu}_0^m|^2 + (\bar{\lambda}_0^m \tau)^{-1} \right) - 2\tau \text{Re}(\tilde{r}_s^* \bar{\mu}_0^m d_k) \right], \\ &= \frac{\bar{a}_0}{\bar{b}_0} \left[ |\tilde{r}_s|^2 + |d_k|^2 |\bar{\mu}_0^m|^2 - 2\text{Re}(\tilde{r}_s^* \bar{\mu}_0^m d_k) \right] + \frac{|d_k|^2}{\bar{\lambda}_0^m}, \end{aligned} \quad (25)$$

where (23) and (24) are obtained by the properties of the Dirichlet distribution and the Gamma distribution respectively,  $\psi(\cdot)$  is the digamma function,  $\tilde{r}_s^*$  is the conjugate of  $\tilde{r}_s$ , and  $\text{Re}(\cdot)$  takes the real part of a complex number.

With (23)-(25), we can calculate  $e_{sk}$  in (13). After that, we update the mean  $\hat{x}_s$  and variance  $\hat{\tau}_s$  of each element  $x_m^j$  in the intermediate detection target  $\mathbf{X}$ . As illustrated in Fig. 2, the updated mean  $\hat{x}_s$  and variance  $\hat{\tau}_s$  will be fed back to the next iteration of AMP decoupling for further refinement. Specifically, the mean  $\hat{x}_s$  of  $x_m^j$  is updated as

$$\hat{x}_s = \mathbf{E}_{\boldsymbol{\mu}, \tau, \mathbf{e}_s} [x_m^j] = \mathbf{E}_{\boldsymbol{\mu}, \tau, \mathbf{e}_s} [\mu_m d_{[m]}^j] = \bar{\mu}_0^m \sum_{k=1}^K e_{sk} d_k, \quad (26)$$

where  $\mathbf{e}_s = \{e_{s1}, \dots, e_{sK}\}$ , the random variables  $\boldsymbol{\mu}$  and  $\tau$  in (26) takes the updated distribution as in (18), and  $\bar{\mu}_0^m$  is the

updated channel estimate for UE  $m$ , which is derived in (20). Furthermore, the updated variance  $\hat{\tau}_s$  of  $x_m^j$  is derived as

$$\begin{aligned}
\hat{\tau}_s &= \text{Var}_{\mu, \tau, e_s} [x_m^j], \\
&\stackrel{(g)}{=} \text{Var}_{\mu, \tau} [\mu_m] \text{Var}_{e_s} [d_{[m]}^j], \\
&= \text{E}_\tau \left[ \text{Var}_{\mu|\tau} [\mu_m] \right] \left( \sum_{k=1}^K e_{sk} |d_k|^2 - \left| \sum_{k=1}^K e_{sk} d_k \right|^2 \right) \\
&= \text{E}_\tau \left[ (\bar{\lambda}_0^m \tau)^{-1} \right] \left( \sum_{k=1}^K e_{sk} |d_k|^2 - \left| \sum_{k=1}^K e_{sk} d_k \right|^2 \right) \\
&\stackrel{(h)}{=} \frac{\bar{b}_0}{\bar{\lambda}_0^m (\bar{a}_0 - 1)} \left( \sum_{k=1}^K e_{sk} |d_k|^2 - \left| \sum_{k=1}^K e_{sk} d_k \right|^2 \right)
\end{aligned} \tag{27}$$

where equation (g) of (27) is obtained by the mutual independence of  $\mu_m$  and  $d_{[m]}^j$  in  $x_m^j$ , and equation (h) is obtained by the following property for a Gamma-distributed random variable  $\tau \sim \text{Gam}(\tau; \bar{a}_0, \bar{b}_0)$ , i.e.

$$\text{E}[\tau^k] = \frac{(\bar{b}_0)^{-k} \Gamma(\bar{a}_0 + k)}{\Gamma(\bar{a}_0)}. \tag{28}$$

### C. Exploiting Joint Sparsity for Activity Detection

According to the VBI-based data-detection result,  $e_{sk}$  in (13) represents the probability that the observation  $\tilde{r}_s$  belongs to the  $k$ -th cluster. We first ignore the joint sparsity, and the transmitted symbol  $d_{[m]}^j$  should be decided as  $d_k$  if  $e_{sk}$  is the largest element among  $\{e_{s1}, e_{s2}, \dots, e_{sK}\}$ . It is noted that  $d_1 = 0$  is the equivalent null symbol from inactive UEs. We further denote  $p_s^{\text{act}}$  and  $p_s^{\text{ina}}$  as the probability that the symbol  $d_{[m]}^j$  is transmitted from an active UE or an inactive UE, respectively. We have,

$$\begin{aligned}
p_s^{\text{act}} &\propto \max\{e_{s2}, \dots, e_{sK}\}, \\
p_s^{\text{ina}} &\propto e_{s1}.
\end{aligned} \tag{29}$$

Denote  $l_m^{\text{VBI}}$  as the VBI-based log-likelihood ratio (LLR) for the activity of UE  $m$ . Considering the joint sparsity caused by UE activity,  $l_m^{\text{VBI}}$  is obtained from all the observations  $\tilde{r}_s$  of UE  $m$ , i.e.

$$l_m^{\text{VBI}} = \sum_{s=(m-1)J+1}^{mJ} \ln \frac{p_s^{\text{act}}}{p_s^{\text{ina}}} = \sum_{s=(m-1)J+1}^{mJ} \ln \frac{\max\{e_{s2}, \dots, e_{sK}\}}{e_{s1}} \tag{30}$$

If  $l_m^{\text{VBI}}$  is solely adopted for activity detection, the detection accuracy may be significantly undermined by the problem of *false alarm*, which is explained as follows.

For an inactive UE  $m$ , we have  $d_{[m]}^j = 0$ , and therefore  $x_m^j = \mu_m d_{[m]}^j = 0$ . Consequently, the pseudo observation of  $x_m^j$ , i.e.  $\tilde{r}_s$  will also be close to zero. In the VBI clustering module,  $\tilde{r}_s$  is used to jointly estimate the unknown channel gain  $\mu_m$  as in (20) and update the mean  $\hat{x}_s$  as in (26). As a result, both the channel estimate result  $\bar{\mu}_0^m$  and mean  $\hat{x}_s$  will be close to zero. In this case, the VBI module may detect this inactive UE  $m$  as an active UE which has a small channel gain  $\bar{\mu}_0^m$ . To address this problem, we consider an intuitive judgment that large estimate  $\hat{x}_s$  usually comes from active UEs, while the VBI clustering module tends to produce small

estimates  $\hat{x}_s$  for inactive UEs. Then, according to the mean  $\hat{x}_s$  and variance  $\hat{\tau}_s$  in (26) and (27), we compute an *offset LLR* [23], [24] to improve the activity detection accuracy.

Specifically, we characterize the mean  $\hat{x}_s$  as

$$\hat{x}_s = x_m^j + e_m^j, \tag{31}$$

where  $e_m^j$  denotes the estimation error between  $\hat{x}_s$  and  $x_m^j$ , with the distribution  $e_m^j \sim \mathcal{CN}(0, \hat{\tau}_s)$ . For an inactive UE  $m$ , we have  $x_m^j = 0$ , and therefore the *prior* distribution of the mean  $\hat{x}_s$  is  $\hat{x}_s \sim \mathcal{CN}(0, \hat{\tau}_s)$ . For an active UE  $m$ , the prior channel distribution  $\mu_m \sim \mathcal{CN}(0, 1)$  is assumed, and  $\mu_m$  is independent from  $d_{[m]}^j$ . Therefore, the prior mean of  $x_m^j$  is zero, while the prior variance of  $x_m^j$  is calculated as

$$\text{Var}(x_m^j) = \text{Var}(\mu_m) \text{Var}(d_{[m]}^j) = E_{\text{sym}} \triangleq \frac{1}{K-1} \sum_{k=2}^K |d_k|^2. \tag{32}$$

In this way, if UE  $m$  is active, the prior distribution for  $\hat{x}_s$  is  $\hat{x}_s \sim \mathcal{CN}(0, E_{\text{sym}} + \hat{\tau}_s)$ . Based on the above-mentioned prior distribution of  $\hat{x}_s$ , we can calculate an offset LLR  $l_s^{\text{offset}}$  for each observation index  $s$

$$\begin{aligned}
l_s^{\text{offset}} &= \ln \frac{p(\text{mean} = \hat{x}_s | \text{UE } m \text{ is active})}{p(\text{mean} = \hat{x}_s | \text{UE } m \text{ is inactive})}, \\
&= \ln \frac{\mathcal{CN}(\hat{x}_s; 0, E_{\text{sym}} + \hat{\tau}_s)}{\mathcal{CN}(\hat{x}_s; 0, \hat{\tau}_s)}, \\
&= \ln \frac{\hat{\tau}_s}{E_{\text{sym}} + \hat{\tau}_s} + \frac{|\hat{x}_s|^2}{\hat{\tau}_s} - \frac{|\hat{x}_s|^2}{E_{\text{sym}} + \hat{\tau}_s}.
\end{aligned} \tag{33}$$

Then, the decision LLR  $l_m^{\text{dec}}$  for activity detection is obtained by combining the VBI-based LLR  $l_m^{\text{VBI}}$ , the offset LLR  $l_s^{\text{offset}}$ , and the prior LLR  $l_0 = \ln \frac{p_a}{1-p_a}$  for each UE  $m$ ,

$$l_m^{\text{dec}} = l_m^{\text{VBI}} + \sum_{s=(m-1)J+1}^{mJ} l_s^{\text{offset}} + l_0. \tag{34}$$

The data detection result is obtained as

$$\begin{cases} \hat{d}_{[m]}^j = d_{k'} \text{ where } k' = \arg \max_k \{e_{s2}, \dots, e_{sK}\}, \text{ if } l_m^{\text{dec}} > 0, \\ \hat{d}_{[m]}^j = d_1 = 0, \text{ if } l_m^{\text{dec}} \leq 0. \end{cases} \tag{35}$$

After traversing all the UE indexes  $m$  and symbol indexes  $j$  for  $\hat{d}_{[m]}^j$ , we finally obtain the detection result  $\hat{D}$  of the transmitted signal matrix  $D$ .

### D. Algorithm Summary and Complexity Analysis

According to the explanations above, the AMP decoupling module works with the VBI clustering module to jointly detect UE activity and data symbols for the grant-free random access system. This entire framework is termed as the AMP-VBIC algorithm, and summarized as in Algorithm 1.

Specifically,  $\alpha_0, a_0, b_0$  and  $e_{sk}$  are initialized for the VBIC algorithm. The matrices  $S, T_s, P, T_p$  and  $T$  are initialized and updated only within the AMP decoupling module, while their  $j$ -th columns are denoted as  $s^j, \tau_s^j, p^j, \tau_p^j$  and  $\tau^j$ , respectively.  $R = [r^1, \dots, r^J]$  is updated in the AMP decoupling module, and then passed to the VBI clustering module.  $N_{it}$  is the total iteration number, and we omit the iteration index

**Algorithm 1** AMP-VBIC algorithm**Input:** received signal matrix  $\mathbf{Y}$ , spreading matrix  $\mathbf{A}$ **Output:** data detection result  $\hat{d}_{[m]}^j$  for  $\forall j, m$ **Initialize:**

$$\begin{aligned} \alpha_0 &= 0.1, \quad a_0 = 10^{-4}, \quad b_0 = 1, \quad e_{sk} = \frac{1}{K} \text{ for } \forall s, k \\ \mathbf{S} &= \mathbf{0}_{N \times J}, \quad \mathbf{T}_s = \mathbf{0}_{N \times J}, \quad \mathbf{P} = \mathbf{0}_{N \times J}, \quad \mathbf{T}_p = \mathbf{0}_{N \times J}. \\ \mathbf{R} &= \mathbf{0}_{M \times J}, \quad \mathbf{T} = \mathbf{0}_{M \times J}, \quad \hat{\mathbf{X}} = \mathbf{0}_{M \times J}, \quad \hat{\mathbf{T}} = E_{\text{sym}} \mathbf{1}_{M \times J}. \end{aligned}$$

**for**  $l = 1 : N_{it}$  **do****for**  $j = 1 : J$  **do**

$$\begin{aligned} \tau_p^j &= |\mathbf{A}|^2 \hat{\tau}^j \\ \mathbf{p}^j &= \mathbf{A} \hat{\mathbf{x}}^j - \tau_p^j \cdot \mathbf{s}^j \\ \tau_s^j &= \mathbf{1} ./ (\tau_p^j + \sigma_n^2 \mathbf{1}) \\ \mathbf{s}^j &= \tau_s^j \cdot (\mathbf{y}^j - \mathbf{p}^j) \\ \mathbf{1} ./ \tau^j &= |\mathbf{A}^H|^2 \tau_s^j \\ \mathbf{r}^j &= \hat{\mathbf{x}}^j + \tau^j \cdot (\mathbf{A}^H \mathbf{s}^j) \end{aligned}$$

1. Reshape matrix  $\mathbf{R} = [\mathbf{r}^1, \dots, \mathbf{r}^J]$  into vector  $\tilde{\mathbf{r}}$ .
2. Update  $\bar{\alpha}_0^s[k]$  as in (17) for  $\forall s, k$
3. Update  $\bar{\lambda}_0^m$  as in (19) for  $\forall m$
4. Update  $\bar{\mu}_0^m$  as in (20) for  $\forall m$
5. Update  $\bar{a}_0$  and  $\bar{b}_0$  as in (21) and (22)
6. Update  $e_{sk} = \rho_{sk} / (\sum_{k'=1}^K \rho_{sk'})$  with  $\rho_{sk}$  given in (15)
7.  $\alpha_0^s[k] = \bar{\alpha}_0^s[k]$ ,  $\lambda_0^m = \bar{\lambda}_0^m$ ,  $\mu_0^m = \bar{\mu}_0^m$ ,  $a_0 = \bar{a}_0$ ,  $b_0 = \bar{b}_0$
8. Update  $\hat{\mathbf{x}}_s$  in (26) and reorganize  $\hat{\mathbf{x}}_s$  into matrix  $\hat{\mathbf{X}}$ .
9. Update  $\hat{\tau}_s$  in (27) and reorganize  $\hat{\tau}_s$  into matrix  $\hat{\mathbf{T}}$ .

**Data Detection:** Perform final data detection as in (35)

$l$  in the notations for reading clarity. In addition,  $|\mathbf{A}|^2$  returns the square of the modulus for each element in  $\mathbf{A}$ , while  $\cdot$  and  $./$  represent the element-wise multiplication and element-wise division operations, respectively. As shown in Algorithm 1, the  $l$ -th iteration of the AMP-VBIC algorithm starts with the decoupling module, i.e. an inner loop of column-by-column operations. In this way, the AMP decoupling module accomplishes column-wise detection for all the  $J$  columns in  $\mathbf{X}$ , and produces the pseudo observation matrix  $\mathbf{R} = [\mathbf{r}^1, \dots, \mathbf{r}^J]$ .

For the VBI clustering module, we first reshape the pseudo observation matrix  $\mathbf{R}$  into vector  $\tilde{\mathbf{r}} = [r_1, \dots, r_S]^T$ . Then the VBIC algorithm is performed as in line 2 to line 6, while line 7 initializes related parameters for the next VBIC iteration. Next, the mean  $\hat{\mathbf{x}}_s$  and variance  $\hat{\tau}_s$  are updated in line 8 and line 9, and they will be reshaped into matrices  $\hat{\mathbf{X}} = [\hat{\mathbf{x}}^1, \dots, \hat{\mathbf{x}}^J]$  and  $\hat{\mathbf{T}} = [\hat{\tau}^1, \dots, \hat{\tau}^J]$ , which will be fed back to the next AMP decoupling iteration. Finally, the data detection is made according to (35).

We further analyze the computational complexity of the AMP-VBIC algorithm, which is dominated by the number of multiplication/division and exponential/logarithmic operations [22], [23]. Firstly, the AMP decoupling operations are well-known for the low complexity. Considering all the  $J$  columns, the update of  $\tau_p^j$ ,  $\mathbf{p}^j$ ,  $\tau_s^j$ , and  $\mathbf{r}^j$  will introduce  $\mathcal{O}(NMJ)$  multiplications, respectively. In addition, updating  $\tau_s^j$  and  $\mathbf{s}^j$  entails only  $\mathcal{O}(NJ)$  multiplications. For the VBIC algorithm,  $\mathcal{O}(MJK)$  multiplications are required for the update of  $\bar{\lambda}_0^m$ ,  $\bar{\mu}_0^m$ ,  $\bar{b}_0$ ,  $\hat{\mathbf{x}}_s$ , and  $\hat{\tau}_s$ , respectively. Then, the calculation of  $\rho_{sk}$  entails  $\mathcal{O}(MJK)$  multiplications and  $\mathcal{O}(MJK)$  exponential

operations. Finally, it is concluded that the AMP-VBIC algorithm totally needs  $\mathcal{O}(N_{it}NMJ + N_{it}MJK)$  multiplications and  $\mathcal{O}(N_{it}MJK)$  exponential operations. In other words, the total complexity scales only linearly with the system parameters, making the AMP-VBIC algorithm computationally favorable for practical grant-free random access systems.

**Remark 1: (Phase Ambiguity and Correction by Reference Symbol)** The Gaussian mixture model in (4) is employed in the VBIC algorithm, where the channel gain  $\mu_m$  and data symbol  $d_k$  are jointly estimated and detected. Since the modulation constellation is symmetric, we can always find non-zero phase shift  $\theta$  satisfying  $\hat{d}_k = d_k e^{j\theta} \in \Delta_a$ , e.g.  $\theta = \pi$ . In this case, the VBIC algorithm may detect the data symbol as  $\hat{d}_k = d_k e^{j\theta}$  and estimate the channel gain as  $\hat{\mu}_m = \mu_m e^{-j\theta}$  by mistake, since the wrong combination  $(\hat{d}_k, \hat{\mu}_m)$  and the correct one  $(d_k, \mu_m)$  will produce the same probability  $p(\tilde{\mathbf{r}}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\tau})$  in (4). This problem is dubbed as the phase ambiguity problem, and it can be readily addressed by the RS [16]. Specifically, we can take an arbitrary symbol from the modulation symbol alphabet  $\Delta_a$  as the RS, and denote  $d_{[m]}^{\text{RS}}$  as the RS of the  $m$ -th UE. Denote  $\mathbf{d}_{[m]}^{\mathcal{J}} = [d_{[m]}^1, \dots, d_{[m]}^J]$  as the length- $J$  data sequence of the  $m$ -th UE, and  $\mathbf{d}_{[m]}^{\text{RS}}$  will be transmitted along with  $\mathbf{d}_{[m]}^{\mathcal{J}}$  if the  $m$ -th UE is activated. If the  $m$ -th UE is further detected as active at the AP, the AMP-VBIC algorithm will produce the detection results  $\hat{d}_{[m]}^{\text{RS}}$  and  $\hat{\mathbf{d}}_{[m]}^{\mathcal{J}}$  for the RS and the data sequence, respectively. Since  $d_{[m]}^{\text{RS}}$  is predetermined and known to the AP, the phase ambiguity can be corrected as

$$\bar{\mathbf{d}}_{[m]}^{\mathcal{J}} = \hat{\mathbf{d}}_{[m]}^{\mathcal{J}} \frac{d_{[m]}^{\text{RS}}}{\hat{d}_{[m]}^{\text{RS}}}, \quad (36)$$

where  $\bar{\mathbf{d}}_{[m]}^{\mathcal{J}}$  is the corrected data-detection result. For notation clarity, we assume that the final detection results obtained in (35) have already been corrected by the RS.

## IV. SIMULATIONS

In this section, we evaluate the performances of our proposed AMP-VBIC algorithm for joint UE activity and data detection. To begin with, we define the *CSIR* as the knowledge of the channel matrix  $\mathbf{U}$  at the AP, the *support* of  $\mathbf{X}$  as the *exact identity* of active UEs, and the *sparsity level* as the *exact number* of active UEs. Due to the massiveness and sporadic activity of UEs, these three types of information defined above are unavailable to the AP. However, the spreading matrix  $\mathbf{A}$  is assumed predetermined, and thus known to the AP. In addition, pseudo-random Gaussian sequences are adopted as spreading sequences for each UE, i.e. the elements in  $\mathbf{A}$  are independently and identically distributed with distribution  $\mathcal{CN}(0, 1)$ . Furthermore, we adopt the 16-Quadrature Amplitude Modulation (16-QAM) for transmitted data symbols, and we consider the detection performance for *uncoded* data sequences<sup>1</sup>. Specifically, three performance

<sup>1</sup>The proposed VBI-based data detection can also work with coded data sequences. If UE  $m$  is detected as active, the probability that  $d_{[m]}^j = d_k \in \Delta_a$  is proportional to  $e_{sk}$  in line 6 of Algorithm 1. In this way, we can compute the LLR for each transmitted bit according to the 16-QAM constellation, and the LLR is output from the VBI module to the soft-decision decoder.

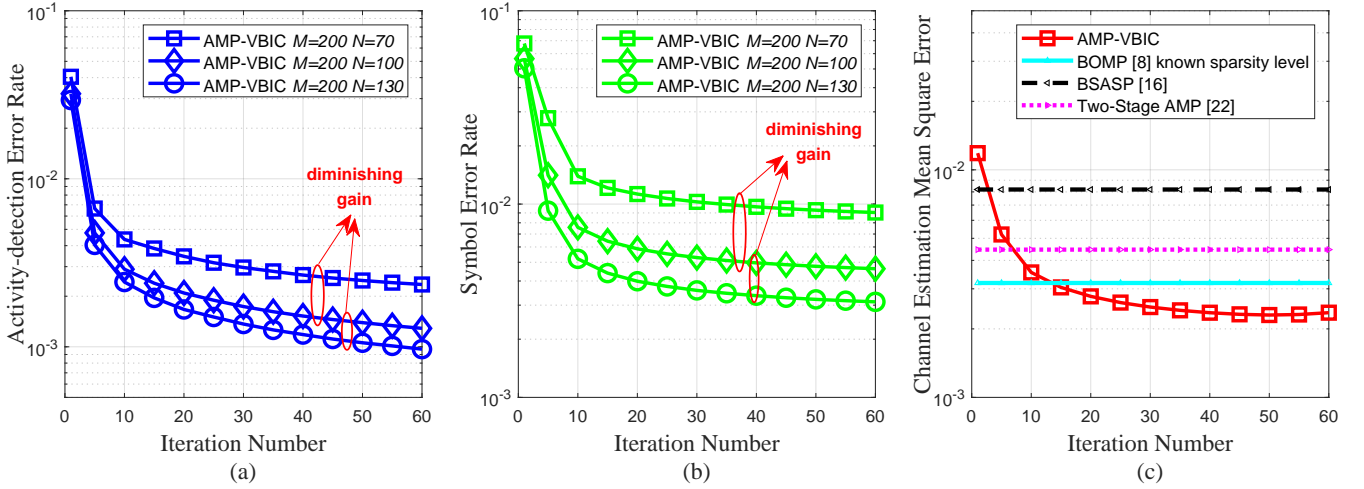


Fig. 3. Joint detection performances under different iteration number  $N_{it}$  with configurations  $M = 200$ ,  $p_a = 0.1$ ,  $J = 10$ , and  $\text{SNR} = 5\text{dB}$ .

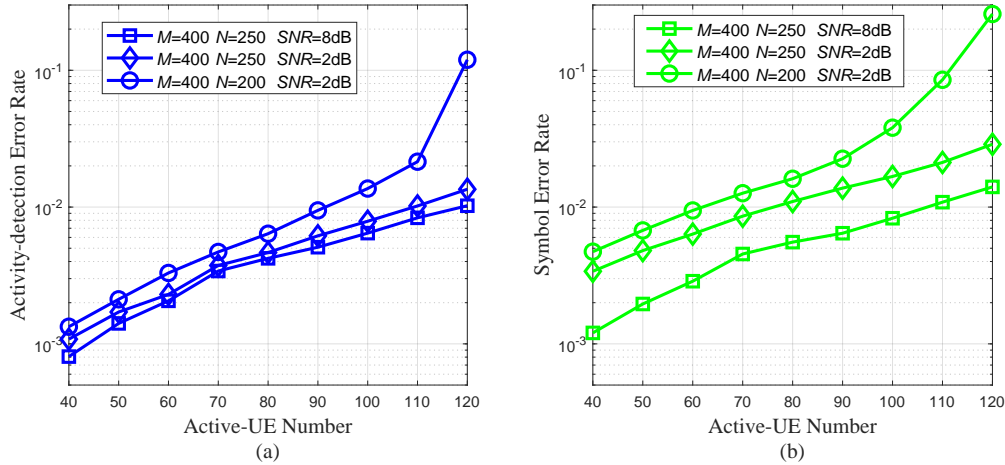


Fig. 4. Joint detection performances under different number of active UEs with configurations  $M = 400$  and  $J = 10$ .

metrics are considered in the following simulations, i.e., the activity-detection error rate (AER), the symbol error rate (SER), and the channel estimation mean square error (CE-MSE), which are defined as follows

$$\begin{aligned} \text{AER} &\triangleq \sum_{m=1}^M |\delta_m - \hat{\delta}_m|/M, \\ \text{SER} &\triangleq \frac{1}{MJ} \|\mathbf{D}_{M \times J} - \hat{\mathbf{D}}_{M \times J}\|_0, \\ \text{CE-MSE} &\triangleq \sum_{m=1}^M |\mu_m - \bar{\mu}_0^m|^2/M, \end{aligned}$$

where the activity indicator  $\delta_m = 1$  if the  $m$ -th UE is activated. Otherwise,  $\delta_m = 0$ .  $\hat{\delta}_m$  is the detection result of  $\delta_m$ . The  $l_0$  norm of a matrix, i.e.  $\|\cdot\|_0$  returns the number of non-zero elements, and we set  $\mu_m = 0$  for inactive UEs. In addition, we define the signal-to-noise ratio (SNR) as  $\text{SNR} \triangleq 10 \ln \frac{E_{\text{Sym}}}{\sigma_n^2}$ .

#### A. Convergence Performance

Firstly, we investigate the convergence performance of our proposed AMP-VBIC algorithm under different number of

iterations  $N_{it}$ , and the simulation results are illustrated in Fig. 3. It is shown in Fig. 3(a) and Fig. 3(b) that both the AER and the SER performances get rapidly improved in the first 20 iterations, then the detection performances tend to converge afterwards. Furthermore, increasing the spreading length  $N$  could effectively lower the AER and SER. However, we also observe a *diminishing gain*, i.e., increasing  $N$  from 70 to 100 contributes to a more prominent performance gain than further increasing  $N$  from 100 to 130. In addition, the spreading length  $N$  is fixed as 70 in Fig. 3(c), and it is shown that we can gradually improve the CE accuracy with iterations. As shown in (20), the CE update of  $\bar{\mu}_0^m$  exploits the clustering results of *all the data symbols*. Therefore, the AMP-VBIC algorithm outperforms the other existing solutions [8], [16], [22], which only employ one reference symbol for CE.

#### B. Performance with Different Active-UE Number

We further investigate the SER and AER performances of the AMP-VBIC algorithm with different number of active UEs, and the simulation results are illustrated in Fig. 4. A general observation is that more active UEs will lead



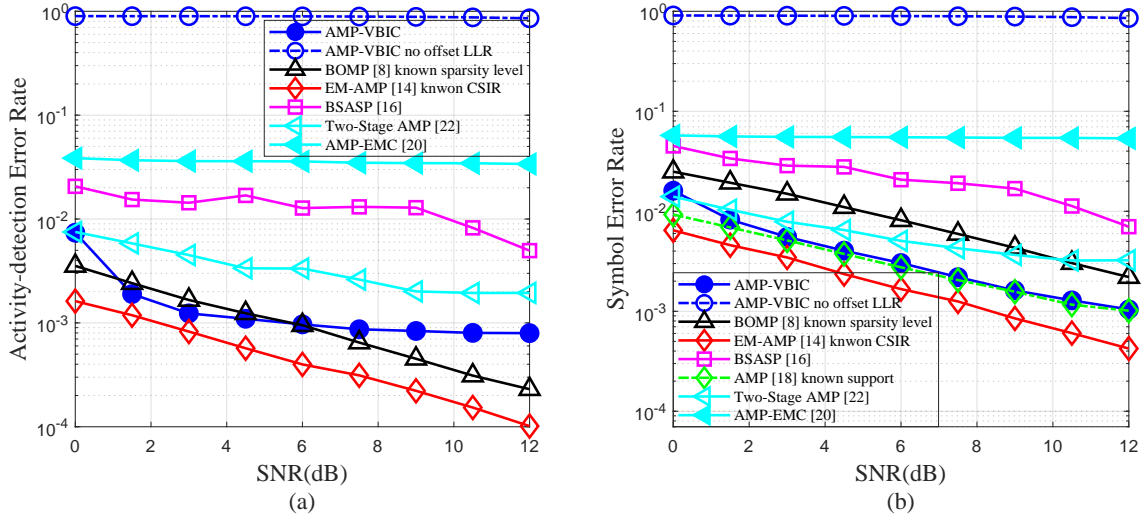


Fig. 5. Comparison on joint detection performances under different SNR with configurations  $M = 200$ ,  $N = 120$ ,  $p_a = 0.1$ , and  $J = 20$ .

to deteriorated SER and AER performances. Then, we fix  $N = 250$  and reduce SNR from 8dB to 2dB. In this case, the performance loss on AER is almost negligible, while that on the SER performance is more obvious. After that, we fix SNR as 2dB, and reduce  $N$  from 250 to 200. It is shown that, with only a small number of active UEs, the performance loss caused by reducing  $N$  is not obvious on both AER and SER. However, given  $N = 200$  and SNR = 2 dB, if we further increase the active-UE number from 90 to 120, both the SER and the AER performances will deteriorate drastically, which indicates the detection failure. In other words, to avoid detection failure caused by RA congestion, we need to increase the SNR or the spreading length  $N$ .

### C. Performance Comparison with Different SNR

Finally, we simulate the AER and SER performances under different SNR, and compare the proposed AMP-VBIC algorithm with different state-of-art solutions. These solutions include the BOMP algorithm [8] with *known sparsity level*, EM-AMP algorithm [14] with *known CSIR*, BSASP algorithm [16], AMP algorithm [18] with *known support*, and the two-stage AMP algorithm where the UAD problem is firstly addressed by the BGMP algorithm [22] and the MUD problem is secondly addressed by the AMP algorithm. In addition, we try to replace our proposed VBI clustering module with the EM-based clustering (EMC) approach in [20], and establish an AMP-EMC algorithm for comparison. Recall that the sparsity level, the CSIR, and the support information are actually unavailable to the AP. Therefore, the solutions aided by such ideal knowledge can provide some performance lower bounds. The simulation results are illustrated in Fig. 5

It is shown in Fig. 5(a) that the proposed AMP-VBIC algorithm exhibits superior AER performance to most solutions, except for those aided by known CSIR or sparsity level. In addition, the AER performance of the AMP-VBIC algorithm will not be further improved with higher SNR, which is consistent with the results in Fig. 4(a). This observation can be explained by the fact that the channel noise has much smaller

impacts on AER than the MUI in the high-SNR regime, while increasing  $N$  is an effective method to mitigate the MUI. We can observe from Fig. 5(b) that the AMP-VBIC algorithm still outperforms most state-of-art solutions, and its SER performance could closely approach the performance lower bounds within a wide range of SNR. In contrast to the AER performance in Fig. 5(a), the SER of the AMP-VBIC algorithm could be effectively improved with higher SNR, since weaker channel noise is beneficial to the data-detection accuracy for active UEs. In addition, the BOMP algorithm [8] and the BSASP [16] algorithm are shown to exhibit inferior SER performances, since they adopt the least square principle for data detection, which neglects the inherent modulation constraints of data symbols and thus undermines the data detection accuracy.

We also demonstrate the effectiveness of including the offset LLR for activity detection in (34). It is shown in Fig. 5 that the AER and SER performances will approach  $1 - p_a$  if the offset LLR is not included for the AMP-VBIC algorithm. This result supports our claim that the problem of false alarm will significantly undermine the activity-detection accuracy if  $l_{m}^{VBI}$  is solely adopted for activity detection. In addition, the simulation results demonstrate that EMC approach [20] fails to work for the clustering problem under the AMP framework. One possible reason is that the centroid of each cluster, i.e. the term  $\mu_m d_k$  in (4), is estimated *independently* in the EMC approach. In other words, the EMC approach neglects the inherent constraint that different cluster centroids of UE  $m$  share the same channel-gain term  $\mu_m$ . Consequently, the clustering accuracy of the EMC approach is significantly undermined when the pseudo-observations  $\mathbf{R}$  are contaminated by MUI, or when the symbol alphabet  $\Delta$  is composed of high-order modulation symbols.

## V. CONCLUSIONS

In order to address the joint UAD and MUD problem in grant-free random access, we formulated this joint detection problem as a clustering problem under the Gaussian mixture

model. In conjunction with the AMP decoupling module, we developed a VBIC algorithm to solve this clustering problem. Compared with the state-of-art algorithms, our proposed AMP-VBIC algorithm demonstrated a significant performance gain.

## REFERENCES

- [1] W. Zhan, C. Xu, X. Sun, and J. Zou, "Toward optimal connection management for massive machine-type communications in 5G system," *IEEE Internet of Things Journal*, vol. 8, no. 17, pp. 13237-13250, Sept. 2021.
- [2] Y. Ma, Z. Yuan, W. Li, and Z. Li, "Novel solutions to NOMA-based modern random access for 6G-enabled IoT," *IEEE Internet of Things Journal*, vol. 8, no. 20, pp. 15382-15395, Oct. 2021.
- [3] J. Gao, W. Zhuang, M. Li, X. Shen and X. Li, "MAC for machine-type communications in industrial IoT—Part I: protocol design and analysis," *IEEE Internet of Things Journal*, vol. 8, no. 12, pp. 9945-9957, June 2021.
- [4] J. Hong, W. Choi, and B. D. Rao, "Sparsity controlled random multiple access with compressed sensing," *IEEE Transactions on Wireless Communications*, vol. 14, no. 2, pp. 998-1010, Feb. 2015.
- [5] D. L. Donoho, "Compressed sensing," *IEEE Trans. Inf. Theory.*, vol. 52, no. 4, pp. 1289-1306, Apr. 2006.
- [6] B. Shim and B. Song, "Multiuser detection via compressive sensing," *IEEE Commun. Lett.*, vol. 16, no. 7, pp. 972-974, Jul. 2012.
- [7] X. Xu, X. Rao, and V. K. N. Lau, "Active user detection and channel estimation in uplink CRAN systems," in *Proc. IEEE Int. Conf. Commun.*, Jun. 2015, pp. 2727-2732.
- [8] Y. Zhang, Q. Guo, Z. Wang, J. Xi, and N. Wu, "Block sparse Bayesian learning based joint user activity detection and channel estimation for grant-free NOMA systems," *IEEE Trans. Veh. Technol.*, vol. 67, no. 10, pp. 9631-9640, Oct. 2018.
- [9] 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, Jun. 2018.
- [10] L. Liu and W. Yu, "Massive connectivity with massive MIMO-Part II: Achievable rate characterization," *IEEE Trans. Signal Process.*, vol. 66, no. 11, pp. 2947-2959, Jun. 2018.
- [11] Z. Chen and W. Yu, "Massive device activity detection by approximate message passing," *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 3514-3518.
- [12] Z. Zhang, Y. Li, C. Huang, Q. Guo, C. Yuen, and Y. L. Guan, "DNN-aided block sparse Bayesian learning for user activity detection and channel estimation in grant-free non-orthogonal random access," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 12, pp. 12000-12012, Dec. 2019.
- [13] A. T. Abebe and C. G. Kang, "Iterative order recursive least square estimation for exploiting frame-wise sparsity in compressive sensing-based MTC," *IEEE Commun. Lett.*, vol. 20, no. 5, pp. 1018-1021, May 2016.
- [14] C. Wei, H. Liu, Z. Zhang, J. Dang, and L. Wu, "Approximate message passing-based joint user activity and data detection for NOMA," *IEEE Commun. Lett.*, vol. 21, no. 3, pp. 640-643, Mar. 2017.
- [15] B. Wang, L. Dai, T. Mir and Z. Wang, "Joint user activity and data detection based on structured compressive sensing for NOMA," *IEEE Commun. Lett.*, vol. 20, no. 7, pp. 1473-1476, Jul. 2016.
- [16] Y. Du *et al.*, "Joint channel estimation and multiuser detection for uplink grant-free NOMA," *IEEE Wireless Commun. Lett.*, vol. 7, no. 4, pp. 1473-1476, Aug. 2018.
- [17] S. Jiang, X. Yuan, X. Wang, C. Xu, and W. Yu, "Joint user identification, channel estimation, and signal detection for grant-free NOMA" *IEEE Transactions on Wireless Communications*, vol. 19, no. 10, pp. 6960-6976, 2020.
- [18] D. L. Donoho, A. Maleki, and A. Montanari, "Message-passing algorithms for compressed sensing," *Proc. Nat. Acad. Sci. USA*, vol. 106, no. 45, pp. 18914-18919, Nov. 2009.
- [19] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *J. Roy. Statist. Soc. B. Methodol.*, vol. 39, no. 1, pp. 1-38, 1977.
- [20] A. Salari, M. Shirvanimoghaddam, M. B. Shahab, R. Arablouei, and S. Johnson, "Clustering-based joint channel estimation and signal detection for grant-free NOMA" in *Proc. 2020 IEEE Globecom Workshops*, 2020, pp. 1-6.
- [21] Z. Zhang, X. Cai, C. Li, C. Zhong, and H. Dai, "One-bit quantized massive MIMO detection based on variational approximate message passing," *IEEE Transactions on Signal Processing*, vol. 66, no. 9, pp. 2358-2373, 2018.
- [22] L. Liu, C. Huang, Y. Chi, C. Yuen, Y. L. Guan, and Y. Li, "Sparse vector recovery: Bernoulli-Gaussian message passing" in *Proc. 2020 IEEE Globecom*, 2017, pp. 1-6.
- [23] Z. Zhang *et al.*, "User activity detection and channel estimation for grant-free random access in LEO satellite-enabled Internet of Things," *IEEE Internet of Things Journal*, vol. 7, no. 9, pp. 8811-8825, Sept. 2020.
- [24] J. Ahn, B. Shim, and K. B. Lee, "EP-based joint active user detection and channel estimation for massive machine-type communications," *IEEE Transactions on Communications*, vol. 67, no. 7, pp. 5178-5189, July 2019.