

Improved Bayesian Learning Detectors for Uplink Grant-Free MIMO-NOMA

Boran Yang, Xiaoxu Zhang, *Member, IEEE*,
Li Hao, *Member, IEEE*, and George K. Karagiannidis, *Fellow, IEEE*

Abstract—Grant-free non-orthogonal multiple access (GF-NOMA) and multiple-input multiple-output (MIMO) techniques are key enablers for massive machine-type communications (mMTC) in 5G cellular Internet of Things. On the other hand, compressed sensing (CS) is widely accepted for multiuser detection, due to the sporadic traffic of mMTC. In this letter, we propose two Bayesian-based CS algorithms for uplink grant-free MIMO-NOMA. Exploiting the iterative expectation maximization (EM) and maximum ratio combining techniques, the spatially enhanced sparse Bayesian learning (SE-SBL) algorithm is developed to alternately update hyperparameter values and combined observation signals. Furthermore, we embed the generalized approximate message passing technique into the SE-SBL and propose a low-complexity computational approach. In particular, the aforementioned Bayesian algorithms do not require any prior knowledge of user activity level and noise power. Simulation results show that the proposed Bayesian methods exhibit a superior performance gain over the state-of-the-art.

Index Terms—Massive machine-type communications, grant-free, non-orthogonal multiple access, multiple-input multiple-output, Bayesian compressed sensing

I. INTRODUCTION

Massive machine-type communications (mMTC) focuses on uplink access of large-scale users with bursty short-time transmissions, where conventional wireless communication networks cannot adequately meet the performance requirements [1]. Non-orthogonal multiple access (NOMA) [2], which provides massive connectivity through non-orthogonal resource allocation among users with limited resources, has been extensively studied. In addition, a grant-free NOMA (GF-NOMA) scheme [3], in which users can arbitrarily transmit data to the base station (BS) without a complicated interactive handshaking process, is widely accepted to mitigate transmission delay and signaling overhead. mMTC has an inherent sporadic traffic characteristic, i.e., only a small percentage of

potential users transmit data to the BS [4]. However, the BS cannot identify which users are active at any given time, which motivates additional challenges for multiuser detection (MUD) [5]. Fortunately, sparse activity detection can be formulated as a compressed sensing (CS) problem and efficiently solved by using sophisticated sparse reconstruction algorithms [1]–[5].

Ridge detector (RD) and Lasso detector (LD) proposed in [6] for MUD, by exploiting the sparsity of the transmitted signal. In terms of detection accuracy, RD and LD significantly surpass linear least squares (LS) and minimum mean square error (MMSE) methods [3]. However, both works have not carried a further step toward mining the underlying spatial structure of user activity. Inspired by the above observations, the authors in [7] explored the spatial correlation of user activity in the multiple input multiple output (MIMO) case and proposed an effective MUD approach for multiantenna reception, based on approximate message passing (AMP). The authors in [8] and [9] proposed the variant forms of bilinear generalized AMP algorithm for joint channel estimation and MUD in MIMO-enabled GF-NOMA systems. Block greedy-type algorithms, such as block orthogonal matching pursuit (B-OMP) [10], block compressive sampling matching pursuit (B-CoSaMP) [11], and block subspace pursuit (B-SP) [12], have been proposed to perform MUD by exploiting the block sparsity of a set of sparse vectors. However, the proposed algorithms require accurate knowledge of user activity or noise power, which is generally not practical in mMTC scenarios.

In this letter, we formulate the MUD problem of uplink grant-free MIMO-NOMA system as a block compressed sensing one, by integrating the spatial structure of user activity. The main contributions of this letter are summarized as follows: 1) Bayesian CS algorithms, i.e., spatially enhanced sparse Bayesian learning (SE-SBL) and spatially enhanced generalized approximate message passing sparse Bayesian learning (SE-GAMP-SBL) algorithms, are proposed to enable effective MUD for multiantenna reception; 2) the SE-SBL algorithm alternately obtains estimated values for hyperparameters and inferred value for combined observation signal via iterative expectation maximization (EM) [13] and combining [14] techniques; 3) the SE-GAMP-SBL algorithm significantly reduces the computational complexity of SE-SBL by approximating the posterior distribution without performance penalty; 4) the proposed two algorithms do not require user sparsity and noise power as auxiliary prior information.

II. SYSTEM MODEL

We consider a typical uplink grant-free MIMO-NOMA system, where K single-antenna users transmit data to the

Manuscript received July 25, 2023; revised August 21, 2023; accepted September 14, 2023. Date of publication September 14, 2023; date of current version September 14, 2023. This work was supported in part by National Key R&D Program of China under Grant 2018YFB1801104, in part by NSFC Project under Grant 62001399, and in part by Sichuan Science and Technology Program under Grant 2022NSFSC0910. The associate editor coordinating the review of this letter and approving it for publication was Lukas T. N. Landau. (Corresponding author: Xiaoxu Zhang.)

B. Yang, X. Zhang, and L. Hao are with the School of Information Science and Technology, Southwest Jiaotong University, Chengdu 611756, China (e-mail: bryang@my.swjtu.edu.cn; xiaoxuzhang@swjtu.edu.cn; l-hao@swjtu.edu.cn).

G. K. Karagiannidis is with the Department of Electrical and Computer Engineering, Aristotle University of Thessaloniki, 541-24 Thessaloniki, Greece, and also with the Cyber Security Systems and Applied AI Research Center, Lebanese American University (LAU), Beirut 1102-2801, Lebanon (e-mail: geokarag@auth.gr).

BS equipped with M antennas. For a certain time slot, K_a users are activated to transmit data while others remain silent, $K_a \ll K$. The transmitted symbol of the k -th user is denoted as z_k . If the k -th user is active, $z_k \in \mathcal{A}$, otherwise, $z_k = 0$, where \mathcal{A} denotes a constellation set. The spreading sequence for the k -th user is $\Psi_k = (\Psi_{1k}, \Psi_{2k}, \dots, \Psi_{Nk})^T$ with N subcarriers. To gain insight into the fundamental limit of MUD performance, we assume synchronous reception and perfect channel estimation [3], [4], [6]. In this case, the received signal from the m -th antenna of the BS can be expressed as

$$\mathbf{y}^m = \sum_{k=1}^K \text{diag}(\mathbf{h}_k^m) \Psi_k z_k + \mathbf{w}^m = \mathbf{H}^m \mathbf{z} + \mathbf{w}^m, \quad (1)$$

where $\mathbf{h}_k^m = (h_{1k}^m, h_{2k}^m, \dots, h_{Nk}^m)^T$ is the flat Rayleigh fading channel coefficient between user k and BS antenna m , $\mathbf{z} = (z_1, z_2, \dots, z_K)^T$ is the transmitted signal vector for all K users, $\mathbf{H}^m \in \mathbb{C}^{N \times K}$ is the equivalent channel matrix, $\mathbf{w}^m \in \mathbb{C}^{N \times 1}$ is the noise vector with complex Gaussian distribution $\mathcal{CN}(0, \sigma^2 \mathbf{I})$, $\text{diag}(\cdot)$ denotes regular diagonalization.

Although the channel coefficients between each user and BS exhibit non-correlation, the user activities are associated for different receiving antennas, as illustrated in Fig. 1. To exploit this special structure, we extend (1) to the block CS model for multiantenna reception,

$$\mathbf{y} = \mathbf{H} \mathbf{x} + \mathbf{w}, \quad (2)$$

where $\mathbf{H} = \text{blkdiag}(\mathbf{H}^1, \mathbf{H}^2, \dots, \mathbf{H}^M) \in \mathbb{C}^{MN \times MK}$ is the block diagonal matrix, $\mathbf{x} = \text{vec}(\mathbf{z}, \mathbf{z}, \dots, \mathbf{z}) \in \mathbb{C}^{MK \times 1}$ is the block sparse vector, $\mathbf{y} = \text{vec}(\mathbf{y}^1, \mathbf{y}^2, \dots, \mathbf{y}^M) \in \mathbb{C}^{MN \times 1}$ is the linear observation vector, $\mathbf{w} = \text{vec}(\mathbf{w}^1, \mathbf{w}^2, \dots, \mathbf{w}^M) \in \mathbb{C}^{MN \times 1}$ is the linear noise vector, $\text{vec}(\cdot)$ denotes column-wise vectorization, $\text{blkdiag}(\cdot)$ denotes block diagonalization.

A key observation of this paper is to detect the active users and the symbols transmitted, which is equivalent to estimating \mathbf{z} from the block sparse signal \mathbf{x} recovered from the block CS model (2). Note that \mathbf{x} can be decoupled into M blocks with a common sparse pattern. Fortunately, this unique property can be extracted to enhance MUD performance.

III. THE SE-SBL ALGORITHM

The classical SBL algorithm [3] does not leverage the high correlation of user activity in the spatial domain induced by multiantenna reception. To this end, we introduce the diversity combining technique into the SBL framework to develop a spatially enhanced SBL algorithm named SE-SBL.

Similar to [3], ignoring the finite-alphabet constraints, SE-SBL employs a standard two-layer hierarchical prior model dominated by a set of empirical hyperparameters. Since most of the elements in \mathbf{x} are zero for small user activity factor p_a , in the first layer, a zero-mean Gaussian prior distribution is assigned to each element of \mathbf{x} ,

$$p(\mathbf{x}; \boldsymbol{\alpha}) = \prod_{m=1}^M \prod_{k=1}^K \mathcal{CN}(x_{mk} | 0, \alpha_{mk}^{-1}) \quad (3)$$

$$= (2\pi)^{-MK} |\boldsymbol{\Lambda}| \exp(-\mathbf{x}^T \boldsymbol{\Lambda} \mathbf{x}),$$

where x_{mk} is the k -th element of the m -th block in \mathbf{x} ,

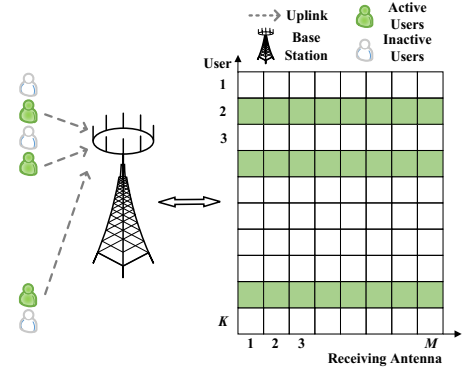


Fig. 1. A typical uplink grant-free MIMO-NOMA system and the spatial structure of user activity.

$\boldsymbol{\alpha} = (\alpha_{11}, \dots, \alpha_{mk}, \dots, \alpha_{MK})^T$ is a non-negative hyperparameter designed to regulate the block sparse signal \mathbf{x} , $\boldsymbol{\Lambda} = \text{diag}(\boldsymbol{\alpha})$. Specifically, the element x_{mk} tends to zero when α_{mk} approaches positive infinity. In the second layer, the hyperparameter $\boldsymbol{\alpha}$ follows Gamma prior distribution dominated by the parameters $\{a, b\}$,

$$p(\boldsymbol{\alpha}) = \prod_{m=1}^M \prod_{k=1}^K \Gamma(a)^{-1} b^a \alpha_{mk}^{a-1} e^{-b\alpha_{mk}}. \quad (4)$$

For simplicity, we consider β as the inverse of the noise variance, i.e., $\beta \triangleq 1/\sigma^2$. The Gaussian noise generally satisfies $p(\mathbf{w}) = \mathcal{CN}(\mathbf{w} | 0, \beta^{-1} \mathbf{I})$. Similarly, the Gamma prior distribution dominated by the parameters $\{c, d\}$ is assigned to β , i.e., $p(\beta) = \Gamma(c)^{-1} d^c \beta^{c-1} e^{-d\beta}$. The likelihood function of the received signal \mathbf{y} can be expressed as

$$p(\mathbf{y} | \mathbf{x}; \beta) = \mathcal{CN}(\mathbf{y} | \mathbf{H} \mathbf{x}, \beta^{-1} \mathbf{I}). \quad (5)$$

Based on Bayesian principle, the posterior distribution of the block sparse signal \mathbf{x} is calculated as follows

$$p(\mathbf{x} | \mathbf{y}; \boldsymbol{\alpha}, \beta) = \frac{p(\mathbf{y} | \mathbf{x}; \beta) p(\mathbf{x}; \boldsymbol{\alpha})}{\int p(\mathbf{y} | \mathbf{x}; \beta) p(\mathbf{x}; \boldsymbol{\alpha}) d\mathbf{x}}. \quad (6)$$

Inserting (3) and (5) into (6), and performing several approximation operations [3], the posterior distribution of \mathbf{x} can be simplified as

$$p(\mathbf{x} | \mathbf{y}; \boldsymbol{\alpha}, \beta) = \mathcal{CN}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad (7)$$

where $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ are the posterior mean and variance as follows

$$\mathbf{x}^{\text{SBL}} = \boldsymbol{\mu} = \beta \boldsymbol{\Sigma} \mathbf{H}^T \mathbf{y}, \quad (8)$$

$$\boldsymbol{\Sigma} = (\beta \mathbf{H}^T \mathbf{H} + \boldsymbol{\Lambda})^{-1}. \quad (9)$$

The posterior mean is treated as the estimated block sparse signal. After obtaining \mathbf{x}^{SBL} , we can acquire the observation signals recovered from each receiving antenna according to

$$\hat{\mathbf{X}} = (\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_M) = \text{vec}^{-1}(\mathbf{x}^{\text{SBL}}, M). \quad (10)$$

Exploiting the spatial diversity of interference pattern suffered by the transmitted signal and the spatial correlation of user activity induced by multiantenna reception, the combining technique [14] is applied to $\hat{\mathbf{X}}$, for creating a combined

observation, on which updating the hyperparameters is attempted. Benefiting from the spatial structure of user activity, the combining process of multiple observed versions can be interpreted as an improvement of the posterior distribution of the transmitted signal \mathbf{z} . Consider M observations of the transmitted signal, $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_M$, the combined observation signal can be expressed as

$$\hat{\mathbf{z}} = \sum_{m=1}^M \eta_m \mathbf{X}_m, \quad (11)$$

where η_m is the combination weight for the m -th branch.

Any of the combining techniques can be applied to $\hat{\mathbf{X}}$, we focus on the maximum ratio combining (MRC) scheme, since it is optimal when the interference on each observed signal is independent [14]. Assuming that the local recovery noises have the same power, the weighting $\boldsymbol{\eta} = (\eta_1, \eta_2, \dots, \eta_M)^T$ is proportional to the amplitude ratio. In the MRC scheme, each branch is weighted proportional to its root mean square of the observed signal level. Furthermore, the combined signal is expanded into a block-wise structure serving as the improved posterior mean of \mathbf{x} ,

$$\mathbf{x}^{\text{SE-SBL}} = \hat{\boldsymbol{\mu}} = \text{vec}(\hat{\mathbf{z}}, \hat{\mathbf{z}}, \dots, \hat{\mathbf{z}}). \quad (12)$$

The iterative EM algorithm is implemented to update the hyperparameters [13]. In the E-step, we first fixed the current values of hyperparameters $\{\alpha^{(t-1)}, \beta^{(t-1)}\}$, followed by calculating the expectation value of the maximum complete likelihood function incorporating the extended value of the combined signal $\mathbf{x}^{\text{SE-SBL}}$, also known as the Q -function. In the M-step, we derive the partial derivatives of Q -function with respect to α and β , and take its values zero to obtain the updated hyperparameters as follows

$$\alpha_{mk}^{(t)} = \frac{2a - 1}{2b + \hat{\boldsymbol{\mu}}_{mk}^2 + \Sigma_{mk, mk}}, \quad (13)$$

$$\beta^{(t)} = \frac{NM + 2c - 2}{2d + \|\mathbf{y} - \mathbf{H}\hat{\boldsymbol{\mu}}\|_2^2 + \text{tr}[\boldsymbol{\Sigma}\mathbf{H}^T\mathbf{H}]}. \quad (14)$$

The proposed SE-SBL approach dedicated to MUD with multiantenna reception is summarized in Algorithm 1.

Algorithm 1 SE-SBL

Initializing: $\forall m, k: \alpha_{mk}^{(0)} = 1, \beta^{(0)} = 1, \hat{\boldsymbol{\mu}}^{(0)} = \mathbf{0}$, and $t = 1$.
 Step 1. Calculating the posterior distribution of block sparse signal \mathbf{x} according to (8) and (9);
 Step 2. Creating the combined observation signal $\hat{\mathbf{z}}$ according to (11), and extending it into a block-wise pattern according to (12);
 Step 3. Updating the hyperparameters according to (13) and (14).
 Repeating the above steps until $\|\hat{\boldsymbol{\mu}}^{(t)} - \hat{\boldsymbol{\mu}}^{(t-1)}\|_2 \leq \varepsilon$, where ε is a prescribed threshold value of 10^{-6} .

IV. THE SE-GAMP-SBL ALGORITHM

SE-SBL is an effective MUD algorithm for multiantenna reception and possesses excellent recovery performance without prior knowledge of user activity factor. However, the inverse operation of the $KM \times KM$ covariance matrix is computed in the process of iteratively updating the block sparse signal \mathbf{x} . The SE-SBL algorithm has a computational complexity of

$\mathcal{O}(K^3M^3)$, which raises new challenges for MUD in the case of massive connectivity.

Motivated by GAMP-SBL [15], we embed the GAMP technique into the SE-SBL framework to replace matrix inversion and develop a computationally efficient MUD method entitled SE-GAMP-SBL. GAMP is an approximate implementation of loop belief propagation (BP) via central limit theorem and Taylor series expansion.

SE-GAMP-SBL also features a two-layer hierarchical prior model. Let $\Theta = \{\alpha, \beta\}$ denote the known hyperparameters and it will be updated in M-step. GAMP provides the approximated posterior distributions of \mathbf{x} and $\mathbf{u} = \mathbf{H}\mathbf{x}$. First, GAMP approximates the posterior distribution of x_{mk} as

$$\begin{aligned} p(x_{mk} | \mathbf{y}, \hat{r}_{mk}, \tau_{mk}^r, \Theta) &= \mathcal{CN}(x_{mk} | \mu_{mk}^x, \phi_{mk}^x) \\ &= \frac{p(x_{mk} | \Theta) \mathcal{CN}(x_{mk} | \hat{r}_{mk}, \tau_{mk}^r)}{\int p(x_{mk} | \Theta) \mathcal{CN}(x_{mk} | \hat{r}_{mk}, \tau_{mk}^r) dx_{mk}}, \end{aligned} \quad (15)$$

where

$$\mu_{mk}^x \triangleq \frac{\hat{r}_{mk}}{\alpha_{mk} \tau_{mk}^r + 1}, \quad (16)$$

$$\phi_{mk}^x \triangleq \frac{\tau_{mk}^r}{\alpha_{mk} \tau_{mk}^r + 1}. \quad (17)$$

in which \hat{r}_{mk} denotes the approximated Gaussian noise corrupted version of x_{mk} and τ_{mk}^r denotes the variance of noise. Besides, GAMP approximates the posterior distribution of u_{mn} as

$$\begin{aligned} p(u_{mn} | \mathbf{y}, \hat{p}_{mn}, \tau_{mn}^p, \Theta) &= \mathcal{CN}(u_{mn} | \mu_{mn}^u, \phi_{mn}^u) \\ &= \frac{p(y_{mn} | u_{mn}, \Theta) \mathcal{CN}(u_{mn} | \hat{p}_{mn}, \tau_{mn}^p)}{\int p(y_{mn} | u_{mn}, \Theta) \mathcal{CN}(u_{mn} | \hat{p}_{mn}, \tau_{mn}^p) du_{mn}}, \end{aligned} \quad (18)$$

where

$$\mu_{mn}^u \triangleq \frac{\tau_{mn}^p \beta y_{mn} + \hat{p}_{mn}}{\beta \tau_{mn}^p + 1}, \quad (19)$$

$$\phi_{mn}^u \triangleq \frac{\tau_{mn}^p}{\beta \tau_{mn}^p + 1}, \quad (20)$$

in which \hat{p}_{mn} denotes the approximated Gaussian noise corrupted version of u_{mn} and τ_{mn}^p denotes the variance of noise. Consequently, we harvest the approximated posterior distributions of \mathbf{x} and \mathbf{u} replacing the matrix inversion with the GAMP algorithm. The computational complexity is reduced to the order of $\mathcal{O}(NK M^2)$.

SE-GAMP-SBL algorithm performs weighted sum of multiple observed signals by using diversity combining technique. Since multiple observation signals inherently share the same sparse pattern, the combined signal is much closer to the original signal \mathbf{z} . Following the SE-SBL algorithm, SE-GAMP-SBL adopts the MRC scheme to achieve the diversity gain, i.e., $\hat{\mathbf{z}} = (\hat{z}_1, \hat{z}_2, \dots, \hat{z}_K)^T$ with $\hat{z}_k = \sum_{m=1}^M \eta_m \mu_{mk}^x$. Then, the combined observation symbols are expanded into block-wise patterns, $\hat{\mathbf{x}}_{mk} = \hat{\mu}_{mk}^x = \hat{z}_k, \forall m, \forall k$.

The SE-GAMP-SBL algorithm is also concerned with updating the hyperparameters Θ . The Q -function indicating the log-likelihood expectation is calculated and split into two parts based on the required estimated hyperparameters,

$$\begin{aligned} Q(\alpha, \beta | \alpha^{(t-1)}, \beta^{(t-1)}) &= \mathbb{E}[\ln p(\mathbf{y}, \mathbf{x}, \alpha, \beta)] \\ &= \mathbb{E}[\ln p(\mathbf{x}; \alpha) p(\alpha)] + \mathbb{E}[\ln p(\mathbf{y} | \mathbf{x}; \beta) p(\beta)]. \end{aligned} \quad (21)$$

For the Q -function containing only the α term, we have

$$\begin{aligned} \mathbb{E}[\ln p(\mathbf{x}; \boldsymbol{\alpha})p(\boldsymbol{\alpha})] &= \sum_{m=1}^M \sum_{k=1}^K \{(a-1) \ln \alpha_{mk} - b\alpha_{mk} \\ &+ \frac{1}{2} \ln \alpha_{mk} - \frac{1}{2} \alpha_{mk} \mathbb{E}[(\hat{x}_{mk})^2]\} + \text{const.} \end{aligned} \quad (22)$$

Since $p(\hat{x}_{mk} | \mathbf{y}, \hat{r}_{mk}, \tau_{mk}^r, \Theta)$ approximates the Gaussian distribution with posterior mean $\hat{\mu}_{mk}^x$ and variance ϕ_{mk}^x , $\mathbb{E}[(\hat{x}_{mk})^2]$ can be expressed as

$$\begin{aligned} \mathbb{E}[(\hat{x}_{mk})^2] &= (\hat{\mu}_{mk}^x)^2 + \phi_{mk}^x \\ &= \left(\sum_{m=1}^M \eta_m \mu_{mk}^x \right)^2 + \frac{\tau_{mk}^r}{\alpha_{mk} \tau_{mk}^r + 1}. \end{aligned} \quad (23)$$

Combining (22) and (23), taking the first-order derivative of $\mathbb{E}[\ln p(\mathbf{x}; \boldsymbol{\alpha})p(\boldsymbol{\alpha})]$ over α and setting it to zero, we obtain the updated α ,

$$\alpha_{mk}^{(t)} = \frac{2a-1}{2b + \mathbb{E}[(\hat{x}_{mk})^2]}. \quad (24)$$

For the Q -function containing only the β term, we have

$$\begin{aligned} \mathbb{E}[\ln p(\mathbf{y} | \mathbf{x}; \beta)p(\beta)] &= \frac{NM}{2} \ln \beta + (c-1) \ln \beta \\ &- \frac{\beta}{2} \sum_{m=1}^M \sum_{n=1}^N \mathbb{E}[(y_{mn} - u_{mn})^2] - d\beta + \text{const.} \end{aligned} \quad (25)$$

Since $p(u_{mn} | \mathbf{y}, \hat{p}_{mn}, \tau_{mn}^p, \Theta)$ approximates the Gaussian distribution with posterior mean μ_{mn}^u and variance ϕ_{mn}^u , $\mathbb{E}[(y_{mn} - u_{mn})^2]$ can be expressed as

$$\begin{aligned} \mathbb{E}[(y_{mn} - u_{mn})^2] &= (y_{mn} - \mu_{mn}^u)^2 + \phi_{mn}^u \\ &= \left(y_{mn} - \frac{\tau_{mn}^p \beta y_{mn} + \hat{p}_{mn}}{\beta \tau_{mn}^p + 1} \right)^2 + \frac{\tau_{mn}^p}{\beta \tau_{mn}^p + 1}. \end{aligned} \quad (26)$$

Combining (25) and (26), taking the first-order derivative of $\mathbb{E}[\ln p(\mathbf{y} | \mathbf{x}; \beta)p(\beta)]$ over β and setting it to zero, we obtain the updated β ,

$$\beta^{(t)} = \frac{NM + 2c - 2}{2d + \sum_{m=1}^M \sum_{n=1}^N \mathbb{E}[(y_{mn} - u_{mn})^2]}. \quad (27)$$

The proposed SE-GAMP-SBL approach dedicated to MUD with multiantenna reception is summarized in Algorithm 2.

Remark 1: The proposed algorithms use diversity combining technique to improve the posterior distribution without modifying the prior distribution, which amounts to integrating the spatial structure of user activity from a novel perspective.

Remark 2: The proposed algorithms do not require knowledge of user activity factor and noise power.

V. SIMULATION RESULTS

In this section, the performance of the proposed SE-SBL and SE-GAMP-SBL is analyzed and compared with the classical LS, MMSE, RD [6], LD [6], SBL [3], GAMP-SBL [15], B-OMP [10], B-CoSaMP [11], B-SP [12], and B-SBL [10] algorithms. Oracle-LS is introduced as a benchmark [4], where the actual active users are known at the BS. Similar to

Algorithm 2 SE-GAMP-SBL

Initialization: $\forall m, k: \alpha_{mk}^{(0)} = 1, \beta^{(0)} = 1, \hat{\mu}_{mk}^{x(0)} = 1, \phi_{mk}^{x(0)} = 1, \hat{s}_{mn}^{(0)} = 0$, and $t = 1$.

Repeat:

Step 1. $\forall m \in \{1, 2, \dots, M\}$ and $\forall n \in \{1, 2, \dots, N\}$:

$$u_{mn}^{(t)} = \sum_{m=1}^M \sum_{k=1}^K H_{mn, mk} \hat{\mu}_{mk}^{x(t-1)}$$

$$\tau_{mn}^{p(t)} = \sum_{m=1}^M \sum_{k=1}^K H_{mn, mk}^2 \phi_{mk}^{x(t-1)}$$

$$\hat{p}_{mn}^{(t)} = u_{mn}^{(t)} - \tau_{mn}^{p(t)} \hat{s}_{mn}^{(t-1)}$$

Step 2. $\forall m \in \{1, 2, \dots, M\}$ and $\forall n \in \{1, 2, \dots, N\}$:

$$\hat{s}_{mn}^{(t)} = \frac{1}{\tau_{mn}^{p(t)} \beta^{(t-1)} + 1} (\tau_{mn}^{p(t)} \beta^{(t-1)} y_{mn} + \hat{p}_{mn}^{(t)} - \hat{p}_{mn}^{(t-1)})$$

$$\tau_{mn}^{s(t)} = \frac{\beta^{(t-1)}}{\beta^{(t-1)} \tau_{mn}^{p(t)} + 1}$$

Step 3. $\forall m \in \{1, 2, \dots, M\}$ and $\forall k \in \{1, 2, \dots, K\}$:

$$\tau_{mk}^{r(t)} = \left(\sum_{m=1}^M \sum_{n=1}^N H_{mn, mk}^2 \tau_{mn}^{s(t)} \right)^{-1}$$

$$\hat{r}_{mk}^{(t)} = \hat{\mu}_{mk}^{x(t-1)} + \tau_{mk}^{r(t)} \sum_{m=1}^M \sum_{n=1}^N H_{mn, mk} \hat{s}_{mn}^{(t)}$$

Step 4. $\forall m \in \{1, 2, \dots, M\}$ and $\forall k \in \{1, 2, \dots, K\}$:

$$\mu_{mk}^{x(t)} = \frac{\hat{r}_{mk}^{(t)}}{\alpha_{mk}^{(t-1)} \tau_{mk}^{r(t)} + 1}$$

$$\phi_{mk}^{x(t)} = \frac{\tau_{mk}^{r(t)}}{\alpha_{mk}^{(t-1)} \tau_{mk}^{r(t)} + 1}$$

Step 5. $\forall m \in \{1, 2, \dots, M\}$ and $\forall k \in \{1, 2, \dots, K\}$:

$$\hat{z}_k^{(t)} = \sum_{m=1}^M \eta_m^{(t)} \mu_{mk}^{x(t)}$$

$$\hat{x}_{mk}^{(t)} = \hat{\mu}_{mk}^{x(t)} = \hat{z}_k^{(t)}$$

Step 6. Updating of α :

$$\alpha_{mk}^{(t)} = \frac{2a-1}{2b + \mathbb{E}[(\hat{x}_{mk}^{(t)})^2]}$$

Step 7. Updating of β :

$$\beta^{(t)} = \frac{NM + 2c - 2}{2d + \sum_{m=1}^M \sum_{n=1}^N \mathbb{E}[(y_{mn} - u_{mn}^{(t)})^2]}$$

Repeating the above steps until $\sqrt{\sum_{m=1}^M \sum_{k=1}^K |\hat{\mu}_{mk}^{x(t)} - \hat{\mu}_{mk}^{x(t-1)}|^2} \leq \varepsilon$, where ε is a prescribed threshold value of 10^{-6} .

the existing literatures [3], [4], [6], the adaptive gain control scheme is adopted by users to compensate the path loss, i.e., $\mathbf{h}_k^m \sim \mathcal{CN}(0, \mathbf{I})$. The number of potential users is $K = 200$, the user activity factor is $p_a = 0.1$, and the length of the spreading sequence is $N = 120$. Binary phase shift keying is considered in this simulation. The parameter values for the Gamma distribution are $a = 1.5$ and $b = c = d = 10^{-8}$.

Fig. 2 illustrates the effect of receiving antenna number on the SER performance of the proposed algorithms. For comparison, the case of single antenna $M = 1$ is also included. We can observe that increasing the number of antennas M brings significant performance improvement. Moreover, SE-SBL and SE-GAMP-SBL perform approximately the same, indicating that although these two algorithms employ different strategies to compute the posterior distribution, they both approach the true posterior distribution.

Fig. 3 shows the SER performance comparison of the considered algorithms. The proposed SE-SBL and SE-GAMP-SBL significantly outperform the block greedy-type CS algorithms. Moreover, in high SNR regime, the proposed algorithms exhibit a lower SER than the B-SBL and Oracle-

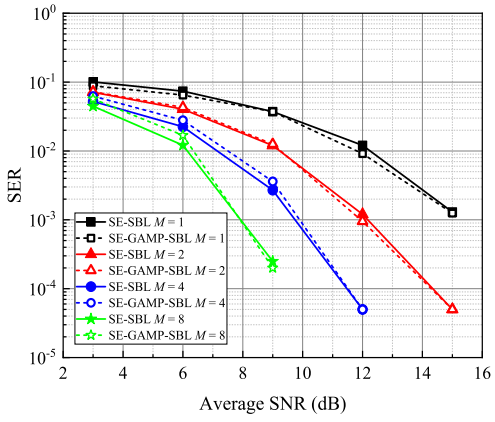


Fig. 2. SER performance of the proposed algorithms using different number of receiving antennas.

TABLE I
COMPLEXITY COMPARISON

Algorithm	Number of multiplications in each iteration	Runtime (ms)
LS	$NKM^2 + 2NK^2M^3 + K^3M^3$	226.539
MMSE	$KM + NK^2M^2 + 2NK^2M^3 + K^3M^3$	237.013
RD	$4KM + NK^2M^2 + 2NK^2M^3 + K^3M^3$	242.793
LD	polynomial time	3238.326
SBL	$NM + 2KM + 2NK^2M^2 + 4NK^2M^3 + K^3M^3$	288.708
GAMP-SBL	$9NM + 8KM + 6NK^2M^2$	10.818
B-OMP	$K + 2KM + N(2K + t)M^2 + 2Nt^2M^3 + t^3M^3$	3.476
B-CoSaMP	$3K + 4KM + N(2K + 3K_a)M^2 + 10NK_a^2M^3 + 9K_a^3M^3$	18.406
B-SP	$3K + 4KM + 2N(K + K_a)M^2 + 4NK_a^2M^3 + 2K_a^3M^3$	14.061
B-SBL	$NM + KM + (2N + 1)KM^2 + 4NK^2M^3 + K^3M^3$	295.056
SE-SBL	$(N + 1)M + 4KM + 2NK^2M^2 + 4NK^2M^3 + K^3M^3$	292.131
SE-GAMP-SBL	$(9N + 1)M + 10KM + 6NK^2M^2$	12.025
Oracle-LS	$KM + NK_aM^2 + 2NK_a^2M^3 + K_a^3M^3$	5.036

LS counterparts. This implies that improving the posterior distribution brings more benefits than modifying the prior distribution and knowing the actual active users. Particularly, SBL and GAMP-SBL can be regarded as the special case of the proposed algorithms by setting $M = 1$ as shown in Fig. 2 and Fig. 3.

Finally, we compare the computational complexity of the considered algorithms. Table I presents the computational cost and running time required to execute single iteration for $M = 4$. LD is a polynomial time algorithm that can be readily accomplished using the available CVX toolbox. The computational cost of the proposed SE-SBL is higher than the block greedy-type CS algorithms due to $K_a \ll K$. However, the proposed SE-GAMP-SBL has the ability to obtain superiority in computational complexity mitigation. Moreover, the computational costs of SE-SBL and SE-GAMP-SBL are in the same order as those of SBL and GAMP-SBL respectively, because the complexity of Bayesian algorithms mainly stems from computing the posterior distribution.

VI. CONCLUSION

In this letter, we proposed two spatially enhanced Bayesian CS algorithms for efficient MUD in uplink grant-free MIMO-NOMA system. The proposed SE-SBL and SE-GAMP-SBL algorithms exploit the full potential of shared sparsity induced by the spatial structure of user activity, and do not require prior information about user sparsity level and noise power.

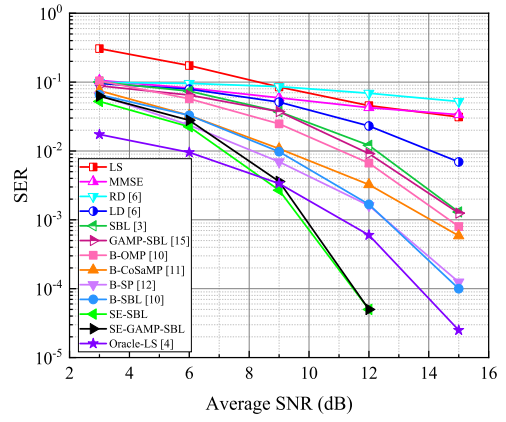


Fig. 3. SER performance comparison of the considered algorithms under receiving antenna number of $M = 4$.

Simulation results show that the proposed algorithms outperform off-the-shelf methods both in detection performance and computational complexity.

REFERENCES

- [1] X. Chen, D. W. K. Ng, W. Yu, E. G. Larsson, N. A. Dahir, and R. Schober, "Massive access for 5G and beyond," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 3, pp. 615-636, Mar. 2021.
- [2] M. Ke, Z. Gao, Y. Wu, X. Gao, and R. Schober, "Compressive sensing-based adaptive active user detection and channel estimation: Massive access meets massive MIMO," *IEEE Trans. Signal Process.*, vol. 68, pp. 764-779, 2020.
- [3] X. Zhang, P. Fan, J. Liu, and L. Hao, "Bayesian learning-based multiuser detection for grant-free NOMA systems," *IEEE Trans. Wireless Commun.*, vol. 21, no. 8, pp. 6317-6328, Aug. 2022.
- [4] L. Wu, Z. Wang, P. Sun, and Y. Yang, "Temporal correlation enhanced sparse activity detection in MIMO enabled grant-free NOMA," *IEEE Trans. Veh. Technol.*, vol. 71, no. 3, pp. 2887-2899, Mar. 2022.
- [5] T. Hara and K. Ishibashi, "Grant-free non-orthogonal multiple access with multiple-antenna base station and its efficient receiver design," *IEEE Access.*, vol. 7, pp. 175717-175726, 2019.
- [6] H. Zhu and G. B. Giannakis, "Exploiting sparse user activity in multiuser detection," *IEEE Trans. Commun.*, vol. 59, no. 2, pp. 454-465, Feb. 2011.
- [7] L. Liu and W. Yu, "Massive connectivity with massive MIMO-Part I: Device activity detection and channel estimation," *IEEE Trans. Signal Process.*, vol. 66, no. 11, pp. 2933-2946, Jun. 2018.
- [8] S. Zhang, Y. Cui, and W. Chen, "Joint device activity detection, channel estimation and signal detection for massive grant-free access via BiGAMP," *IEEE Trans. Signal Process.*, vol. 71, pp. 1200-1215, 2023.
- [9] T. Ding, X. Yuan, and S. C. Liew, "Sparsity learning-based multiuser detection in grant-free massive-device multiple access," *IEEE Trans. Wireless Commun.*, vol. 18, no. 7, pp. 3569-3582, Jul. 2019.
- [10] 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.
- [11] X. Zhang, W. Xu, Y. Cui, L. Lu, and J. Lin, "On recovery of block sparse signals via block compressive sampling matching pursuit," *IEEE Access.*, vol. 7, pp. 175554-175563, 2019.
- [12] Y. Du *et al.*, "Joint channel estimation and multiuser detection for uplink grant-free NOMA," *IEEE Wireless Commun. Lett.*, vol. 7, no. 4, pp. 682-685, Aug. 2018.
- [13] D. G. Tzikas, A. C. Likas, and N. P. Galatsanos, "The variational approximation for Bayesian inference," *IEEE Signal Process. Mag.*, vol. 25, no. 6, pp. 131-146, Nov. 2008.
- [14] F. Clazzer, C. Kissling, and M. Marchese, "Enhancing contention resolution ALOHA using combining techniques," *IEEE Trans. Commun.*, vol. 66, no. 6, pp. 2576-2587, Jun. 2018.
- [15] M. Al-Shoukairi, P. Schniter, and B. D. Rao, "A GAMP-based low complexity sparse Bayesian learning algorithm," *IEEE Trans. Signal Process.*, vol. 66, no. 2, pp. 294-308, Jan. 2018.