

On Massive IoT Connectivity with Temporally-Correlated User Activity

Qipeng Wang, Liang Liu, Shuowen Zhang, and Francis C. M. Lau
Department of Electronic and Information Engineering
The Hong Kong Polytechnic University
Emails: qipeng.wang@connect.polyu.hk
{liang-eie.liu,shuowen.zhang,francis-cm.lau}@polyu.edu.hk

Abstract—This paper considers joint device activity detection and channel estimation in Internet of Things (IoT) networks, where a large number of IoT devices exist but merely a random subset of them become active for short-packet transmission at each time slot. In particular, to improve the detection performance, we propose to leverage the *temporal correlation* in user activity, i.e., a device active at the previous time slot is more likely to be still active at the current time slot. Despite the appealing temporal correlation feature, it is challenging to unveil the connection between the estimated activity pattern for the previous time slot (which may be imperfect) and the true activity pattern at the current time slot due to the unknown estimation error. In this paper, we manage to tackle this challenge under the framework of approximate message passing (AMP). Specifically, thanks to the state evolution, the correlation between the activity pattern estimated by AMP at the previous time slot and the real activity pattern at the previous and current time slot is quantified explicitly. Based on the well-defined temporal correlation, we further manage to embed this useful SI into the design of the minimum mean-squared error (MMSE) denoisers and log-likelihood ratio (LLR) test based activity detectors under the AMP framework. Theoretical comparison between the SI-aided AMP algorithm and its counterpart without utilizing temporal correlation is provided. Moreover, numerical results are given which show the significant gain in activity detection accuracy brought by the SI-aided algorithm.

I. INTRODUCTION

A typical massive Internet of Things (IoT) connectivity system consists of a large number of low-cost devices, each of which stays in the silence mode for a long period to save the energy and becomes active merely when triggered by the unusual events. Under such a setting, one key challenge lies in how to jointly identify the randomly active IoT devices and estimate their channels in a fast and accurate manner [1]. Recently, it was shown that the above job can be accomplished by utilizing the compressed sensing technique thanks to the sparse user activity [2]–[7]. In particular, under the framework of multiple measurement vector (MMV) based approximate message passing (AMP) [8], [9], it has been shown in [2] that the activity detection error probability decreases significantly with the number of antennas at the base station (BS). Such an exciting result arises from reaping the spatial correlation in user activity: if one device is active for one antenna, it is also active for all the other antennas. However, this theoretical performance gain is achieved at the cost of high computational complexity in practice: in an IoT system with a large number

of devices and BS antennas, the dimension of the data to be processed is tremendous. A nature question is: if only a small number of antennas are utilized to reduce the computational complexity, is it still possible to achieve high-quality device activity detection and channel estimation?

This paper provides an affirm answer to the above question. The core is to utilize the temporal correlation in user activity to compensate for the spatial correlation gain of the MMV-based AMP algorithm. In practice, temporally-correlated user activity may come from the fact that if an abnormal event is detected by some sensor at a moment, then this device is more likely to be still activated by this event in the near future. To fully take advantage of this temporal correlation, this paper aims to design a side information (SI) aided AMP framework.

Note that at each time slot, the available information at the BS is the imperfect device activity pattern estimated at the previous time slot, whose connection to the real device activity pattern at the previous or current time slot is unclear in general due to the unknown estimation error. As a result, despite the existence of temporal correlation in user activity, it is a challenging task to utilize this as SI to improve the performance of activity detection. In this work, we point out that under the framework of AMP, the correlation between the device activity pattern estimated at the previous time slot and the real one at the current time slot can be explicitly quantified thanks to the state evolution. Based on this correlation, we further manage to design the minimum mean-squared error (MMSE) denoisers and log-likelihood ratio (LLR) based detectors in AMP with SI taken into consideration. The impact of using SI in AMP is illustrated both theoretically and numerically.

In the literature, temporal correlation has been also utilized in [10], [11] to design the AMP algorithm. However, different from [10] that considers a Turbo extension of the AMP algorithm [12] based on the idea of factor graph, our approach provides a framework to incorporate SI into the MMV-AMP algorithm without needing to craft the graph model for each new signal. Note that [10] merely works for the single measurement vector (SMV) case, i.e., the BS has a single antenna. More importantly, even for the special case of the SMV problem, our approach can be shown to be Bayes-optimal, thus yielding improved detection performance. Moreover, compared with [11] whose emphasize is on estimating the sparse channels, our paper places more focus on device

activity detection. As a result, dedicated activity detectors are proposed based on the SI and the LLR test.

II. SYSTEM MODEL

Baseband Model: This paper considers the uplink communication in a massive IoT connectivity system consisting of one BS equipped with M antennas and N single-antenna IoT devices. We assume quasi-static block-fading channels, in which all user channels remain approximately constant in each coherence block, but vary independently from block to block. Let J denote the number of consecutive coherence blocks considered in this work. At coherence block j , the channel from device n to the BS is denoted by $\mathbf{h}_n^{(j)} \in \mathbb{C}^{M \times 1}$, $j = 1, \dots, J$, $n = 1, \dots, N$. It is assumed that the user channels follow the independent and identically distributed (i.i.d.) Rayleigh fading channel model, i.e., $\mathbf{h}_n^{(j)} \in \mathcal{CN}(\mathbf{0}, \gamma_n \mathbf{I})$, $\forall j, n$, where γ_n is the path loss of device n . Note that $\mathbf{h}_n^{(j)}$'s are independent over n and j .

Due to the sporadic data traffic in IoT networks, only a small set of devices become active in each coherence block. We define the user activity indicator functions as follows:

$$\delta_n^{(j)} = \begin{cases} 1, & \text{if user } n \text{ is active at coherence block } j, \\ 0, & \text{otherwise,} \end{cases} \quad \forall j, n, \quad (1)$$

so that $\delta_n^{(j)}$ is a Bernoulli random variable with

$$Pr(\delta_n^{(j)} = 1) = \lambda, \quad Pr(\delta_n^{(j)} = 0) = 1 - \lambda, \quad \forall n, j. \quad (2)$$

In this work, we consider the grant-free random access scheme [1] in our interested IoT system, where at the beginning of each coherence block, the active devices transmit their pilot sequences to the BS to perform joint device activity detection and channel estimation. Let $\mathbf{s}_n = [s_{n,1}, \dots, s_{n,L}]^T \in \mathbb{C}^{L \times 1}$ denote the pilot sequence with length L assigned to device n , $\forall n$. Similar to [1]–[3], it is assumed that all the entries in \mathbf{s}_n are generated from i.i.d. complex Gaussian distribution with zero mean and variance $1/L$, $\forall n$. Then, the BS received signal at coherence block j is expressed as

$$\mathbf{Y}^{(j)} = \sum_{n=1}^N \delta_n^{(j)} \mathbf{h}_n^{(j)} \mathbf{s}_n + \mathbf{Z}^{(j)} = \mathbf{S} \mathbf{X}^{(j)} + \mathbf{Z}^{(j)}, \quad \forall j, \quad (3)$$

where $\mathbf{Z}^{(j)} \in \mathbb{C}^{L \times M} \in \mathcal{CN}(\mathbf{0}, \sigma_z^2 \mathbf{I})$ is the additive white Gaussian noise (AWGN) of the BS at coherence block j , $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_N] \in \mathbb{C}^{L \times N}$, and $\mathbf{X}^{(j)} = [\mathbf{x}_1^{(j)}, \dots, \mathbf{x}_N^{(j)}]^T \in \mathbb{C}^{N \times M}$ with $\mathbf{x}_n^{(j)} = \delta_n^{(j)} \mathbf{h}_n^{(j)}$ denoting the effective channel of device n at coherence block j , $\forall j, n$. At each coherence block j , the job of the BS is to jointly detect the active devices and estimate their channels by estimating $\mathbf{X}^{(j)}$ based on its received signal $\mathbf{Y}^{(j)}$ and its knowledge of the user pilots \mathbf{S} .

Temporally-Correlated User Activity Model: This paper considers the case of temporally-correlated user activity, which

is modeled by a Markov chain with the following transition probabilities:

$$\begin{aligned} Pr(\delta_n^{(j)} = 1 \mid \delta_n^{(j-1)} = 1) &= \alpha, \\ Pr(\delta_n^{(j)} = 0 \mid \delta_n^{(j-1)} = 1) &= 1 - \alpha, \\ Pr(\delta_n^{(j)} = 1 \mid \delta_n^{(j-1)} = 0) &= \beta, \\ Pr(\delta_n^{(j)} = 0 \mid \delta_n^{(j-1)} = 0) &= 1 - \beta, \end{aligned} \quad \forall j, n. \quad (4)$$

In other words, if user n is active in coherence block $j - 1$, then with probability α , it is still active in coherence block j ; if user n is inactive in coherence block $j - 1$, then with probability β , it is active in coherence block j . Given the above temporal correlation, over two consecutive coherence blocks $j - 1$ and j , we have the following four cases to model each device's activity:

Case 1: An user is active for both coherence blocks $j - 1$ and j , i.e., $\mathbf{x}_n^{(j-1)} = \mathbf{h}_n^{(j-1)}$ and $\mathbf{x}_n^{(j)} = \mathbf{h}_n^{(j)}$, with probability $\alpha\lambda$.

Case 2: An user is active at coherence block $j - 1$, but becomes inactive at coherence block j , i.e., $\mathbf{x}_n^{(j-1)} = \mathbf{h}_n^{(j-1)}$ and $\mathbf{x}_n^{(j)} = \mathbf{0}$, with probability $(1 - \alpha)\lambda$.

Case 3: An user is inactive at coherence block $j - 1$, but becomes active at coherence block j , i.e., $\mathbf{x}_n^{(j-1)} = \mathbf{0}$ and $\mathbf{x}_n^{(j)} = \mathbf{h}_n^{(j)}$, with probability $\beta(1 - \lambda)$.

Case 4: An user is inactive for both coherence blocks $j - 1$ and j , i.e., $\mathbf{x}_n^{(j-1)} = \mathbf{0}$ and $\mathbf{x}_n^{(j)} = \mathbf{0}$, with probability $(1 - \beta)(1 - \lambda)$.

Similar to [10], [11], we assume that each Markov chain operates in steady-state such that the probability that a device becomes active is λ over all the J coherence blocks, i.e., (2). Under this condition, the relation between α and β is given by $\alpha\lambda + \beta(1 - \lambda) = \lambda$. Due to this relation, the Markov chains are completely characterized by two parameters λ and α .

Under the temporal correlation modeled by (4), we should not detect the user activity over consecutive coherence blocks in an independent manner as in [2], [3], since the user activity at the previous coherence block can provide SI for improving the estimation accuracy at the current coherence block. However, at each coherence block j , only an imperfect estimation of the device activity at coherence block $j - 1$, denoted by $\hat{\delta}_n^{(j-1)}$, $\forall n$, is available at the BS. Despite the temporal correlation shown in (4), it is non-trivial to model a precise statistical relation between $\delta_n^{(j)}$ and $\hat{\delta}_n^{(j-1)}$, $\forall n$, since the connection between $\delta_n^{(j-1)}$ and $\hat{\delta}_n^{(j-1)}$, $\forall n$, is in general unknown. Without such a relation characterization, it is possible that the imperfect estimation at the previous coherence block is not properly utilized, which may even degrade the estimation performance at the current coherence block. This motivates us to study a systematic approach that is able to always leverage SI to improve the performance of activity detection and channel estimation.

Note that in the case without using SI, [2] and [3] showed that the estimation of $\mathbf{X}^{(j)}$ based on (3) is a compressed sensing problem, since many rows in $\mathbf{X}^{(j)}$ are zero vectors due to the sparse user activity. Moreover, the MMV-AMP

algorithm has been used to estimate the row-sparse matrix $\mathbf{X}^{(j)}$ at each coherence block. In the rest of this paper, we study connection of SI to the current estimation and the method to embed SI into the AMP algorithm design.

III. LEVERAGING SI IN AMP

This paper adopts the framework proposed in [11] to integrate SI into the MMV-AMP algorithm. At coherence block j , the SI-aided MMV-AMP algorithm will generate an estimation of $\mathbf{X}^{(j)}$, denoted by $\hat{\mathbf{X}}^{(j)} = [\hat{\mathbf{x}}_1^{(j)}, \dots, \hat{\mathbf{x}}_N^{(j)}]^T$, based on the signal received at the current time slot (3) and the estimation made by SI-aided MMV-AMP algorithm at the previous coherence block, i.e., $\hat{\mathbf{X}}^{(j-1)}$. Specifically, at coherence block j , the SI-aided MMV-AMP algorithm starts from $\mathbf{X}_0^{(j)} = \mathbf{0}$ and $\mathbf{R}_0^{(j)} = \mathbf{Y}^{(j)}$ and iterates as follows:

$$\mathbf{x}_{n,t+1}^{(j)} = \eta_{n,t}^{(j)} \left(\mathbf{x}_{n,t}^{(j)} + \left(\mathbf{R}_t^{(j)} \right)^H \mathbf{s}_n, f_{n,j} \left(\hat{\mathbf{x}}_n^{(j-1)} \right) \right), \quad (5)$$

$$\begin{aligned} \mathbf{R}_{t+1}^{(j)} &= \mathbf{Y}^{(j)} - \mathbf{S} \mathbf{X}_{t+1}^{(j)} \\ &+ \frac{N}{L} \mathbf{R}_t^{(j)} \left\langle \eta_{n,t}^{(j)'} \left(\mathbf{x}_{n,t}^{(j)} + \left(\mathbf{R}_t^{(j)} \right)^H \mathbf{s}_n, f_{n,j} \left(\hat{\mathbf{x}}_n^{(j-1)} \right) \right) \right\rangle. \end{aligned} \quad (6)$$

In (5) and (6), $t = 0, 1, \dots$ denotes the index of algorithm iteration, $\mathbf{X}_t^{(j)} = [\mathbf{x}_{1,t}^{(j)}, \dots, \mathbf{x}_{N,t}^{(j)}]^T$ denotes the estimation of $\mathbf{X}^{(j)}$ at the t -th iteration of the AMP algorithm, $f_{n,j}(\hat{\mathbf{x}}_n^{(j-1)})$ is a function of $\hat{\mathbf{x}}_n^{(j-1)}$ which is used as the SI for device n , $\mathbf{R}_t^{(j)}$ is the corresponding residual at iteration t , $\eta_{n,t}^{(j)}(\cdot, \diamond) \in \mathbb{C}^{M \times 1}$ is the denoising function for device n , $\eta_{n,t}^{(j)'}(\cdot, \diamond)$ is the first-order derivative of $\eta_{n,t}^{(j)}(\cdot, \diamond)$ with respect to the first variable \cdot , and $\langle \cdot \rangle$ is the averaging operation over all entries of $\eta_{n,t}^{(j)'}(\cdot, \diamond)$. Let $\mathbf{X}_\infty^{(j)} = [\mathbf{x}_{1,\infty}^{(j)}, \dots, \mathbf{x}_{N,\infty}^{(j)}]^T$ and $\mathbf{R}_\infty^{(j)}$ denote the estimation of $\mathbf{x}_n^{(j)}$ and the corresponding residual after the convergence of the SI-aided MMV-AMP algorithm at coherence block j . Then, we have $\hat{\mathbf{x}}_n^{(j-1)} = \mathbf{x}_{n,\infty}^{(j-1)}$, $\forall j, n$.

Different from the conventional MMV-AMP algorithm [8], [9] used in [2], [3], the estimation at the previous coherence block is utilized for denoiser design as given by (5) under our considered SI-aided MMV-AMP algorithm. To implement this algorithm in IoT systems with temporally-correlated activity, in the following, we introduce what SI should be extracted from previous estimation, i.e., the design of $f_{n,j}(\hat{\mathbf{x}}_n^{(j-1)})$, and how to design the denoisers based on this SI.

A. Identifying SI from State Evolution

According to [11], with the SI-aided MMV-AMP algorithm shown in (5) and (6), there exists the state evolution in the asymptotic regime where $N, L \rightarrow \infty$ with fixed N/L . Specifically, at each iteration t of AMP to estimate $\mathbf{X}^{(j)}$, $\mathbf{x}_{n,t}^{(j)} + \left(\mathbf{R}_t^{(j)} \right)^H \mathbf{s}_n$ is statistically equivalent to:

$$\tilde{\mathbf{x}}_{n,t}^{(j)} = \mathbf{x}_n^{(j)} + \left(\Sigma_t^{(j)} \right)^{\frac{1}{2}} \mathbf{v}_n^{(j)}, \quad \forall n, j, \quad (7)$$

where $\mathbf{v}_n^{(j)} \in \mathbb{C}^{M \times 1} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I})$ is the noise independent of $\mathbf{x}_n^{(j)}$ and $\Sigma_t^{(j)} \in \mathbb{C}^{M \times M}$ is the *state*. Define a set of random

vectors $\mathbf{X}_n^{(j)} \in \mathbb{C}^{M \times 1}$, $\mathbf{V}_n^{(j)} \in \mathbb{C}^{M \times 1}$, and $\hat{\mathbf{X}}_n^{(j-1)} \in \mathbb{C}^{M \times 1}$ which capture the distribution of $\mathbf{x}_n^{(j)}$, $\mathbf{v}_n^{(j)}$, and $\hat{\mathbf{x}}_n^{(j-1)}$, respectively, $\forall j, n$. The state evolution is given by:

$$\begin{aligned} \Sigma_{t+1}^{(j)} &= \sigma_z^2 \mathbf{I} + \\ &\frac{N}{L} \mathbb{E} \left[\left(\eta_{n,t}^{(j)} \left(\mathbf{X}_n^{(j)} + \left(\Sigma_t^{(j)} \right)^{\frac{1}{2}} \mathbf{V}_n^{(j)}, f_{n,j} \left(\hat{\mathbf{X}}_n^{(j-1)} \right) \right) - \mathbf{X}_n^{(j)} \right)^H \right. \\ &\quad \left. \left(\eta_{n,t}^{(j)} \left(\mathbf{X}_n^{(j)} + \left(\Sigma_t^{(j)} \right)^{\frac{1}{2}} \mathbf{V}_n^{(j)}, f_{n,j} \left(\hat{\mathbf{X}}_n^{(j-1)} \right) \right) - \mathbf{X}_n^{(j)} \right) \right]. \end{aligned} \quad (8)$$

Note that at coherence block j , we already have the estimation of $\mathbf{X}^{(j-1)}$, i.e., $\hat{\mathbf{x}}_n^{(j-1)} = \mathbf{x}_{n,\infty}^{(j-1)}$, $\forall n$. According to (7), $\hat{\mathbf{x}}_n^{(j-1)} + \left(\mathbf{R}_\infty^{(j-1)} \right)^H \mathbf{s}_n$ is statistically equivalent to:

$$\tilde{\mathbf{x}}_{n,\infty}^{(j-1)} = \mathbf{x}_n^{(j-1)} + \left(\Sigma_\infty^{(j-1)} \right)^{\frac{1}{2}} \mathbf{v}_n^{(j-1)}, \quad \forall n, j, \quad (9)$$

where $\Sigma_\infty^{(j-1)}$ denotes the state of the AMP algorithm (8) after it converges at coherence block $j-1$. It is worth noting that (9) reveals the correlation between $\tilde{\mathbf{x}}_{n,\infty}^{(j-1)}$ and $\mathbf{x}_n^{(j-1)}$, while (4) reveals the correlation between $\tilde{\mathbf{x}}_{n,\infty}^{(j-1)}$ and $\mathbf{x}_n^{(j)}$. Thus, the correlation between $\tilde{\mathbf{x}}_{n,\infty}^{(j-1)}$ and $\mathbf{x}_n^{(j)}$ can be built for all the devices. This motivates us to adopt the following SI to design the denoisers in the AMP algorithm:

$$f_{n,j} \left(\hat{\mathbf{x}}_n^{(j-1)} \right) = \hat{\mathbf{x}}_n^{(j-1)} + \left(\mathbf{R}_\infty^{(j-1)} \right)^H \mathbf{s}_n, \quad \forall n, j. \quad (10)$$

The next question is how to utilize the correlation between $\hat{\mathbf{x}}_n^{(j-1)} + \left(\mathbf{R}_\infty^{(j-1)} \right)^H \mathbf{s}_n$'s and $\mathbf{x}_n^{(j)}$'s to design the denoisers (5).

B. Leveraging SI for MMSE Denoiser Design

In this paper, we adopt the Bayesian approach to design the MMSE denoisers $\eta_{n,t}^{(j)}(\cdot, \diamond)$'s for signal recovery at each coherence block. At the $(t+1)$ -th iteration of the AMP algorithm at coherence block j , the available information includes $\mathbf{x}_{n,t}^{(j)} + \left(\mathbf{R}_t^{(j)} \right)^H \mathbf{s}_n$ from the current coherence block whose distribution is modeled by (7) and the SI from the previous coherence block $\hat{\mathbf{x}}_n^{(j-1)} + \left(\mathbf{R}_\infty^{(j-1)} \right)^H \mathbf{s}_n$ whose distribution is modeled by (9). Based on the above information, the MMSE denoisers can be expressed as

$$\begin{aligned} \mathbb{E}[\mathbf{X}_n^{(j)} | \mathbf{X}_n^{(j)} + \left(\Sigma_t^{(j)} \right)^{\frac{1}{2}} \mathbf{V}_n^{(j)} = \tilde{\mathbf{x}}_{n,t}^{(j)}, \\ \mathbf{X}_n^{(j-1)} + \left(\Sigma_\infty^{(j-1)} \right)^{\frac{1}{2}} \mathbf{V}_n^{(j-1)} = \tilde{\mathbf{x}}_{n,\infty}^{(j-1)}], \quad \forall n, j. \end{aligned} \quad (11)$$

Based on the similar approach used in [2, Appendix B], it can be shown that with the above MMSE denoisers, the matrix $\Sigma_t^{(j)}$ generated by the state evolution (8) always stays as a scaled version of the identity matrix, i.e.,

$$\Sigma_t^{(j)} = \left(\tau_t^{(j)} \right)^2 \mathbf{I}, \quad \forall t, j. \quad (12)$$

With this result, the MMSE denoisers (11) can be explicitly characterized by the following theorem.

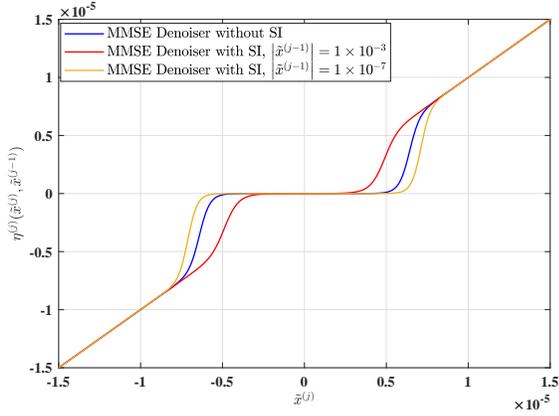


Fig. 1. Comparison of MMSE denoisers with/without using SI.

Theorem 1: Consider the SI-aided MMV-AMP algorithm given by (5) and (6) under the temporal correlation model for user activity shown in (4). Define

$$\Delta_{n,t}^{(j)} = \left(\tau_t^{(j)}\right)^{-2} - \left(\left(\tau_t^{(j)}\right)^2 + \gamma_n\right)^{-1}, \quad \forall n, t, j. \quad (13)$$

Under the asymptotic regime where $N, L \rightarrow \infty$ with fixed N/L , the MMSE denoisers (11) at coherence block j with the SI given in (10) are expressed as:

$$\eta_{t,n}^{(j)}(\tilde{\mathbf{x}}_{n,t}^{(j)}, \tilde{\mathbf{x}}_{n,\infty}^{(j-1)}) = \frac{\gamma_n \left(\gamma_n + \left(\tau_t^{(j)}\right)^2\right)^{-1} \tilde{\mathbf{x}}_{n,t}^{(j)}}{1 + \frac{1-\lambda}{\lambda} \mu_{n,t}^{(j)} \times \frac{\beta + (1-\beta)\mu_{n,\infty}^{(j-1)}}{\alpha + (1-\alpha)\mu_{n,\infty}^{(j-1)}}, \quad \forall n, t, j, \quad (14)$$

where

$$\mu_{n,t}^{(j)} = \left(\frac{\left(\tau_t^{(j)}\right)^2 + \gamma_n}{\left(\tau_t^{(j)}\right)^2}\right)^M \exp\left(-\Delta_{n,t}^{(j)} \|\tilde{\mathbf{x}}_{n,t}^{(j)}\|^2\right), \quad (15)$$

and $(\tau_\infty^{(j-1)})^2$ can be obtained from the state evolution (8) and (12) after AMP converges in coherence block $j-1$.

Proof: Please refer to Appendix A. ■

To gain insights from Theorem 1, we consider some special cases. First, if user activity is independent over different coherence blocks, i.e., $\alpha = \beta = \lambda$ such that $Pr(\delta_n^{(j)} | \delta_n^{(j-1)}) = Pr(\delta_n^{(j)})$, $\forall n, j$, the denoisers (14) will reduce to

$$\eta_{n,t}^{(j)}(\tilde{\mathbf{x}}_{n,t}^{(j)}, \tilde{\mathbf{x}}_{n,\infty}^{(j-1)}) = \frac{\gamma_n \left(\gamma_n + \left(\tau_t^{(j)}\right)^2\right)^{-1} \tilde{\mathbf{x}}_{n,t}^{(j)}}{1 + \frac{1-\lambda}{\lambda} \mu_{n,t}^{(j)}}, \quad \forall n, j, \quad (16)$$

which are the MMSE denoisers proposed in [2] [3] without taking SI into account. This is because if there is no temporal correlation in user activity, SI will have no effect on the MMSE denoiser design. Second, if $(\tau_\infty^{(j-1)})^2 \rightarrow \infty$, it can be shown from (15) that $\mu_{n,\infty}^{(j-1)} = 1$, $\forall n$. Then, the MMSE denoisers shown in (14) will also reduce to the MMSE denoisers (16) proposed in [2] [3] without taking SI into account. This is because according to (9) and (12), $(\tau_\infty^{(j-1)})^2$ can be viewed as the equivalent noise power for estimating $\mathbf{X}^{(j-1)}$ by AMP.

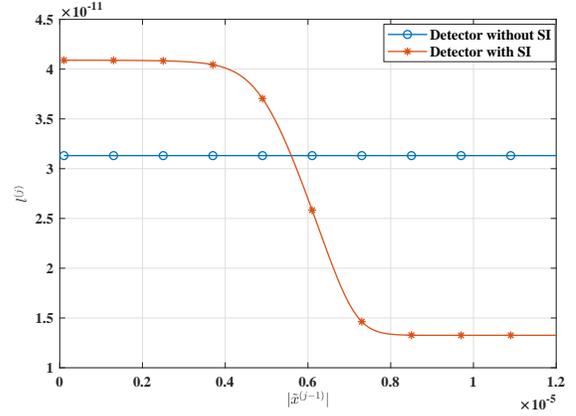


Fig. 2. Comparison of threshold for activity detectors with/without using SI.

If this noise power is infinite but the power of each row in $\mathbf{X}^{(j-1)}$ is finite, then the estimation does not provide any useful information for the estimation in the next block, despite the existence of temporal correlation in activity.

Next, we provide a numerical example to compare the denoisers using SI, i.e., (14), and without using SI, i.e., (16). In this example, we set $M = 1$, $\lambda = 0.1$, $\alpha = 0.91$, $\beta = 0.01$, $\tau_j = \tau^{(j-1)} = 2 \times 10^{-6}$, and $\gamma = \gamma = 1 \times 10^{-8}$. Fig. 1 shows the SI-aided MMSE denoisers when $|\tilde{x}^{(j-1)}| = 1 \times 10^{-3}$ and $|\tilde{x}^{(j-1)}| = 1 \times 10^{-7}$ as well as the denoiser without using SI. Compared to the denoiser without using SI, it is observed that when $|\tilde{x}^{(j-1)}|$ is larger/smaller, i.e., the user tends to be detected as an active/inactive device at the previous block, the SI-aided MMSE denoiser estimates $x^{(j)}$ as zero over a smaller/larger range of $\tilde{x}^{(j)}$, i.e., the user tends to be detected as an active/inactive device at the current time slot.

C. Leveraging SI for Activity Detector Design

After the convergence of the SI-aided MMV-AMP algorithm, the LLR test is applied to conduct device activity detection. For the hypothesis detection problem H_0 : user is inactive; H_1 : user is active, the LLR-based detector is

$$LLR_{n,t}^{(j)} = \log \left(\frac{p(\tilde{\mathbf{x}}_{n,t}^{(j)}, \tilde{\mathbf{x}}_{n,\infty}^{(j-1)} | \mathbf{x}_n^{(j)} \neq \mathbf{0})}{p(\tilde{\mathbf{x}}_{n,t}^{(j)}, \tilde{\mathbf{x}}_{n,\infty}^{(j-1)} | \mathbf{x}_n^{(j)} = \mathbf{0})} \right) \underset{H_0}{\overset{H_1}{\gtrless}} l, \quad (17)$$

where l is a common threshold of LLR for all the users over all the coherence blocks.

Theorem 2: Consider the SI-aided MMV-AMP algorithm given by (5) and (6) under the temporal correlation model for user activity shown in (4). Under the asymptotic regime where $N, L \rightarrow \infty$ with fixed N/L , the LLR-based decision rule (17) can be expressed as:

$$\|\tilde{\mathbf{x}}_{n,t}^{(j)}\|^2 \underset{H_0}{\overset{H_1}{\gtrless}} l_{n,t}^{(j)} \triangleq \frac{1}{\Delta_{n,t}^{(j)}} \left(l + M \log \left(\frac{\left(\tau_t^{(j)}\right)^2 + \gamma_n}{\left(\tau_t^{(j)}\right)^2} \right) + \log \left(\frac{\beta + (1-\beta)\mu_{n,\infty}^{(j-1)}}{\alpha + (1-\alpha)\mu_{n,\infty}^{(j-1)}} \right) \right), \quad \forall n, t, j, \quad (18)$$

where $\Delta_{n,t}^{(j)}$ and $\mu_{n,\infty}^{(j-1)}$ are given by (13) and (15), respectively.

Proof: Please refer to Appendix B. ■

Theorem 2 states that a device is detected to be active if $\|\tilde{\mathbf{x}}_{n,t}^{(j)}\|^2$ is larger than a threshold, which depends on the SI from the previous block. Similar to the MMSE denoisers in Theorem 1, if $\alpha = \beta = \lambda$, or if $\tau_\infty^{(j-1)} \rightarrow \infty$, then the LLR-based detectors (18) will reduce to:

$$\|\tilde{\mathbf{x}}_{n,t}^{(j)}\|^2 \underset{H_0}{\overset{H_1}{\gtrless}} \frac{l + M \log \left(\frac{(\tau_t^{(j)})^2 + \gamma_n}{(\tau_t^{(j)})^2} \right)}{\Delta_{n,t}^{(j)}}, \quad \forall n, t, j, \quad (19)$$

which are the detectors without using SI [3].

Next, we provide a numerical example to compare the LLR-based activity detectors using SI, i.e., (18), and those without using SI, i.e., (19). The setup is the same as that for Fig. 1. Moreover, we set $l = 0$. Fig. 2 shows the threshold in the activity detectors, i.e., $l^{(j)}$, versus different values of $|\tilde{x}^{(j-1)}|$, which is the SI from the previous time slot. Compared to the case without using SI, it is observed that if $|\tilde{x}^{(j-1)}|$ is larger/smaller, i.e., the user tends to be detected as an active/inactive user previously, the threshold in the SI-aided detectors (18) becomes smaller/larger, i.e., this user tends to be detected as an active/inactive user at the current time slot.

IV. NUMERICAL RESULTS

In this section, we provide numerical results to evaluate the performance of the proposed SI-aided MMV-AMP algorithm in massive IoT connectivity systems. We assume that there are $N = 4000$ devices randomly located in a cell of radius $R = 1000$ m. The path loss model is $-128.1 - 36.7 \log_{10}(d_n)$ in dB, where d_n in km denotes the distance from device n to the BS. We consider the communication over $J = 10$ coherence blocks, while at each time slot, we have $\lambda = 0.1$, $\alpha = 0.46$, and $\beta = 0.06$. Next, the user transmit power is 23 dBm. Last, the power spectrum density of the noise is -169 dBm/Hz, while the bandwidth of the channel is assumed to be 10 MHz.

First, we consider the case when the BS is equipped with one antenna, i.e., $M = 1$. In this case, we call our proposed algorithm as SMV-AMP with SI. In Fig. 3, we show the tradeoff between the probabilities of false alarm P_{FA} and missed detection P_{MD} , which is obtained by varying the value of l in the activity detectors. In this numerical example, we set $L = 600$. Moreover, we consider the Dynamic Compressed Sensing via Approximate Message Passing (DCS-AMP) algorithm proposed in [10] as the benchmark scheme. The DCS-AMP is implemented in filtering mode to match our setting. It is observed from Fig. 3 that under our proposed SMV-AMP algorithm with SI, the activity detection performance improves over time. This shows that the proposed algorithm is capable of intelligently exploiting the SI obtained in the previous time slots for improving the detection performance. Moreover, at time slot 5, our proposed scheme can achieve much lower detection error probability than the conventional AMP algorithm without utilizing SI. It is also observed that the proposed SI-aided SMV-AMP algorithm outperforms the DCS-AMP algorithm significantly at time slots 3 and 5 because our proposed scheme is built on the true statistical

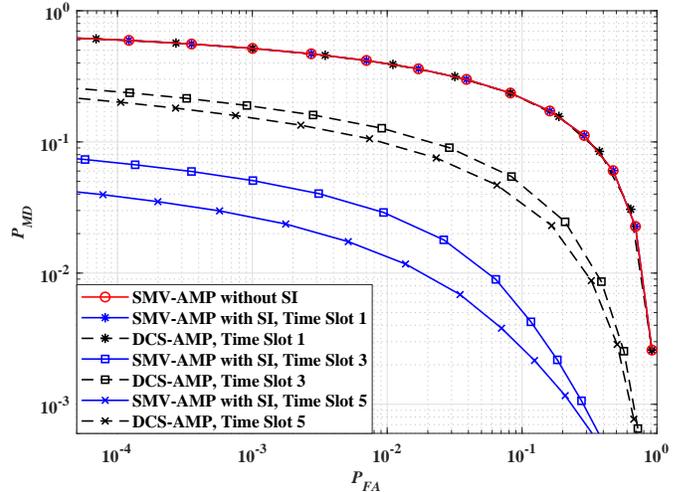


Fig. 3. Activity detection under SMV-AMP algorithm with SI.

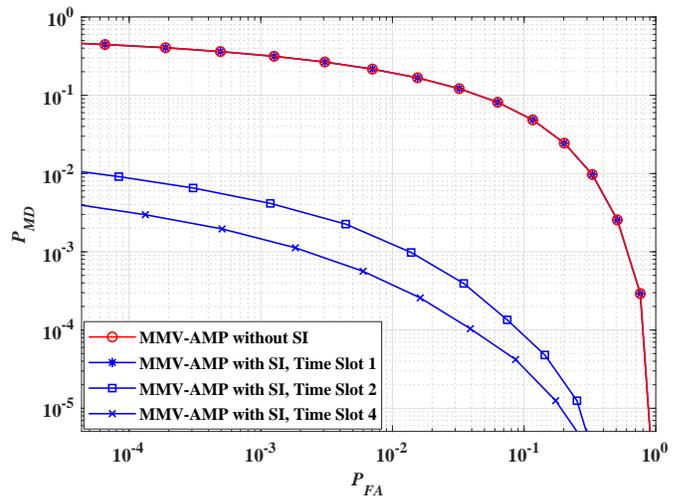


Fig. 4. Activity detection under MMV-AMP algorithm with SI.

correlation between the estimated activity in the last time slot and true activity in the current time slot.

Next, we consider the case when the BS is equipped with multiple antennas. In this case, we term our proposed algorithm as MMV-AMP with SI. Fig. 4 shows the tradeoff between the false alarm probability P_{FA} and missed detection probability P_{MD} when $M = 2$ and $L = 500$. It is observed that the proposed algorithm with SI achieves significant performance gain over the one without using SI; moreover, the performance improves as the time slot index increases. These results in Fig. 3 and Fig. 4 show that by smartly leveraging SI, the proposed algorithm is able to achieve satisfactory detection performance with only a small number of BS antennas (e.g., even 1 or 2), thus being a promising cost-effective solution for future massive IoT systems.

V. CONCLUSION

In this paper, we utilized the temporal correlation in user activity for device activity detection in IoT systems. The main motivation is to achieve satisfactory detection performance with a smaller number of BS antennas and thus lower computational complexity. Along this line, we designed a framework of SI-aided MMV-AMP, where the estimation at the previous time slot was leveraged as SI to devise better MMSE denoisers and activity detectors at the current time slot.

APPENDIX A PROOF OF THEOREM 1

In this proof, for simplicity, we omit the subscripts t and n in all the notations. Then, in the MMSE denoisers (11), the conditional expectation can be given by

$$\begin{aligned} & E[\mathbf{X}^{(j)} \mid \tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)}] \\ \stackrel{(a)}{=} & E[\mathbf{X}^{(j)} \mid \tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)}, \text{Case 1}] p(\text{Case 1} \mid \tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)}) \\ & + E[\mathbf{X}^{(j)} \mid \tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)}, \text{Case 3}] p(\text{Case 3} \mid \tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)}), \end{aligned} \quad (20)$$

where (a) is because $\mathbf{X}^{(j)} = \mathbf{0}$ for *Case 2* and *Case 4* according to Section II. In the following, we focus on *Case 1* and *Case 3* to characterize (20).

Case 1: According to Section II, under *Case 1*, it follows that $\mathbf{x}_n^{(j-1)} = \mathbf{h}_n^{(j-1)}$ and $\mathbf{x}_n^{(j)} = \mathbf{h}_n^{(j)}$. Based on (7), (9), and (12), we have $\tilde{\mathbf{x}}^{(j)} = \mathbf{h}^{(j)} + \tau^{(j)}\mathbf{v}^{(j)}$, $\tilde{\mathbf{x}}_\infty^{(j-1)} = \mathbf{h}^{(j-1)} + \tau_\infty^{(j-1)}\mathbf{v}^{(j-1)}$. Thus, $E[\mathbf{X}^{(j)} \mid \tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)}, \text{Case 1}]$ in (20) can be given by

$$\begin{aligned} E[\mathbf{h}^{(j)} \mid \tilde{\mathbf{x}}^{(j)} = \mathbf{h}^{(j)} + \tau^{(j)}\mathbf{v}^{(j)}, \\ \tilde{\mathbf{x}}_\infty^{(j-1)} = \mathbf{h}^{(j-1)} + \tau_\infty^{(j-1)}\mathbf{v}^{(j-1)}]. \end{aligned} \quad (21)$$

Because $\mathbf{h}^{(j)}$, $\mathbf{h}^{(j-1)}$, $\mathbf{v}^{(j)}$, and $\mathbf{v}^{(j-1)}$ are independent with each other, we have

$$\begin{aligned} & E[\mathbf{h}^{(j)} \mid \tilde{\mathbf{x}}^{(j)} = \mathbf{h}^{(j)} + \tau^{(j)}\mathbf{v}^{(j)}, \\ & \tilde{\mathbf{x}}_\infty^{(j-1)} = \mathbf{h}^{(j-1)} + \tau_\infty^{(j-1)}\mathbf{v}^{(j-1)}] \\ = & E[\mathbf{h}^{(j)} \mid \tilde{\mathbf{x}}^{(j)} = \mathbf{h}^{(j)} + \tau^{(j)}\mathbf{v}^{(j)}] \end{aligned} \quad (22)$$

$$= \gamma \left(\gamma + \left(\tau^{(j)} \right)^2 \right)^{-1} \tilde{\mathbf{x}}^{(j)}. \quad (23)$$

Next, we calculate $p(\text{Case 1} \mid \tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)})$ in (20) as follows

$$\begin{aligned} & p(\text{Case 1} \mid \tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)}) \\ = & \frac{p(\tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)}, \text{Case 1})}{p(\tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)})} \end{aligned} \quad (24)$$

$$= \frac{P(\text{Case 1}) p(\tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)} \mid \text{Case 1})}{p(\tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)})} \quad (25)$$

$$= \frac{\alpha \lambda \psi_{\gamma + (\tau^{(j)})^2}(\tilde{\mathbf{x}}^{(j)}) \psi_{\gamma + (\tau_\infty^{(j-1)})^2}(\tilde{\mathbf{x}}_\infty^{(j-1)})}{p(\tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)})}, \quad (26)$$

where

$$\psi_{\sigma^2}(\mathbf{x}) = \frac{1}{\pi^{|\sigma^2 \mathbf{I}|}} \exp\left(-\frac{\mathbf{x}^H \mathbf{x}}{\sigma^2}\right), \quad (27)$$

is the power density function (pdf) of a complex Gaussian random vector with zero mean and covariance $\sigma^2 \mathbf{I}$. We will derive the joint pdf $p(\tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)})$ in (26) later.

Case 3: Similar to (23) in *Case 1*, under *Case 3*, it can be shown that

$$\begin{aligned} & E[\mathbf{X}^{(j)} \mid \tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)}, \text{Case 3}] \\ = & \gamma \left(\gamma + \left(\tau^{(j)} \right)^2 \right)^{-1} \tilde{\mathbf{x}}^{(j)}. \end{aligned} \quad (28)$$

Moreover, similar to (26), it follows that

$$\begin{aligned} & p(\text{Case 3} \mid \tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)}) \\ = & \frac{p(\tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)}, \text{Case 3})}{p(\tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)})} \\ = & \frac{\beta(1-\lambda) \psi_{\gamma + (\tau^{(j)})^2}(\tilde{\mathbf{x}}^{(j)}) \psi_{(\tau_\infty^{(j-1)})^2}(\tilde{\mathbf{x}}_\infty^{(j-1)})}{p(\tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)})}. \end{aligned} \quad (29)$$

To derive $p(\text{Case 1} \mid \tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)})$ in (26) and $p(\text{Case 3} \mid \tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)})$ in (30), the last step is to characterize $p(\tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)})$. Similar to $p(\tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)}, \text{Case 1})$ shown in (26) and $p(\tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)}, \text{Case 3})$ shown in (30), it can be shown that

$$\begin{aligned} & p(\tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)}, \text{Case 2}) \\ = & (1-\alpha) \lambda \psi_{(\tau^{(j)})^2}(\tilde{\mathbf{x}}^{(j)}) \psi_{\gamma + (\tau_\infty^{(j-1)})^2}(\tilde{\mathbf{x}}_\infty^{(j-1)}), \end{aligned} \quad (31)$$

$$\begin{aligned} & p(\tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)}, \text{Case 4}) \\ = & (1-\beta)(1-\lambda) \psi_{(\tau^{(j)})^2}(\tilde{\mathbf{x}}^{(j)}) \psi_{(\tau_\infty^{(j-1)})^2}(\tilde{\mathbf{x}}_\infty^{(j-1)}). \end{aligned} \quad (32)$$

Thus, it follows that

$$\begin{aligned} & p(\tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)}) \\ = & p(\tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)}, \text{Case 1}) + p(\tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)}, \text{Case 2}) \\ & + p(\tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)}, \text{Case 3}) + p(\tilde{\mathbf{x}}^{(j)}, \tilde{\mathbf{x}}_\infty^{(j-1)}, \text{Case 4}) \end{aligned} \quad (33)$$

$$\begin{aligned} = & \alpha \lambda \psi_{\gamma + (\tau^{(j)})^2}(\tilde{\mathbf{x}}^{(j)}) \psi_{\gamma + (\tau_\infty^{(j-1)})^2}(\tilde{\mathbf{x}}_\infty^{(j-1)}) \\ & + (1-\alpha) \lambda \psi_{(\tau^{(j)})^2}(\tilde{\mathbf{x}}^{(j)}) \psi_{\gamma + (\tau_\infty^{(j-1)})^2}(\tilde{\mathbf{x}}_\infty^{(j-1)}) \\ & + \beta(1-\lambda) \psi_{\gamma + (\tau^{(j)})^2}(\tilde{\mathbf{x}}^{(j)}) \psi_{(\tau_\infty^{(j-1)})^2}(\tilde{\mathbf{x}}_\infty^{(j-1)}) \\ & + (1-\beta)(1-\lambda) \psi_{(\tau^{(j)})^2}(\tilde{\mathbf{x}}^{(j)}) \psi_{(\tau_\infty^{(j-1)})^2}(\tilde{\mathbf{x}}_\infty^{(j-1)}). \end{aligned} \quad (34)$$

By plugging (23), (26), (28), (30), and (34) into (20), it can be shown that the MMSE denoisers by taking SI into account are expressed by (14). Theorem 1 is thus proved.

APPENDIX B
PROOF OF THEOREM2

It can be shown that

$$\begin{aligned} & p(\tilde{\mathbf{x}}_{n,t}^{(j)}, \tilde{\mathbf{x}}_{n,\infty}^{(j-1)} \mid \mathbf{x}_n^{(j)} \neq \mathbf{0}) \\ &= \frac{p(\tilde{\mathbf{x}}_{n,t}^{(j)}, \tilde{\mathbf{x}}_{n,\infty}^{(j-1)}, \text{Case 1}) + p(\tilde{\mathbf{x}}_{n,t}^{(j)}, \tilde{\mathbf{x}}_{n,\infty}^{(j-1)}, \text{Case 3})}{p(\mathbf{x}_n^{(j)} \neq \mathbf{0})} \quad (35) \end{aligned}$$

$$\begin{aligned} &= \frac{p(\tilde{\mathbf{x}}_{n,t}^{(j)}, \tilde{\mathbf{x}}_{n,\infty}^{(j-1)} \mid \text{Case 1})p(\text{Case 1})}{p(\mathbf{x}_n^{(j)} \neq \mathbf{0})} \\ &+ \frac{p(\tilde{\mathbf{x}}_{n,t}^{(j)}, \tilde{\mathbf{x}}_{n,\infty}^{(j-1)} \mid \text{Case 3})p(\text{Case 3})}{p(\mathbf{x}_n^{(j)} \neq \mathbf{0})} \quad (36) \\ &= \psi_{\gamma_n + (\tau_t^{(j)})^2}(\tilde{\mathbf{x}}_{n,t}^{(j)}) \times \\ &\quad \left(\alpha \psi_{\gamma_n + (\tau_\infty^{(j-1)})^2}(\tilde{\mathbf{x}}_{n,\infty}^{(j-1)}) + (1 - \alpha) \psi_{(\tau_\infty^{(j-1)})^2}(\tilde{\mathbf{x}}_{n,\infty}^{(j-1)}) \right), \quad (37) \end{aligned}$$

where (37) is due to (26) and (30). Similarly, it can be shown that

$$\begin{aligned} & p(\tilde{\mathbf{x}}_{n,t}^{(j)}, \tilde{\mathbf{x}}_{n,\infty}^{(j-1)} \mid \mathbf{x}_n^{(j)} = \mathbf{0}) = \psi_{(\tau_t^{(j)})^2}(\tilde{\mathbf{x}}_{n,t}^{(j)}) \times \\ & \left(\beta \psi_{\gamma_n + (\tau_\infty^{(j-1)})^2}(\tilde{\mathbf{x}}_{n,\infty}^{(j-1)}) + (1 - \beta) \psi_{(\tau_\infty^{(j-1)})^2}(\tilde{\mathbf{x}}_{n,\infty}^{(j-1)}) \right). \quad (38) \end{aligned}$$

Therefore, by taking (37) and (38) into (17), it follows that

$$\begin{aligned} LLR_{n,t}^{(j)} &= \log \left(\frac{\psi_{\gamma_n + (\tau_t^{(j)})^2}(\tilde{\mathbf{x}}_{n,t}^{(j)})}{\psi_{(\tau_t^{(j)})^2}(\tilde{\mathbf{x}}_{n,t}^{(j)})} \right) \\ &+ \log \left(\frac{\alpha + (1 - \alpha) \frac{\psi_{(\tau_\infty^{(j-1)})^2}(\tilde{\mathbf{x}}_{n,\infty}^{(j-1)})}{\psi_{\gamma_n + (\tau_\infty^{(j-1)})^2}(\tilde{\mathbf{x}}_{n,\infty}^{(j-1)})}}{\beta + (1 - \beta) \frac{\psi_{(\tau_\infty^{(j-1)})^2}(\tilde{\mathbf{x}}_{n,\infty}^{(j-1)})}{\psi_{\gamma_n + (\tau_\infty^{(j-1)})^2}(\tilde{\mathbf{x}}_{n,\infty}^{(j-1)})}} \right). \quad (39) \end{aligned}$$

With (39), it can be shown that the detection rule (17) is equivalent to (18). Theorem 2 is thus proved.

REFERENCES

- [1] L. Liu, E. G. Larsson, W. Yu, P. Popovski, C. Stefanovic, and E. de Carvalho, "Sparse signal processing for grant-free massive connectivity: A future paradigm for random access protocols in the Internet of Things," *IEEE Signal Process. Mag.*, vol. 35, no. 5, pp. 88-99, Sep. 2018.
- [2] 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.
- [3] Z. Chen, F. Sahrabi, and W. Yu, "Sparse activity detection for massive connectivity," *IEEE Trans. Signal Process.*, vol. 66, no. 7, pp. 1890-1904, Apr. 2018.
- [4] 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.
- [5] 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.
- [6] T. Jiang, Y. Shi, J. Zhang, and K. B. Letaief, "Joint activity detection and channel estimation for IoT networks: Phase transition and computation-estimation tradeoff," *IEEE Internet Things J.*, vol. 6, no. 4, pp. 6212-6225, Aug. 2019.

- [7] 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.
- [8] D. Donoho, A. Maleki, and A. Montanari, "Message passing algorithms for compressed sensing," *Proc. Nat. Acad. Sci.*, vol. 106, no. 45, pp. 18 914-18919, Nov. 2009.
- [9] J. Kim, W. Chang, B. Jung, D. Baron, and J. C. Ye, "Belief propagation for joint sparse recovery," Feb. 2011, [Online] Available: <http://arxiv.org/abs/1102.3289>.
- [10] J. Ziniel and P. Schniter, "Dynamic compressive sensing of time-varying signals via approximate message passing," *IEEE Trans. Signal Process.*, vol. 61, no. 21, pp. 5270-5284, Nov. 2013.
- [11] A. Ma, C. Rush, D. Baron, and D. Needell "An approximate message passing framework for side information," *IEEE Trans. Signal Process.*, vol. 67, no. 7, pp. 1875-1888, Apr. 2019.
- [12] P. Schniter, "Turbo reconstruction of structured sparse signals," in *Proc. Conf. Inf. Sci. Syst. (CISS)*, Mar. 2010, pp. 1-6.