

MIMO-OFDM-Based Massive Connectivity With Frequency Selectivity Compensation

Wenjun Jiang, Mingyang Yue, Xiaojun Yuan, *Senior Member, IEEE*, and Yong Zuo

Abstract—In this paper, we study how to efficiently and reliably detect active devices and estimate their channels in a multiple-input multiple-output (MIMO) orthogonal frequency-division multiplexing (OFDM) based grant-free non-orthogonal multiple access (NOMA) system to enable massive machine-type communications (mMTC). First, by exploiting the correlation of the channel frequency responses in narrow-band mMTC, we propose a block-wise linear channel model. Specifically, the continuous OFDM subcarriers in the narrow-band are divided into several sub-blocks and a linear function with only two variables (mean and slope) is used to approximate the frequency-selective channel in each sub-block. This significantly reduces the number of variables to be determined in channel estimation and the sub-block number can be adjusted to reliably compensate the channel frequency-selectivity. Second, we formulate the joint active device detection and channel estimation in the block-wise linear system as a Bayesian inference problem. By exploiting the block-sparsity of the channel matrix, we develop an efficient turbo message passing (Turbo-MP) algorithm to resolve the Bayesian inference problem with near-linear complexity. We further incorporate machine learning approaches into Turbo-MP to learn unknown prior parameters. Numerical results demonstrate the superior performance of the proposed algorithm over state-of-the-art algorithms.

Index Terms—grant-free non-orthogonal multiple access, orthogonal frequency division multiplexing, turbo message passing.

I. INTRODUCTION

MASSIVE machine-type communications (mMTC) has been envisioned as one of the three key application scenarios of fifth-generation (5G) wireless communications. To support massive connectivity of machine-type devices, mMTC typically has the feature of sporadic transmission with short packets, which is different from conventional human-type communications [1]. This implies that only a small subgroup of devices are active in any time instance of mMTC. As such, in addition to channel estimation and signal detection, a fundamentally new challenge for the design of an mMTC receiver is to reliably and efficiently identify which subgroup of devices are actively engaged in packet transmission.

Recently, a new random access protocol termed grant-free non-orthogonal multiple access (NOMA) has been evaluated and highlighted for mMTC [2]. In specific, in grant-free NOMA, the devices share the same time and frequency resources for signal transmission, and the signals can be transmitted without the scheduling grant from the base station (BS). The receiver at the BS is then required to perform active device detection (ADD), channel estimation (CE), and signal detection (SD), either separately or jointly. The earlier work [3]–[6] assumed that full channel state information (CSI) can

be acquired at the BS and studied joint CE and SD. However, the assumption of full CSI availability is not practical since it will cause a huge overhead to estimate the CSI of all access devices. The follow-up work [7] proposed to divide the process at the BS into two phases, namely, the joint ADD and CE phase and the SD phase. Since the BS only needs to estimate the CSI of active devices, the pilot overhead is significantly reduced. In addition, the channel sparsity in the device domain enables the employment of compressed sensing (CS) algorithms [8] to solve the joint ADD and CE problem. For example, the authors in [9] considered the multiple-input multiple-output (MIMO) transmission and leveraged a multiple measurement vector (MMV) CS technique termed vector approximate message passing (Vector AMP) [10] to achieve asymptotically perfect ADD. In [11], the authors further considered a mixed analog-to-digital converters (ADCs) architecture at the BS antennas and proposed a CS algorithm based on the turbo compressive sensing (Turbo-CS) [12]. It is known from [12] that Turbo-CS outperforms approximate message passing (AMP) [13] both in convergence performance and computational complexity. Another line of research considered the more challenging joint ADD, CE, and SD problem [14], [15]. As compared to the two-phase approach, these joint schemes can achieve significant performance improvement but at the expense of higher computational complexity.

Orthogonal frequency division multiplexing (OFDM) is a mature and enabling technology for 5G to provide high spectral efficiency. As such, the design of OFDM-based grant-free NOMA has attracted much research interest in recent years [16]–[18]. In [16], the authors exploited the block-sparsity of the channel responses on OFDM subcarriers to design a message-passing-based iterative algorithm. Besides, it has been demonstrated in [17] that the message-passing-based iterative algorithm can be unfolded into a deep neural network. By training the parameters of the neural network, the convergence and performance of the algorithm are improved. Furthermore, OFDM-based grant-free NOMA with massive MIMO has been considered in [18]. By leveraging the sparsity both in the device domain and the virtual angular domain, the authors utilized the AMP algorithm to achieve the joint ADD and CE. One issue in [16]–[18] is that the frequency-domain channel estimation on every subcarrier requires a high pilot overhead, which is inefficient for the short data packets transmission in mMTC.

To reduce the pilot overhead, a common strategy is to transform the frequency-domain channel into the time-domain channel by inverse discrete Fourier transform (IDFT) [19]. Due to the limited delay spread in the time-domain channel, the

time-domain channel is sparse, thereby requiring fewer pilots for CE. Furthermore, by exploiting the sparsity of both the time-domain channel and the device activity pattern, some state-of-the-art CS algorithms such as Turbo-CS [12], [20] and Vector AMP [9], [10] can be applied to the considered systems with some straightforward modifications. However, there exists an energy leakage problem caused by the IDFT transform to obtain the time-domain channel. The energy leakage compromises the channel sparsity in the time domain. In addition, the power delay profile (PDP) is generally difficult for the BS to acquire, and thus cannot be exploited as prior information to improve the system performance.

Motivated by the bottleneck of the existing channel models when applied to the MIMO-OFDM-based grant-free NOMA system, we aim to construct a channel model to enable efficient massive random access. Due to the short packets in mMTC, the bandwidth for packet transmission is usually narrow, e.g., 1MHz for 10^5 access devices [21]. Then the variations of the channel frequency responses across the subcarriers are limited and slow. By leveraging this fact, we propose a block-wise linear channel model. Specifically, the continuous subcarriers are divided into several sub-blocks. In each sub-block, the frequency-selective channel is approximated by a linear function. We demonstrate that the number of variables in the block-wise linear channel model is typically much less than the number of non-zero delay taps in the time-domain channel. Moreover, the sub-block number can be modified to strike the trade-off between the model accuracy and the number of the channel variables to be estimated.

Based on the block-wise linear system model, we aim to design a CS algorithm to solve the joint ADD and CE problem. We first introduce a probability model to characterize the block-sparsity of the channel matrix. Then the joint ADD and CE is formulated as a Bayesian inference problem. Inspired by the success of Turbo-CS [12], [20] in sparse signal recovery, we design a message passing algorithm termed turbo message passing (Turbo-MP) to resolve the Bayesian inference problem. The message passing processes are derived based on the principle of Turbo-CS and the sum-product rule [22]. By designing a partial orthogonal pilot matrix with fast transform, Turbo-MP achieves near-linear complexity.

Furthermore, we show that machine learning methods can be incorporated into Turbo-MP to learn the unknown prior parameters. We adopt the expectation maximization (EM) algorithm [23] to learn the prior parameters. We then show how to unfold Turbo-MP into a neural network (NN), where the prior parameters are seen as the learnable parameters of the neural network. Numerical results show that NN-based Turbo-MP has a faster convergence rate than EM-based Turbo-MP. More importantly, we show that Turbo-MP designed for the propose frequency-domain block-wise linear model significantly outperforms the state-of-the-art counterparts, especially those message passing algorithms designed for the time-domain sparse channel model [9], [11].

A. Notation and Organization

We use bold capital letters like \mathbf{X} for matrices and bold lowercase letters \mathbf{x} for vectors. $(\cdot)^T$ and $(\cdot)^H$ are used to

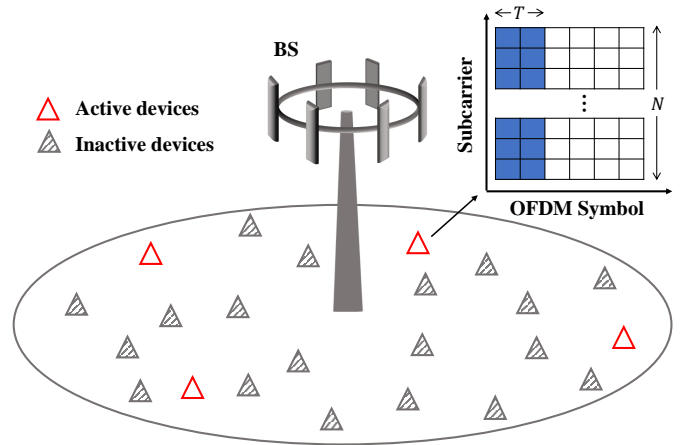


Fig. 1. An illustration of MIMO-OFDM-based mMTC, where a small subgroup of devices are active in each time instance and share N continuous subcarriers for pilot transmission within T OFDM symbols.

denote the transpose and the conjugate transpose, respectively. We use $\text{diag}(\mathbf{x})$ for the diagonal matrix created from vector \mathbf{x} , $\text{diag}(\mathbf{x}_1, \dots, \mathbf{x}_N)$ for the block diagonal matrix with the n -th block being vector \mathbf{x}_n , and $\text{diag}(\mathbf{X}_1, \dots, \mathbf{X}_N)$ for the block diagonal matrix with the n -th block being matrix \mathbf{X}_n . We use $\text{vec}(\mathbf{X})$ for the vectorization of matrix \mathbf{X} and \otimes for the Kronecker product. $\|\mathbf{X}\|_F$ and $\|\mathbf{x}\|_2$ are used to denote the Frobenius norm of matrix \mathbf{X} and the l_2 norm of vector \mathbf{x} , respectively. Matrix \mathbf{I} denotes the identity matrix with an appropriate size. For a random vector \mathbf{x} , we denote its probability density function (pdf) by $p(\mathbf{x})$. $\delta(\cdot)$ denotes the Dirac delta function and $\delta[\cdot]$ denotes the Kronecker delta function. The pdf of a complex Gaussian random vector $\mathbf{x} \in \mathbb{C}^N$ with mean \mathbf{m} and covariance \mathbf{C} is denoted by $\mathcal{CN}(\mathbf{x}; \mathbf{m}, \mathbf{C}) = \exp(-(\mathbf{x} - \mathbf{m})^H \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m})) / (\pi^N |\mathbf{C}|)$.

The remainder of this paper is organized as follows. In Section II, we introduce the existing system models. Furthermore, we propose the block-wise linear system model and demonstrate its superiority. In Section III, we formulate a Bayesian inference problem to address the ADD and CE problem. In Section IV, we propose the Turbo-MP algorithm, describe the pilot design, and analyze the algorithm complexity. In Section V, we apply the EM algorithm to learn the prior parameters. Moreover, we show how Turbo-MP is unfolded into a neural network. In Section VI, we present the numerical results. In Section VII, we conclude this paper.

II. SYSTEM MODELING

A. MIMO-OFDM-Based Grant-free NOMA Model

Consider a MIMO-OFDM-based grant-free NOMA system as shown in Fig. 1, where a frequency band of N adjacent OFDM subcarriers are allocated for mMTC. N is usually much less than the total number of available subcarriers in the considered OFDM system. This frequency allocation strategy follows the idea of narrow-band internet-of-things [24] and network slicing [25]. Then the allocated bandwidth is used to support K single-antenna devices to randomly access an M -antenna base station (BS). In each time instance, only a small

subset of devices are active. To characterize such sporadic transmission, the device activity is described by an indicator function α_k as

$$\alpha_k = \begin{cases} 1, & \text{device } k \text{ is active} \\ 0, & \text{device } k \text{ is inactive,} \end{cases} \quad k = 1, \dots, K \quad (1)$$

with $p(\alpha_k = 1) = \lambda$ where $\lambda \ll 1$.

We adopt a multipath block-fading channel, i.e., the multipath channel response remain constant within the coherence time. Denote the channel frequency response on the n -th subcarrier from the k -th device at m -th BS antenna by

$$g_{k,m,n} = \sum_{l=1}^{L_k} \sqrt{\rho_{k,l}} \beta_{k,m,l} e^{-j2\pi\Delta f \tau_{k,l} n}, \quad k = 1, \dots, K; \quad m = 1, \dots, M; \quad n = 1, \dots, N \quad (2)$$

where Δf is the OFDM subcarrier spacing; L_k is the number of channel taps of the k -th device; $\rho_{k,l}$ and $\tau_{k,l}$ are respectively the l -th tap power and tap delay of the k -th device; $\beta_{k,m,l} \sim \mathcal{CN}(\beta_{k,m,l}; 0, 1)$ is the normalized complex gain and assumed to be independent for any k, m, l [26]. Then the channel frequency response can be expressed in a matrix form as

$$\mathbf{G}_k = \alpha_k \begin{pmatrix} g_{k,1,1} & \cdots & g_{k,m,1} & \cdots & g_{k,M,1} \\ \vdots & & \vdots & & \vdots \\ g_{k,1,N} & \cdots & g_{k,m,N} & \cdots & g_{k,M,N} \end{pmatrix} \quad (3)$$

Let $a_{k,n}^{(t)}$ be the pilot symbol of the k -th device transmitted on the n -th subcarrier at the t -th OFDM symbol, and T be the number of OFDM symbols for pilot transmission. Then we construct a block diagonal matrix $\mathbf{\Lambda}_k \in \mathbb{C}^{TN \times TN}$ with the n -th diagonal block being $[a_{k,n}^{(1)}, \dots, a_{k,n}^{(T)}]^T$, i.e.,

$$\mathbf{\Lambda}_k = \text{diag}([a_{k,1}^{(1)}, \dots, a_{k,1}^{(T)}]^T, \dots, [a_{k,N}^{(1)}, \dots, a_{k,N}^{(T)}]^T) \quad (4)$$

Assume the cyclic-prefix (CP) length $L_{cp} > \tau_{k,l}, \forall k, l$. After removing the CP and applying the discrete Fourier transform (DFT), the system model in the frequency domain is described as

$$\mathbf{Y} = \sum_{k=1}^K \mathbf{\Lambda}_k \mathbf{G}_k + \mathbf{N} \quad (5)$$

where $\mathbf{Y} \in \mathbb{C}^{TN \times M}$ is the observation matrix; \mathbf{N} is an additive white Gaussian noise (AWGN) matrix with its elements independently drawn from $\mathcal{CN}(0, \sigma_N^2)$. In [16]–[18], CS algorithms were proposed based on (5) to achieve the CE on every OFDM subcarrier.

Define the time-domain channel matrix of the k -th device as $\tilde{\mathbf{H}}_k \in \mathbb{C}^{N \times M}$. Note that $\tilde{\mathbf{H}}_k$ can be represented as the IDFT of \mathbf{G}_k , i.e.,

$$\tilde{\mathbf{H}}_k = \mathbf{F}^H \mathbf{G}_k \quad (6)$$

where $\mathbf{F} \in \mathbb{C}^{N \times N}$ is the DFT matrix with the (n_1, n_2) -th element being $1/\sqrt{N} \cdot \exp(-j2\pi n_1 n_2 / N)$. Similarly, \mathbf{G}_k is

¹Strictly speaking, the differences of the tap delay at different BS antennas are absorbed in $\beta_{k,m,l}$

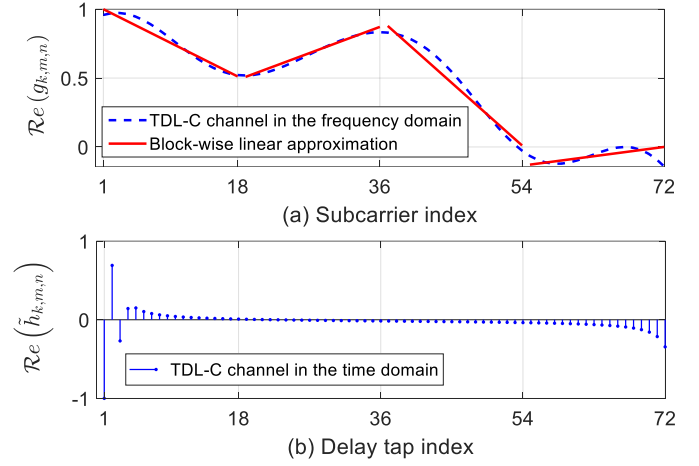


Fig. 2. An example of channel model comparison. TDL-C multi-path channel in TR 38.901 R14 [27]. The number of OFDM subcarriers $N = 72$. (a) frequency-domain channel response and its block-wise linear approximation with sub-block number $Q = 4$. (b) corresponding time-domain channel response where the increase of the channel response at the tail is caused by the energy leakage of the IDFT transform.

represented as $\mathbf{G}_k = \mathbf{F} \tilde{\mathbf{H}}_k$. Substituting $\mathbf{G}_k = \mathbf{F} \tilde{\mathbf{H}}_k$ into (5), we obtain

$$\begin{aligned} \mathbf{Y} &= \sum_{k=1}^K \mathbf{\Lambda}_k \mathbf{F} \tilde{\mathbf{H}}_k + \mathbf{N} \\ &= [\mathbf{\Lambda}_1 \mathbf{F}, \dots, \mathbf{\Lambda}_K \mathbf{F}] \tilde{\mathbf{H}} + \mathbf{N} \end{aligned} \quad (7)$$

where $\tilde{\mathbf{H}} = [\tilde{\mathbf{H}}_1^T, \dots, \tilde{\mathbf{H}}_K^T]^T$ is the channel matrix in the time domain. It is known that the channel delay spread is usually much smaller than N , i.e., some rows of $\tilde{\mathbf{H}}_k$ are zeros. Besides, due to the sporadic transmission of the devices, $\tilde{\mathbf{H}}_k$ is an all-zero matrix if device k is inactive. In this case, CS algorithms such as Vector AMP [9], [10] and Turbo-CS [12], [20] can be used to recover $\tilde{\mathbf{H}}$ from observation \mathbf{Y} by exploiting the sparsity of $\tilde{\mathbf{H}}$. However, path delay is generally not an integer multiple of the sampling interval, resulting in the energy leakage problem of the IDFT transform (6) which severely compromises the sparsity of $\tilde{\mathbf{H}}$. An illustration of the energy leakage problem is given in Fig. 2(b).

B. Proposed Block-Wise Linear System Model

For narrow-band mMTC [1], the variations of the channel frequency response across the subcarriers are typically slow and limited, which inspires us to develop an alternative channel model to efficiently leverage the correlation in the frequency domain. In specific, we propose a block-wise linear channel model. We divide the N continuous subcarriers into Q sub-blocks. In each sub-block q , a linear function is used as the approximation of the channel frequency response:

$$g_{k,n_q,m} = h_{k,q,m} + (n_q - l_q) c_{k,q,m} + \Delta_{k,n_q,m} \quad (8)$$

$$n_q = (q-1)N/Q + 1, \dots, qN/Q$$

where $h_{k,q,m}$ and $c_{k,q,m}$ represent the mean and slope of the linear function in the q -th sub-block, respectively; $\Delta_{k,n_q,m}$ is the error term due to model mismatch; and $l_q = (q - \frac{1}{2})N/Q$ is

the midpoint of n_q . Intuitively, $h_{k,q,m}$ can be seen as the mean-value of the channel response in the q -th sub-block, and $c_{k,q,m}$ is used to characterize the change of the channel response for the compensation of the frequency-selectivity. For the k -th device, we define the matrix $\mathbf{H}_k \in \mathbb{C}^{Q \times M}$ and $\mathbf{C}_k \in \mathbb{C}^{Q \times M}$ as

$$\mathbf{H}_k = \alpha_k \begin{pmatrix} h_{k,1,1} & \cdots & h_{k,1,m} & \cdots & h_{k,1,M} \\ \vdots & & \vdots & & \vdots \\ h_{k,Q,1} & \cdots & h_{k,Q,m} & \cdots & h_{k,Q,M} \end{pmatrix} \quad (9)$$

$$\mathbf{C}_k = \alpha_k \begin{pmatrix} c_{k,1,1} & \cdots & c_{k,1,m} & \cdots & c_{k,1,M} \\ \vdots & & \vdots & & \vdots \\ c_{k,Q,1} & \cdots & c_{k,Q,m} & \cdots & c_{k,Q,M} \end{pmatrix}. \quad (10)$$

The reason for introducing α_k in (9)-(10) is that the channel estimation and device activity detection can be jointly achieved by recovering \mathbf{H}_k and \mathbf{C}_k .

Define $\mathbf{E}_1 = \text{diag}(\mathbf{1}_{N/Q}, \dots, \mathbf{1}_{N/Q}) \in \mathbb{R}^{N \times Q}$ with $\mathbf{1}_{N/Q}$ being an all-one vector of length N/Q and $\mathbf{E}_2 = \text{diag}(\mathbf{d}, \dots, \mathbf{d}) \in \mathbb{R}^{N \times Q}$ with $\mathbf{d} = [-\frac{N}{2Q} + 1, \dots, \frac{N}{2Q}]^T$. Then the block-wise linear approximation of the channel frequency response \mathbf{G}_k is given by

$$\mathbf{G}_k = \mathbf{E}_1 \mathbf{H}_k + \mathbf{E}_2 \mathbf{C}_k + \mathbf{\Delta}_k \quad (11)$$

where $\mathbf{\Delta}_k \in \mathbb{C}^{N \times M}$ is the error matrix from the k -th device with the (n, m) -th element $\Delta_{k,n,m}$ in (8). Substituting (11) into (5) with some manipulations, we obtain the block-wise linear system model as

$$\mathbf{Y} = \mathbf{A}\mathbf{H} + \mathbf{B}\mathbf{C} + \mathbf{W} \quad (12)$$

where $\mathbf{A} = [\mathbf{\Lambda}_1 \mathbf{E}_1, \dots, \mathbf{\Lambda}_K \mathbf{E}_1] \in \mathbb{C}^{TN \times QK}$ and $\mathbf{B} = [\mathbf{\Lambda}_1 \mathbf{E}_2, \dots, \mathbf{\Lambda}_K \mathbf{E}_2] \in \mathbb{C}^{TN \times QK}$ are the pilot matrices; $\mathbf{H} = [\mathbf{H}_1^T, \dots, \mathbf{H}_K^T]^T \in \mathbb{C}^{QK \times M}$ is the *channel mean matrix*; $\mathbf{C} = [\mathbf{C}_1^T, \dots, \mathbf{C}_K^T]^T \in \mathbb{C}^{QK \times M}$ is the *channel compensation matrix*; \mathbf{W} is the summation of the AWGN and the error terms from model mismatch as

$$\mathbf{W} = \sum_{k=1}^K \mathbf{\Lambda}_k \mathbf{\Delta}_k + \mathbf{N}. \quad (13)$$

Following the central limit theorem, for a large devices number K , \mathbf{W} can be modeled as an AWGN matrix with mean zero and variance σ_w^2 . We note that with $Q = N$ and $\mathbf{C} = \mathbf{0}$, system model (12) reduces to the model (5) in [16]–[18] where the channel response on every subcarrier needs to be estimated exactly. Furthermore, we adjust the sub-block number Q to strike a balance between the number of channel variables to be estimated and the model accuracy. An example is shown in Fig. 2 where the channel response in the frequency domain and its block-wise linear approximation is given in Fig. 2(a). It is seen that sub-block number $Q = 4$ is sufficient to ensure that the block-wise linear model approximates the frequency-domain channel response very well. This implies that, with model (12), only $2Q$ variables need to be estimated for each device at each BS antenna. Compared to the channel response in the time domain as shown in Fig. 2(b), it is clear that the number of unknown variables $2Q$ is much less than the number of corresponding non-zero taps. In the remainder of the paper,

we focus on the estimation algorithm design based on model (12). In the simulation section, we will show that our algorithm designed based on (12) significantly outperforms the existing approaches based on (5) and (7).

III. PROBLEM STATEMENT

With model (12), our goal is to recover the channel mean matrix \mathbf{H} and channel compensation matrix \mathbf{C} based on the noisy observation \mathbf{Y} . This task can be constructed as a Bayesian inference problem. In the following, we first introduce the probability model of \mathbf{H} and \mathbf{C} , and then describe the statistical inference problem.

Due to the sporadic transmission of the devices, matrices \mathbf{H} and \mathbf{C} have a structured sparsity referred to as block-sparsity. In specific, if the k -th device is inactive, we have $\mathbf{H}_k = \mathbf{0}$ and $\mathbf{C}_k = \mathbf{0}$. With some abuse of notation, we utilize a conditional Bernoulli-Gaussian (BG) distribution [19] to characterize the block-sparsity as

$$\begin{aligned} p(\mathbf{H}_k | \alpha_k) &\sim \delta[\alpha_k] \delta(\mathbf{H}_k) + \delta[1 - \alpha_k] \mathcal{CN}(\mathbf{H}_k; \mathbf{0}; \vartheta_{\mathbf{H}} \mathbf{I}) \\ p(\mathbf{C}_k | \alpha_k) &\sim \delta[\alpha_k] \delta(\mathbf{C}_k) + \delta[1 - \alpha_k] \mathcal{CN}(\mathbf{C}_k; \mathbf{0}; \vartheta_{\mathbf{C}} \mathbf{I}) \end{aligned} \quad (14)$$

where $\delta(\cdot)$ is the Dirac delta function and $\delta[\cdot]$ is the Kronecker delta function. With indicator function $\alpha_k = 0$, \mathbf{H}_k and \mathbf{C}_k are both zeros. With $\alpha_k = 1$, the elements of \mathbf{H}_k and \mathbf{C}_k are independent and identically distributed (i.i.d.) Gaussian with variances $\vartheta_{\mathbf{H}}$ and $\vartheta_{\mathbf{C}}$, respectively. We further assume that each device accesses the BS in an i.i.d. manner. Then the indicator function α_k is drawn from the Bernoulli distribution as

$$p(\alpha_k) = (1 - \lambda) \delta[\alpha_k] + \lambda \delta[1 - \alpha_k] \quad (15)$$

where λ is the device activity rate.

Consider an estimator to minimize the mean-square error (MSE) of \mathbf{H} and \mathbf{C} . It is known that the estimator which minimizes the MSE is the posterior expectation with respect to the posterior distribution [28]. Define $\mathbf{h}_m \in \mathbb{C}^{QK}$ and $\mathbf{c}_m \in \mathbb{C}^{QK}$ as the m -th column of \mathbf{H} and \mathbf{C} , respectively. Define $\mathbf{y}_m \in \mathbb{C}^{TN}$ as the m -th column of \mathbf{Y} . Then the posterior distribution $p(\mathbf{H}, \mathbf{C}, \boldsymbol{\alpha} | \mathbf{Y})$ is described as

$$\begin{aligned} p(\mathbf{H}, \mathbf{C}, \boldsymbol{\alpha} | \mathbf{Y}) &\propto p(\mathbf{Y} | \mathbf{H}, \mathbf{C}) p(\mathbf{H}, \mathbf{C}, \boldsymbol{\alpha}) \\ &\propto \prod_m p(\mathbf{y}_m | \mathbf{h}_m, \mathbf{c}_m) \prod_k p(\mathbf{H}_k | \alpha_k) p(\mathbf{C}_k | \alpha_k) p(\alpha_k) \end{aligned} \quad (16)$$

where $\prod_m p(\mathbf{y}_m | \mathbf{h}_m, \mathbf{c}_m)$ and $\prod_k p(\mathbf{H}_k | \alpha_k) p(\mathbf{C}_k | \alpha_k) p(\alpha_k)$ are the likelihood and the prior, respectively; vector $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_K]^T$. In mMTC with a large device number K , it is computationally intractable to obtain the minimum mean-square error (MMSE) estimator. In the following section, we propose a low-complexity algorithm termed turbo message passing (Turbo-MP) to obtain an approximate MMSE solution.

IV. TURBO MESSAGE PASSING

A. Algorithm Framework

The factor graph representation of $p(\mathbf{H}, \mathbf{C}, \boldsymbol{\alpha} | \mathbf{Y})$ is shown in Fig. 3, based on which Turbo-MP is established. In the factor graph, the likelihood and prior as in (16) are treated as the factor nodes (grey rectangles), while the random variables are

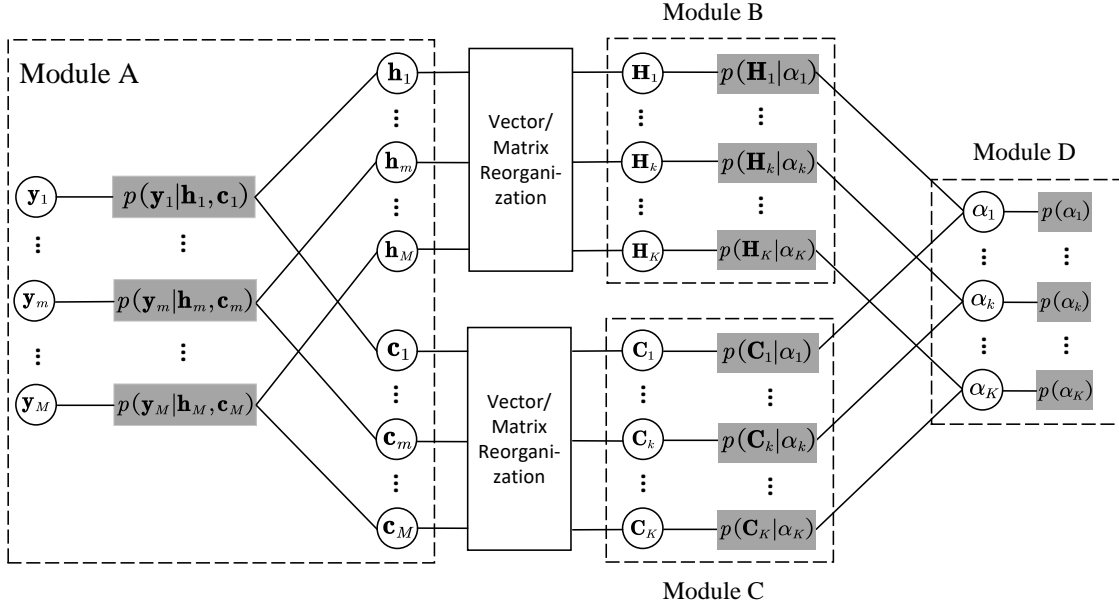


Fig. 3. The block diagram of the proposed turbo message passing (Turbo-MP) algorithm.

treated as the variable nodes (blank circles). In specific, Turbo-MP consists of four modules referred to as Module A, B, C, and D, respectively. Module A is to obtain the estimates of the channel mean matrix \mathbf{H} and the channel compensation matrix \mathbf{C} by exploiting the likelihood $\prod_m p(\mathbf{y}_m | \mathbf{h}_m, \mathbf{c}_m)$. Modules B and C are to obtain the estimates of \mathbf{H} and \mathbf{C} by exploiting the priors $\prod_k p(\mathbf{H}_k | \alpha_k)$ and $p(\mathbf{C}_k | \alpha_k)$, respectively. Module D is to obtain the estimate of α_k by exploiting the prior $p(\alpha_k)$. The estimates are passed between modules in each Turbo-MP iteration like turbo decoding [29]. For example, the output of Module B is used as the input of Module A, and vice versa. The four modules are executed iteratively until convergence. The derivation of Turbo-MP algorithm follows the sum-product rule [22] and the principle of Turbo-CS [12], [20].

B. Module A: Linear Estimation of \mathbf{h}_m and \mathbf{c}_m

In Module A, \mathbf{h}_m and \mathbf{c}_m at antenna $m = 1, \dots, M$ are estimated separately by exploiting the likelihood $\prod_m p(\mathbf{y}_m | \mathbf{h}_m, \mathbf{c}_m)$ and the messages from Modules B and C. In specific, denote the message from variable node \mathbf{h}_m to factor node $p(\mathbf{y}_m | \mathbf{h}_m, \mathbf{c}_m)$ by $\mathcal{M}_{\mathbf{h}_m \rightarrow \mathbf{y}_m}(\mathbf{h}_m)$ and the message from variable node \mathbf{c}_m to factor node $p(\mathbf{y}_m | \mathbf{h}_m, \mathbf{c}_m)$ by $\mathcal{M}_{\mathbf{c}_m \rightarrow \mathbf{y}_m}(\mathbf{c}_m)$. Following the principle of Turbo-CS [12], [20], we assume that the above two messages are Gaussian, i.e., $\mathcal{M}_{\mathbf{h}_m \rightarrow \mathbf{y}_m}(\mathbf{h}_m) \sim \mathcal{CN}(\mathbf{h}_m; \mathbf{h}_m^{A,pri}, v_{\mathbf{h}_m}^{A,pri} \mathbf{I})$ and $\mathcal{M}_{\mathbf{c}_m \rightarrow \mathbf{y}_m}(\mathbf{c}_m) \sim \mathcal{CN}(\mathbf{c}_m; \mathbf{c}_m^{A,pri}, v_{\mathbf{c}_m}^{A,pri} \mathbf{I})$. Note that the algorithm starts with $\mathbf{h}_m^{A,pri} = \mathbf{0}$, $\mathbf{c}_m^{A,pri} = \mathbf{0}$, $v_{\mathbf{h}_m}^{A,pri} = \lambda \vartheta_{\mathbf{H}}$, and $v_{\mathbf{c}_m}^{A,pri} = \lambda \vartheta_{\mathbf{C}}$. From the sum-product rule, the belief of $h_{k,q,m}$ at factor node $p(\mathbf{y}_m | \mathbf{h}_m, \mathbf{c}_m)$ is

$$\mathcal{M}_{\mathbf{y}_m}(h_{k,q,m}) = \int_{/h_{k,q,m}} p(\mathbf{y}_m | \mathbf{h}_m, \mathbf{c}_m) \mathcal{M}_{\mathbf{h}_m \rightarrow \mathbf{y}_m}(\mathbf{h}_m) \times \mathcal{M}_{\mathbf{c}_m \rightarrow \mathbf{y}_m}(\mathbf{c}_m). \quad (17)$$

In the above, \mathbf{y}_m is the subscript is the shorthand for factor node $p(\mathbf{y}_m | \mathbf{h}_m, \mathbf{c}_m)$; $/h_{k,q,m}$ denotes the set includes all elements in \mathbf{h}_m except $h_{k,q,m}$. Then we define the belief of \mathbf{h}_m at factor node $p(\mathbf{y}_m | \mathbf{h}_m, \mathbf{c}_m)$ as

$$\mathcal{M}_{\mathbf{y}_m}(\mathbf{h}_m) = \prod_{k,q} \mathcal{M}_{\mathbf{y}_m}(h_{k,q,m}). \quad (18)$$

Combing (17) and (18), the message $\mathcal{M}_{\mathbf{y}_m}(\mathbf{h}_m)$ is also Gaussian with the mean $\mathbf{h}_m^{A,post}$ and variance $v_{\mathbf{h}_m}^{A,pri}$. From [20, Eq. 11], the mean $\mathbf{h}_m^{A,post}$ is expressed as

$$\mathbf{h}_m^{A,post} = \mathbf{h}_m^{A,pri} + v_{\mathbf{h}_m}^{A,pri} \mathbf{A}^H \Sigma_m^{-1} \times (\mathbf{y}_m - \mathbf{A} \mathbf{h}_m^{A,pri} - \mathbf{B} \mathbf{c}_m^{A,pri}) \quad (19)$$

where Σ_m is the covariance matrix given by

$$\Sigma_m = v_{\mathbf{h}_m}^{A,pri} \mathbf{A} \mathbf{A}^H + v_{\mathbf{c}_m}^{A,pri} \mathbf{B} \mathbf{B}^H + \sigma_w^2 \mathbf{I}. \quad (20)$$

To reduce the computational complexity of channel inverse Σ_m^{-1} , we require that \mathbf{A} is partial orthogonal [12], i.e., $\mathbf{A} \mathbf{A}^H = K P \mathbf{I}$, where P is the average power of the pilot symbol. (The design of \mathbf{A} to guarantee partial orthogonality is presented in IV-F.) In addition, it is easy to verify $\mathbf{B} = \mathbf{D} \mathbf{A}$ where the diagonal matrix $\mathbf{D} = \text{diag}([\mathbf{d}^T, \dots, \mathbf{d}^T]^T \otimes (\mathbf{1}_T)^T) \in \mathbb{R}^{TN \times TN}$ with $\mathbf{1}_T$ being an all-one vector of length T . Then the expression of Σ_m is simplified to

$$\Sigma_m = K P v_{\mathbf{h}_m}^{A,pri} \mathbf{I} + K P v_{\mathbf{c}_m}^{A,pri} \mathbf{D} \mathbf{D}^H + \sigma_w^2 \mathbf{I}. \quad (21)$$

Then the variance $v_{\mathbf{h}_m}^{A,post}$ is given by

$$v_{\mathbf{h}_m}^{A,post} = v_{\mathbf{h}_m}^{A,pri} - \sum_{i=1}^{TN} \frac{P(v_{\mathbf{h}_m}^{A,pri})^2}{\Sigma_{m,i,i}} \quad (22)$$

where $\Sigma_{m,i,i}$ is the (i, i) -th element of Σ_m . We note that (19) and (22) also corresponds to the linear minimum mean-square error (LMMSE) estimator [28, Chap. 11] given the observation \mathbf{y}_m , the mean $\mathbf{h}_m^{A,pri}$, $\mathbf{c}_m^{A,pri}$ and the variance $v_{\mathbf{h}_m}^{A,pri}$, $v_{\mathbf{c}_m}^{A,pri}$.

From the sum-produce rule, the extrinsic message from factor node $p(\mathbf{y}_m|\mathbf{h}_m\mathbf{c}_m)$ to variable node \mathbf{h}_m is calculated as

$$\mathcal{M}_{\mathbf{y}_m \rightarrow \mathbf{h}_m}(\mathbf{h}_m) \propto \frac{\mathcal{M}_{\mathbf{y}_m}(\mathbf{h}_m)}{\mathcal{M}_{\mathbf{h}_m \rightarrow \mathbf{y}_m}(\mathbf{h}_m)}. \quad (23)$$

Given the Gaussian messages $\mathcal{M}(\mathbf{h}_m)$ and $\mathcal{M}_{\mathbf{h}_m \rightarrow \mathbf{y}_m}(\mathbf{h}_m)$, the extrinsic message is also Gaussian as

$$\mathcal{M}_{\mathbf{y}_m \rightarrow \mathbf{h}_m}(\mathbf{h}_m) \sim \mathcal{CN}(\mathbf{h}_m; \mathbf{h}_m^{A,ext}, v_{\mathbf{h}_m}^{A,ext} \mathbf{I}) \quad (24)$$

where the variance $v_{\mathbf{h}_m}^{A,ext}$ and the mean $\mathbf{h}_m^{A,ext}$ are respectively given by

$$\begin{aligned} v_{\mathbf{h}_m}^{A,ext} &= \left(\frac{1}{v_{\mathbf{h}_m}^{A,post}} - \frac{1}{v_{\mathbf{h}_m}^{A,pri}} \right)^{-1} \\ \mathbf{h}_m^{A,ext} &= v_{\mathbf{h}_m}^{A,ext} \left(\frac{\mathbf{h}_m^{A,post}}{v_{\mathbf{h}_m}^{A,post}} - \frac{\mathbf{h}_m^{A,pri}}{v_{\mathbf{h}_m}^{A,pri}} \right). \end{aligned} \quad (25)$$

The calculation to obtain the belief of \mathbf{c}_m is similar. We have the Gaussian belief $\mathcal{M}(\mathbf{c}_m)$ with its mean and variance respectively given by

$$\begin{aligned} \mathbf{c}_m^{A,post} &= \mathbf{c}_m^{A,pri} + v_{\mathbf{c}_m}^{A,pri} \mathbf{B}^H \Sigma_m^{-1} \\ &\quad \times (\mathbf{y}_m - \mathbf{A} \mathbf{h}_m^{A,pri} - \mathbf{B} \mathbf{c}_m^{A,pri}) \\ v_{\mathbf{c}_m}^{A,post} &= v_{\mathbf{c}_m}^{A,pri} - \sum_{i=1}^{TN} \frac{PD_{i,i}^2 (v_{\mathbf{c}_m}^{A,pri})^2}{\Sigma_{m,i,i}} \end{aligned} \quad (26)$$

where $D_{i,i}$ is the (i, i) -th element of \mathbf{D} . Then we obtain the extrinsic message $\mathcal{M}_{\mathbf{y}_m \rightarrow \mathbf{c}_m}(\mathbf{c}_m) \sim \mathcal{CN}(\mathbf{c}_m; \mathbf{c}_m^{A,ext}, v_{\mathbf{c}_m}^{A,ext} \mathbf{I})$ with its mean and variance as follows:

$$\begin{aligned} v_{\mathbf{c}_m}^{A,ext} &= \left(\frac{1}{v_{\mathbf{c}_m}^{A,post}} - \frac{1}{v_{\mathbf{c}_m}^{A,pri}} \right)^{-1} \\ \mathbf{c}_m^{A,ext} &= v_{\mathbf{c}_m}^{A,ext} \left(\frac{\mathbf{c}_m^{A,post}}{v_{\mathbf{c}_m}^{A,post}} - \frac{\mathbf{c}_m^{A,pri}}{v_{\mathbf{c}_m}^{A,pri}} \right). \end{aligned} \quad (27)$$

C. Module B: Denoiser of \mathbf{H}_k

In Module B, each $\mathbf{H}_k, \forall k$ is estimated individually by exploiting the prior $p(\mathbf{H}_k|\alpha_k)$ and the messages from Modules A and D. In specific, given the message $\mathcal{M}_{\mathbf{y}_m \rightarrow \mathbf{h}_m}(\mathbf{h}_m) \sim \mathcal{CN}(\mathbf{h}_m; \mathbf{h}_m^{A,ext}, v_{\mathbf{h}_m}^{A,ext} \mathbf{I})$ from module A, we have

$$\mathcal{M}_{\mathbf{y}_m \rightarrow h_{k,q,m}}(h_{k,q,m}) \sim \mathcal{CN}(h_{k,q,m}; \mathbf{h}_{k,q,m}^{ext}, v_{\mathbf{h}_m}^{A,ext}). \quad (28)$$

For description convenience, the factor node $p(\mathbf{H}_k|\alpha_k)$ is replaced by f_k^B . Then the message from variable node \mathbf{H}_k to factor node $p(\mathbf{H}_k|\alpha_k)$ is expressed as

$$\begin{aligned} \mathcal{M}_{\mathbf{H}_k \rightarrow f_k^B}(\mathbf{H}_k) &= \prod_{q,m} \mathcal{M}_{\mathbf{y}_m \rightarrow h_{k,q,m}}(h_{k,q,m}) \\ &= \mathcal{CN}(\mathbf{H}_k; \mathbf{H}_k^{B,pri}, \mathbf{V}^B) \end{aligned} \quad (29)$$

where $\mathbf{H}_k^{B,pri} \in \mathbb{C}^{Q \times M}$ with the (q, m) -th element $h_{k,q,m}^{B,pri} = h_{k,q,m}^{A,ext}$ and $\mathbf{V}^B = \text{diag}([v_{\mathbf{h}_1}^{B,pri}, \dots, v_{\mathbf{h}_M}^{B,pri}]^T \otimes \mathbf{1}_Q^T)$ with $\mathbf{1}_Q$ being an all-one vector of length Q .

Combing the Bernoulli Gaussian prior $p(\mathbf{H}_k|\alpha_k)$, Gaussian message $\mathcal{M}_{\mathbf{H}_k \rightarrow f_k^B}(\mathbf{H}_k)$ and Bernoulli message

$\mathcal{M}_{\alpha_k \rightarrow f_k^B}(\alpha_k)$ in (53), the belief of \mathbf{H}_k at factor node f_k^B is Bernoulli Gaussian and expressed as

$$\begin{aligned} \mathcal{M}_{f_k^B}(\mathbf{H}_k) &\propto \sum_{\alpha_k=0}^1 p(\mathbf{H}_k|\alpha_k) \mathcal{M}_{\alpha_k \rightarrow f_k^B}(\alpha_k) \mathcal{M}_{\mathbf{H}_k \rightarrow f_k^B}(\mathbf{H}_k) \\ &= (1 - \lambda_k^{B,post}) \delta(\mathbf{H}_k) + \lambda_k^{B,post} \mathcal{CN}(\mathbf{H}_k; \boldsymbol{\mu}_k^B; \boldsymbol{\Phi}_k^B) \end{aligned} \quad (30)$$

where

$$\begin{aligned} \lambda_k^{B,post} &= \left(1 + \frac{(1 - \lambda_k^{B,pri}) \cdot \mathcal{CN}(\mathbf{0}; \mathbf{H}_k^{B,pri}, \mathbf{V}^B)}{\lambda_k^{B,pri} \cdot \mathcal{CN}(\mathbf{0}; \mathbf{H}_k^{B,pri}, \mathbf{V}^B + \vartheta_{\mathbf{H}} \mathbf{I})} \right)^{-1} \\ \boldsymbol{\Phi}_k^B &= (\vartheta_{\mathbf{H}}^{-1} \mathbf{I} + (\mathbf{V}^B)^{-1})^{-1} \\ \boldsymbol{\mu}_k^B &= \vartheta_{\mathbf{H}} (\vartheta_{\mathbf{H}} \mathbf{I} + \mathbf{V}^B)^{-1} \cdot \text{vec}(\mathbf{H}_k^{B,pri}). \end{aligned} \quad (31)$$

We note that the mean of $\text{vec}(\mathbf{H}_k)$ with respect to $\mathcal{M}_{f_k^B}(\mathbf{H}_k)$ is

$$\text{vec}(\mathbf{H}_k^{B,post}) = \lambda_k^{B,post} \boldsymbol{\mu}_k^B. \quad (32)$$

Define $l = (m-1)Q + q$. The variance of the (q, m) -th element of \mathbf{H}_k with respect to $\mathcal{M}_{f_k^B}(\mathbf{H}_k)$ is

$$\vartheta_{h_{k,q,m}}^{B,post} = \lambda_k^{B,post} (|\mu_{k,l}^B|^2 + \Phi_{k,l,l}^B) - |h_{k,q,m}^{B,post}|^2. \quad (33)$$

Recall that the relationship between \mathbf{H}_k and \mathbf{h}_m is described as $[\mathbf{H}_1^T, \dots, \mathbf{H}_K^T]^T = [\mathbf{h}_1, \dots, \mathbf{h}_M]$. Through such relationship, we obtain $\mathbf{h}_m^{B,post}$. Following [12], [20], the variance $v_{\mathbf{h}_m}^{B,post}$ is approximated by

$$v_{\mathbf{h}_m}^{B,post} = \frac{1}{KQ} \sum_{k,q} \vartheta_{h_{k,q,m}}^{B,post}. \quad (34)$$

Then the belief of \mathbf{h}_m at Module B is defined as

$$\mathcal{M}_B(\mathbf{h}_m) \sim \mathcal{CN}(\mathbf{h}_m; \mathbf{h}_m^{B,post}, v_{\mathbf{h}_m}^{B,post} \mathbf{I}). \quad (35)$$

The above Gaussian belief assumption is widely used in message passing based iterative algorithms such as Turbo-CS [12], approximate message passing (AMP) [13] and expectation propagation (EP) [30]. Such treatment may loses some information but facilitates the message updates. (Strictly speaking, the belief of the variable corresponds to a factor node instead of a Module. However, the belief of \mathbf{h}_m is associated with K factor nodes $p(\mathbf{H}_k|\alpha_k), \forall k$, making the expression cumbersome. For convenience, we use the definition $\mathcal{M}_B(\mathbf{h}_m)$.)

From the sum-product rule, we have $\mathcal{M}_{\mathbf{h}_m \rightarrow B}(\mathbf{h}_m) = \mathcal{M}_{\mathbf{y}_m \rightarrow \mathbf{h}_m}(\mathbf{h}_m)$. Then the extrinsic message is calculated as

$$\mathcal{M}_{B \rightarrow \mathbf{h}_m}(\mathbf{h}_m) \propto \frac{\mathcal{M}_B(\mathbf{h}_m)}{\mathcal{M}_{\mathbf{y}_m \rightarrow \mathbf{h}_m}(\mathbf{h}_m)}. \quad (36)$$

Given the Gaussian messages $\mathcal{M}_B(\mathbf{h}_m)$ and $\mathcal{M}_{\mathbf{h}_m \rightarrow B}(\mathbf{h}_m)$, the extrinsic message is also Gaussian with its variance and mean respectively given by

$$\begin{aligned} v_{\mathbf{h}_m}^{B,ext} &= \left(\frac{1}{v_{\mathbf{h}_m}^{B,post}} - \frac{1}{v_{\mathbf{h}_m}^{B,pri}} \right)^{-1} \\ \mathbf{h}_m^{B,ext} &= v_{\mathbf{h}_m}^{B,ext} \left(\frac{\mathbf{h}_m^{B,post}}{v_{\mathbf{h}_m}^{B,post}} - \frac{\mathbf{h}_m^{B,pri}}{v_{\mathbf{h}_m}^{B,pri}} \right) \end{aligned} \quad (37)$$

where $\mathbf{h}_m^{B,ext}$ and $v_{\mathbf{h}_m}^{B,ext}$ are respectively used as the input mean and variance of \mathbf{h}_m for Module A, i.e., $\mathbf{h}_m^{A,pri} = \mathbf{h}_m^{B,ext}$ and $v_{\mathbf{h}_m}^{A,pri} = v_{\mathbf{h}_m}^{B,ext}$. Then we have

$$\mathcal{M}_{\mathbf{h}_m \rightarrow \mathbf{y}_m}(\mathbf{h}_m) = \mathcal{M}_{B \rightarrow \mathbf{h}_m}(\mathbf{h}_m) \sim \mathcal{CN}(\mathbf{h}_m; \mathbf{h}_m^{A,pri}, v_{\mathbf{h}_m}^{A,pri} \mathbf{I}) \quad (38)$$

D. Module C: Denoiser of \mathbf{C}_k

Similarly to the process in Module B, each $\mathbf{C}_k, \forall k$ in module C is estimated individually by exploiting the prior $p(\mathbf{C}_k|\alpha_k)$ and the messages from Modules A and D. For description convenience, f_k^C is used to replace factor node $p(\mathbf{C}_k|\alpha_k)$. The message from variable node \mathbf{C}_k to factor node $p(\mathbf{C}_k|\alpha_k)$ is expressed as

$$\mathcal{M}_{\mathbf{C}_k \rightarrow f_k^C}(\mathbf{C}_k) \sim \mathcal{CN}(\mathbf{C}_k; \mathbf{C}_k^{C,pri}, \mathbf{V}^C) \quad (39)$$

where $\mathbf{C}_k^{C,pri} \in \mathbb{C}^{Q \times M}$ with its (q, m) -th element $c_{k,q,m}^{C,pri} = c_{k,q,m}^{A,ext}$ and $\mathbf{V}^C = \text{diag}([v_{c_1}^{C,pri}, \dots, v_{c_M}^{C,pri}]^T \otimes \mathbf{1}_Q^T)$.

With the Bernoulli Gaussian prior $p(\mathbf{C}_k|\alpha_k)$, Gaussian message $\mathcal{M}_{\mathbf{C}_k \rightarrow f_k^C}(\mathbf{C}_k)$ and message $\mathcal{M}_{\alpha_k \rightarrow f_k^C}(\alpha_k)$ in (49), the belief of \mathbf{C}_k at factor node f_k^C is Bernoulli Gaussian and expressed as

$$\begin{aligned} \mathcal{M}_{f_k^C}(\mathbf{C}_k) &\propto \sum_{\alpha_k=0}^1 p(\mathbf{C}_k|\alpha_k) \mathcal{M}_{\alpha_k \rightarrow f_k^C}(\alpha_k) \mathcal{M}_{\mathbf{C}_k \rightarrow f_k^C}(\mathbf{C}_k) \\ &= (1 - \lambda_k^{C,post}) \delta(\mathbf{C}_k) + \lambda_k^{C,post} \mathcal{CN}(\mathbf{C}_k; \boldsymbol{\mu}_k^C; \boldsymbol{\Phi}_k^C) \end{aligned} \quad (40)$$

where

$$\begin{aligned} \lambda_k^{C,post} &= \left(1 + \frac{(1 - \lambda_k^{C,pri}) \cdot \mathcal{CN}(\mathbf{0}; \mathbf{C}_k^{C,pri}, \mathbf{V}^C)}{\lambda_k^{C,pri} \cdot \mathcal{CN}(\mathbf{0}; \mathbf{C}_k^{C,pri}, \mathbf{V}^C + \vartheta_{\mathbf{C}} \mathbf{I})} \right)^{-1} \\ \boldsymbol{\Phi}_k^C &= (\vartheta_{\mathbf{C}}^{-1} \mathbf{I} + (\mathbf{V}^C)^{-1})^{-1} \\ \boldsymbol{\mu}_k^C &= \vartheta_{\mathbf{C}} (\vartheta_{\mathbf{C}} \mathbf{I} + \mathbf{V}^C)^{-1} \cdot \text{vec}(\mathbf{C}_k^{C,pri}). \end{aligned} \quad (41)$$

Then the mean of \mathbf{C}_k with respect to $\mathcal{M}_{f_k^C}(\mathbf{C}_k)$ is

$$\text{vec}(\mathbf{C}_k^{post}) = \lambda_k^{C,post} \boldsymbol{\mu}_k^C. \quad (42)$$

The variance of the (q, m) -th element of \mathbf{C}_k is given by

$$\vartheta_{c_{k,q,m}}^{C,post} = \lambda_k^{C,post} (|\mu_{k,l}^C|^2 + \Phi_{k,l,l}^C) - |c_{k,q,m}^{C,post}|^2. \quad (43)$$

Similarly to (35), we define the belief of \mathbf{c}_m at Module C as

$$\mathcal{M}_C(\mathbf{c}_m) \sim \mathcal{CN}(\mathbf{c}_m; \mathbf{c}_m^{C,post}, v_{\mathbf{c}_m}^{C,post} \mathbf{I}) \quad (44)$$

with the variance $v_{\mathbf{c}_m}^{C,post}$ given by

$$v_{\mathbf{c}_m}^{C,post} = \frac{1}{KQ} \sum_{k,q} \vartheta_{c_{k,q,m}}^{C,post}. \quad (45)$$

Then we calculate the Gaussian extrinsic message with its variance and mean as follows :

$$\begin{aligned} v_{\mathbf{c}_m}^{C,ext} &= \left(\frac{1}{v_{\mathbf{c}_m}^{C,post}} - \frac{1}{v_{\mathbf{c}_m}^{C,pri}} \right)^{-1} \\ \mathbf{c}_m^{A,pri} &= v_{\mathbf{c}_m}^{C,ext} \left(\frac{\mathbf{c}_m^{C,post}}{v_{\mathbf{c}_m}^{C,post}} - \frac{\mathbf{c}_m^{C,pri}}{v_{\mathbf{c}_m}^{C,pri}} \right). \end{aligned} \quad (46)$$

The input mean and variance of \mathbf{c}_m for module A are respectively set as $\mathbf{c}_m^{A,pri} = \mathbf{c}_m^{C,ext}$ and $v_{\mathbf{c}_m}^{A,pri} = v_{\mathbf{c}_m}^{C,ext}$.

E. Module D: Estimation of α_k

In Module D, we calculate the messages $\mathcal{M}_{\alpha_k \rightarrow f_k^B}(\alpha_k)$ and $\mathcal{M}_{\alpha_k \rightarrow f_k^C}(\alpha_k)$ as the inputs of Modules B and C, respectively. Furthermore, the device activity is detected by combing the messages from Modules B and C. According to the sum-product rule, the message from factor node f_k^B to variable node α_k is

$$\begin{aligned} \mathcal{M}_{f_k^B \rightarrow \alpha_k}(\alpha_k) &\propto \int_{\mathbf{H}_k} f_k^B(\mathbf{H}_k, \alpha_k) \cdot \mathcal{M}_{\mathbf{H}_k \rightarrow f_k^B}(\mathbf{H}_k) \\ &= (1 - \pi_k^B) \delta[\alpha_k] + \pi_k^B \delta[1 - \alpha_k] \end{aligned} \quad (47)$$

where

$$\pi_k^B = \left(1 + \frac{\mathcal{CN}(\mathbf{0}; \mathbf{H}_k^{B,pri}, \mathbf{V}^B)}{\mathcal{CN}(\mathbf{0}; \mathbf{H}_k^{B,pri}, \mathbf{V}^B + \vartheta_{\mathbf{H}} \mathbf{I})} \right)^{-1}. \quad (48)$$

Then the message from variable node α_k to factor node f_k^C is

$$\begin{aligned} \mathcal{M}_{\alpha_k \rightarrow f_k^C}(\alpha_k) &\propto \mathcal{M}_{f_k^B \rightarrow \alpha_k} \cdot p(\alpha_k) \\ &= (1 - \lambda_k^{C,pri}) \delta[\alpha_k] + \lambda_k^{C,pri} \delta[1 - \alpha_k] \end{aligned} \quad (49)$$

with

$$\lambda_k^{C,pri} = \frac{\lambda \cdot \pi_k^B}{\lambda \cdot \pi_k^B + (1 - \lambda) \cdot (1 - \pi_k^B)}. \quad (50)$$

Similarly to (47), the message from factor node f_k^C to variable node α_k is

$$\mathcal{M}_{f_k^C \rightarrow \alpha_k} = (1 - \pi_k^C) \delta[\alpha_k] + \pi_k^C \delta[1 - \alpha_k] \quad (51)$$

where

$$\pi_k^C = \left(1 + \frac{\mathcal{CN}(\mathbf{0}; \mathbf{C}_k^{C,pri}, \mathbf{V}^C)}{\mathcal{CN}(\mathbf{0}; \mathbf{C}_k^{C,pri}, \mathbf{V}^C + \vartheta_{\mathbf{C}} \mathbf{I})} \right)^{-1}. \quad (52)$$

Then the message from variable node α_k to factor node f_k^B is

$$\begin{aligned} \mathcal{M}_{\alpha_k \rightarrow f_k^B}(\alpha_k) &\propto \mathcal{M}_{f_k^C \rightarrow \alpha_k} \cdot p(\alpha_k) \\ &= (1 - \lambda_k^{B,pri}) \delta[\alpha_k] + \lambda_k^{B,pri} \delta[1 - \alpha_k] \end{aligned} \quad (53)$$

with

$$\lambda_k^{B,pri} = \frac{\lambda \cdot \pi_k^C}{\lambda \cdot \pi_k^C + (1 - \lambda) \cdot (1 - \pi_k^C)}. \quad (54)$$

Define the belief of α_k at factor node $p(\alpha_k)$ as $\mathcal{M}_k(\alpha_k)$. From the sum-product rule, $\mathcal{M}_k(\alpha_k)$ is given by

$$\begin{aligned} \mathcal{M}_k(\alpha_k) &\propto \mathcal{M}_{f_k^B \rightarrow \alpha_k}(\alpha_k) \mathcal{M}_{f_k^C \rightarrow \alpha_k}(\alpha_k) p(\alpha_k) \\ &= (1 - \lambda_k^{D,post}) \delta[\alpha_k] + \lambda_k^{D,post} \delta[1 - \alpha_k] \end{aligned} \quad (55)$$

where

$$\lambda_k^{D,post} = \frac{\lambda \cdot \pi_k^B \cdot \pi_k^C}{\lambda \cdot \pi_k^B \cdot \pi_k^C + (1 - \lambda) \cdot (1 - \pi_k^B) \cdot (1 - \pi_k^C)}. \quad (56)$$

Clearly, $\lambda_k^{D,post}$ ($0 \leq \lambda_k^{D,post} \leq 1$) indicates the probability that the k -th device is active. Therefore, we adopt a threshold-based strategy for device activity detection as

$$\hat{\alpha}_k = \begin{cases} 1, & \lambda_k^{D,post} \geq \lambda^{thr} \\ 0, & \lambda_k^{D,post} < \lambda^{thr} \end{cases} \quad k = 1, \dots, K \quad (57)$$

where λ^{thr} is a predetermined threshold.

F. Pilot Design and Complexity Analysis

The overall algorithm is summarized in Algorithm 1. In each iteration, the channel mean matrix \mathbf{H} is first updated (step 2-8) and then the channel compensation matrix \mathbf{C} is updated (step 9-15). This is because the power of the channel mean matrix \mathbf{H} is dominant and the iteration process effectively suppresses error propagation. Note that the prior parameters learning is shown in the next section.

As mentioned in Section IV-B, the pilot matrix \mathbf{A} is required to be partial orthogonal. To fulfill this requirement, the pilot symbols $\{a_{k,n}^{(t)}\}_{k=1}^K$ transmitted on the n -th subcarrier at the t -th OFDM symbol has the following property:

$$[a_{1,n}^{(t)}, \dots, a_{k,n}^{(t)}, \dots, a_{K,n}^{(t)}] = \mathbf{u}_i \quad (58)$$

where \mathbf{u}_i is a row vector randomly selected from an orthogonal matrix $\mathbf{U} \in \mathbb{C}^{K \times K}$ and the selected row is different for different n, t . Combing $\mathbf{A} = [\mathbf{\Lambda}_1 \mathbf{E}_1, \dots, \mathbf{\Lambda}_K \mathbf{E}_1]$ and the definition of $\mathbf{\Lambda}_k$ in (4), it is easy to verify the partial orthogonality of \mathbf{A} .

We next show that when \mathbf{U} is the discrete Fourier transform (DFT) matrix, the algorithm complexity can be further reduced. To enable the fast Fourier transform (FFT), the pilot matrix \mathbf{A} is expressed as

$$\mathbf{A} = \text{diag}(\mathbf{S}_1 \mathbf{U}, \dots, \mathbf{S}_Q \mathbf{U}) \mathbf{P} \quad (59)$$

where $\mathbf{S}_q \in \mathbb{R}^{TN/Q \times K}$ is a row selection matrix consisting of TN/Q randomly selected rows from the $K \times K$ identity matrix. (The selected rows are different for different \mathbf{S}_q .) $\mathbf{P} = [\mathbf{P}_1, \dots, \mathbf{P}_K]$ is a column exchange matrix. In specific, in the q -th column of $\mathbf{P}_k \in \mathbb{R}^{KQ \times Q}$, $\forall k$, only the $k + K(q-1)$ -th row is one while the others are zeros. Note that $|a_{k,n}^{(t)}|^2 = P$ and matrix \mathbf{B} has the same fast transform as \mathbf{A} due to the relationship $\mathbf{B} = \mathbf{D}\mathbf{A}$.

By using the FFT, the estimates of $\mathbf{h}_m \in \mathbb{C}^{QK}$ and $\mathbf{c}_m \in \mathbb{C}^{QK}$, $\forall m$ in Module A has the complexity $\mathcal{O}(QM K \log_2 K)$. In modules B, C and D, the estimates of $\mathbf{H}_k \in \mathbb{C}^{Q \times M}$, $\mathbf{C}_k \in \mathbb{C}^{Q \times M}$, and $\alpha_k, \forall k$, involve vector multiplications with the complexity $\mathcal{O}(QMK)$. As a result, the overall complexity of Turbo-MP is $\mathcal{O}(QM K \log_2 K) + \mathcal{O}(QMK)$ per iteration. It is worth noting that the algorithm complexity is linear to the antenna number M and is approximately linear to the device number K .

V. PARAMETERS LEARNING

The prior parameters $\boldsymbol{\theta} = \{\vartheta_{\mathbf{H}}, \vartheta_{\mathbf{C}}, \sigma_w^2, \lambda\}$ used in Turbo-MP are unknown and required to be estimated in practice. In the following, we utilize two machine learning methods, i.e., the EM and the NN approach, to learn these prior parameters.

A. EM Approach

We first use the EM algorithm [23] to learn $\boldsymbol{\theta}$. The EM process is described as

$$\boldsymbol{\theta}^{(i+1)} = \arg \max_{\boldsymbol{\theta}} \mathbb{E} \left[\ln p(\mathbf{Y}, \mathbf{H}, \mathbf{C}, \boldsymbol{\alpha}; \boldsymbol{\theta}) \mid \mathbf{Y}; \boldsymbol{\theta}^{(i)} \right] \quad (60)$$

where $\boldsymbol{\theta}^{(i)}$ is the estimate of $\boldsymbol{\theta}$ at the i -th EM iteration. $\mathbb{E}[\cdot \mid \mathbf{Y}; \boldsymbol{\theta}^{(i)}]$ represents the expectation over the posterior distribution $p(\mathbf{H}, \mathbf{C}, \boldsymbol{\alpha} \mid \mathbf{Y}; \boldsymbol{\theta}^{(i)})$. Note that it is difficult to obtain

Algorithm 1 Turbo Message Passing (Turbo-MP)

Input: $\mathbf{Y}, \mathbf{A}, \mathbf{B}$.

```

1: Initialize  $\boldsymbol{\theta}$  by the EM or NN approach.
while the stopping criterion is not met do
  % Module A: Linear estimation of  $\mathbf{h}_m$ 
  2: Update  $\mathbf{h}_m^{A,post}, v_{\mathbf{h}_m}^{A,post}$  by (19) and (21)-(22),  $\forall m$ .
  3: Update  $\mathbf{h}_m^{A,ext}, v_{\mathbf{h}_m}^{A,ext}$  by (25),  $\forall m$ .
  % Module B: Denoiser of  $\mathbf{H}_k$ 
  4: Update  $\mathbf{H}_k^{B,pri}, \mathbf{V}^B$  by (29),  $\forall k$ .
  5: Update  $\mathcal{M}_{f_k^C \rightarrow \alpha_k}, \mathcal{M}_{\alpha_k \rightarrow f_k^B}$  by (51)-(54),  $\forall k$ .
  6: Update  $\mathbf{H}_k^{B,post}, v_{\mathbf{H}_k}^{B,post}$  by (30)-(33),  $\forall k$ .
  7: Update  $\mathbf{h}_m^{B,post}, v_{\mathbf{h}_m}^{B,post}$  by (34),  $\forall m$ .
  8: Update  $\mathbf{h}_m^{B,ext}, v_{\mathbf{h}_m}^{B,ext}$  by (37),  $\forall m$ .
  % Module A: Linear estimation of  $\mathbf{c}_m$ 
  9: Update  $\mathbf{c}_m^{A,post}, v_{\mathbf{c}_m}^{A,post}$  by (21) and (26),  $\forall m$ .
  10: Update  $\mathbf{c}_m^{A,ext}, v_{\mathbf{c}_m}^{A,ext}$  by (27),  $\forall m$ .
  % Module C: Denoiser of  $\mathbf{C}_k$ 
  11: Update  $\mathbf{C}_k^{C,pri}, \mathbf{V}^C$  by (39),  $\forall k$ .
  12: Update  $\mathcal{M}_{f_k^B \rightarrow \alpha_k}, \mathcal{M}_{\alpha_k \rightarrow f_k^C}$  by (47)-(50),  $\forall k$ .
  13: Update  $\mathbf{C}_k^{C,post}, v_{\mathbf{C}_k}^{C,post}$  by (40)-(43),  $\forall k$ .
  14: Update  $\mathbf{c}_m^{C,post}, v_{\mathbf{c}_m}^{C,post}$  by (45),  $\forall m$ .
  15: Update  $\mathbf{c}_m^{C,ext}, v_{\mathbf{c}_m}^{C,ext}$  by (46),  $\forall m$ .
  % Parameters learning
  16: Update  $\boldsymbol{\theta}$  by EM or NN approach.
end while
17: Update  $\lambda_k^{D,post}$  by (55) and (56),  $\forall k$ .
Output:  $\mathbf{H}_k^{B,post}, \mathbf{C}_k^{C,post}$ , and  $\lambda_k^{D,post}, \forall k$ .

```

the true posterior distribution. Instead, we utilize the message products $\prod_k \mathcal{M}_{f_k^B}(\mathbf{H}_k) \mathcal{M}_{f_k^C}(\mathbf{C}_k) \mathcal{M}_{\alpha_k}(\alpha_k)$ as an approximation. Then we set the derivatives of $\mathbb{E}[\ln p(\mathbf{H}, \mathbf{C}, \mathbf{Y}; \boldsymbol{\theta}) \mid \mathbf{Y}; \boldsymbol{\theta}^{(i)}]$ (with respect to $\vartheta_{\mathbf{H}}$) to zero and obtain

$$\vartheta_{\mathbf{H}}^{(i+1)} = \frac{\sum_k \lambda_k^{D,post} \left(\|\mathbf{H}_k^{B,post}\|_F^2 + \sum_{q,m} \vartheta_{\mathbf{h}_{k,q,m}}^{B,post} \right)}{QM \sum_k \lambda_k^{D,post}}. \quad (61)$$

Similarly, the EM estimate of $\vartheta_{\mathbf{C}}$ is given by

$$\vartheta_{\mathbf{C}}^{(i+1)} = \frac{\sum_k \lambda_k^{D,post} \left(\|\mathbf{C}_k^{C,post}\|_F^2 + \sum_{q,m} \vartheta_{\mathbf{c}_{k,q,m}}^{C,post} \right)}{QM \sum_k \lambda_k^{D,post}}. \quad (62)$$

The EM estimate of σ_w^2 is given by

$$(\sigma_w^2)^{(i+1)} = \frac{1}{MTN} \left\| \mathbf{Y} - \mathbf{A}\mathbf{H}^{B,post} - \mathbf{B}\mathbf{C}^{C,post} \right\|_F^2 + \frac{K}{M} \sum_m \left(v_{\mathbf{h}_m}^{B,post} + \frac{1}{TN} \sum_{i=1}^{TN} D_{i,i}^2 v_{\mathbf{c}_m}^{C,post} \right). \quad (63)$$

The EM estimate of λ is given by

$$\lambda^{(i+1)} = \frac{1}{K} \sum_k \lambda_k^{D,post}. \quad (64)$$

Turbo-MP algorithm with EM approach to learn $\boldsymbol{\theta}$ is shown in Algorithm 1, which we refer to as Turbo-MP-EM. In practice, we can update \mathbf{H} several times and then update \mathbf{C} once to improve the algorithm stability. Besides,

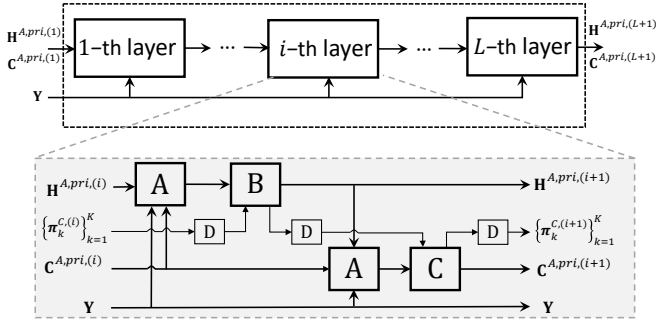


Fig. 4. The block diagram of Turbo-MP-NN. The iterations of Turbo-MP are unfolded into a neural network. Capital letters A, B, C, D denote modules A, B, C, and D, respectively.

it is known the estimation accuracy of $\vartheta_{\mathbf{H}}$, $\vartheta_{\mathbf{C}}$, and λ by EM relies on the accuracy of the posterior approximation in Turbo-MP, and inaccurate estimation may affect the algorithm convergence. Therefore, we recommend to update $\vartheta_{\mathbf{H}}$, $\vartheta_{\mathbf{C}}$, and λ once after \mathbf{H} and \mathbf{C} is updated several times. In the first algorithm iteration, we initialize $\vartheta_{\mathbf{H}}^{(0)} = 1$, $\vartheta_{\mathbf{C}}^{(0)} = 10^{-3}$, and $\lambda^{(0)} = 0.1$. As for the update of σ_w^2 , it is updated in each iteration of Turbo-MP with the initialization $(\sigma_w^2)^{(0)} = \|\mathbf{Y}\|_F^2 / MTN$. However, we find the second part of (63), i.e., $\frac{K}{M} \sum_m (v_{\mathbf{h}_m}^{B,post} + \frac{1}{TN} \sum_{i=1}^{TN} D_{i,i}^2 v_{\mathbf{c}_m}^{C,post})$ is negligible in most cases while in worst case, it continues to rise with the increase of Turbo-MP iterations. Therefore, we delete the second part of (63) in simulation.

B. NN Approach

In deep learning [3], training data can be used to train the parameters of a deep neural network. Inspired by the idea in [31] and [32], we unfold the iterations of Turbo-MP algorithm and regard it as a feed-forward NN termed Turbo-MP-NN. In specific, each iteration represents one layer of the feed-forward NN where θ is seen as the network parameter. We hope that with an appropriately defined loss function, Turbo-MP-NN can adaptively learn θ from the training data.

The structure of Turbo-MP-NN is shown in Fig. 4 which consists of L layers. Each layer has the same structure, i.e., the same linear and non-linear operations following step 2-15 in algorithm 1. To distinguish the estimates at different layer, denote $\mathbf{H}^{A,pri}$, $\mathbf{C}^{A,pri}$ and π_k^C obtained at the $i - 1$ -th layer (iteration) by $\mathbf{H}^{A,pri,(i)}$, $\mathbf{C}^{A,pri,(i)}$, and $\pi_k^{C,(i)}$, respectively. In the i -th layer, the inputs consist of training data \mathbf{Y} and the outputs of the $i - 1$ -th layer including $\mathbf{H}^{A,pri,(i)}$, $\mathbf{C}^{A,pri,(i)}$, and $\{\pi_k^{C,(i)}\}_{k=1}^K$. The loss function is defined as the normalized mean square error (NMSE) of channel estimation given by

$$f(\theta) = \sum_k \frac{\|\mathbf{G}_k - \mathbf{E}_1 \hat{\mathbf{H}}_k^{(L)} - \mathbf{E}_2 \hat{\mathbf{C}}_k^{(L)}\|_F^2}{\|\mathbf{G}_k\|_F^2}, \quad (65)$$

where the estimates $\hat{\mathbf{H}}_k^{(L)}$ and $\hat{\mathbf{C}}_k^{(L)}$ can be $\mathbf{H}_k^{B,post}$ and $\mathbf{C}_k^{C,post}$ or $\mathbf{H}_k^{A,pri}$ and $\mathbf{C}_k^{C,pri}$ obtained in the L -th layer. Note that different from the EM approach, all layers of the NN

have the same θ . To avoid over-fitting, we train the neural network through the layer-by-layer method [32]. Specifically, we begins from the training of the first layer, then first two layers, and finally L layers. θ is initialized following (63). For the first i layers $i = 1, 2, \dots, L$, we optimize θ by using the back-propagation to minimize the loss function

$$f^{(i)}(\theta) = \sum_k \frac{\|\mathbf{G}_k - \mathbf{E}_1 \hat{\mathbf{H}}_k^{(i)} - \mathbf{E}_2 \hat{\mathbf{C}}_k^{(i)}\|_F^2}{\|\mathbf{G}_k\|_F^2}. \quad (66)$$

The detailed training process is shown in Algorithm 2. Once trained, the iteration process of Turbo-MP-NN is illustrated in Algorithm 1. From the simulation, we find that for a fixed sub-block number Q , the change of the learned $\vartheta_{\mathbf{H}}$ and $\vartheta_{\mathbf{C}}$ is linear with respect to the signal-to-ratio $\text{SNR} = P/\sigma_N^2$, and the change of the σ_w^2 minus σ_N^2 is also linear with respect to the SNR. As such, there is no need to train the prior parameters offline at different SNR.

Algorithm 2 Parameter Training of Turbo-MP-NN via Layer-by-Layer Method

Input:

$\mathbf{Y}, \mathbf{A}, \mathbf{B}$.

Output:

- 1: Initialize θ .
 - 2: **for** $i \in [1, L]$ **do**
 - 3: Run i iterations of Turbo-MP algorithm following step 2-15 in Algorithm 1 to obtain $\hat{\mathbf{H}}_k^{(i)}$ and $\hat{\mathbf{C}}_k^{(i)}$.
 - 4: With the loss function $f^{(i)}(\theta)$, use back-propagation to update θ .
 - 5: **end for**
 - 6: **return** θ .
-

VI. NUMERICAL RESULTS

In this section, we evaluate the performance of the proposed Turbo-MP algorithm in channel estimation and device activity detection. The simulation setup is as follows. The BS is equipped with $M = 4$ or 8 antennas. $P = 72$ OFDM subcarriers are allocated for the random access with subcarrier spacing $\Delta f = 15$ kHz. $K = 1000$ devices access the BS and transmit their signals with probability $\lambda = 0.05$ in each timeslot. We adopt the TDL-C channel model with 300 ns r.m.s delay spread and its detailed PDP can be found in TR 38.901 R14 [27]. For Turbo-MP-NN, we randomly generate 10^4 channel realizations for training and the training data is divided into minibatches of size 4. To evaluate the performance of the proposed algorithm, we use the NMSE and the detection error probability $\text{Pe} = \text{P}_{\text{miss}} + \text{P}_{\text{false}}$ as the performance metrics, where $\text{P}_{\text{miss}} = 1/K \sum_k p(\hat{\alpha}_k = 0 | \alpha_k = 1)$ is the probability of miss detection and $\text{P}_{\text{false}} = 1/K \sum_k p(\hat{\alpha}_k = 1 | \alpha_k = 0)$ is the probability of false alarm. Note that for the figures showing CE performance, each data point is averaged over 5000 realizations. For the figures showing ADD performance, at each data point, the accumulative number of the detection errors is larger than 10^3 .

The baseline algorithms are as follows:

- **FD-GMMV-AMP:** The GMMV-AMP algorithm was proposed in [18] to joint estimate the channel response on every subcarrier and detect the active devices based on the frequency-domain system model (5). We assume that the prior parameters including the channel variance on every subcarrier, noise variance σ_N^2 and access probability λ are known for FD-GMMV-AMP. (The following algorithms also use σ_N^2 and λ as known parameters.)
- **TD-Gturbo-MMV:** We adopt the Gturbo-MMV algorithm [11] to achieve the ADD and CE based on the time-domain system model (7), which we refer to as TD-Gturbo-MMV. Note that the denoiser in [11] is extended and applied to $\hat{\mathbf{H}}_k$. We further assume that the BS knows the delay spread of the time-domain channel, and thus the time-domain delay taps can be truncated to reduce the number of channel coefficients to be estimated. Specifically, there are 8 delay taps left with 4 taps at the head and 4 taps at the tail which contains 99% energy of the time-domain channel. As a Bayesian algorithm, TD-Gturbo-MMV requires the PDP of the time-domain channel as prior information. One option is to assume that the exact PDP is known. However, in practice, the exact PDP of each device is difficult to acquire. Therefore, we also consider using a flat PDP with 8 equal-power delay taps as the prior information.
- **TD-Vector AMP:** The Vector AMP [9] algorithm is adopted as another baseline in the time-domain system model (7). Similarly, we extend and apply the denoiser in [9] to $\hat{\mathbf{H}}_k$. TD-Vector AMP with the exact PDP and the flat PDP are both considered.

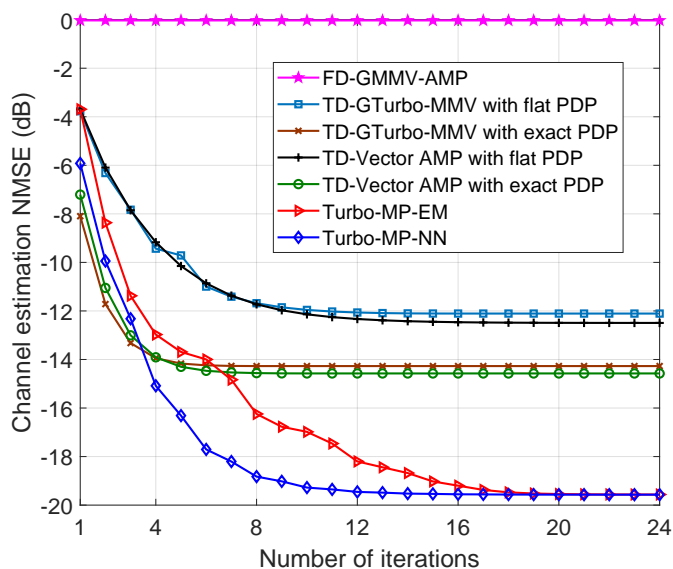


Fig. 5. Channel estimation NMSE versus Iterations number. The number of OFDM symbols $T = 8$ and the sub-block number $Q = 4$. SNR = 10 dB and the number of BS antennas $M = 8$.

Fig. 5 shows the channel estimation NMSE versus the iteration number. Both Turbo-MP-EM and Turbo-MP-NN converge and significantly outperform the other algorithms by more than 6 dB in CE NMSE. Turbo-MP-NN converges faster than Turbo-MP-EM, which implying that the neural

network approach can obtain more accurate prior parameters. In addition, it is seen that FD-GMMV-AMP has a quite poor performance since the average number of the unknown variables on every subcarrier is much larger than the number of the measurements, i.e., $\lambda K \gg T$. Concerning TD-GTurbo-MMV and TD-Vector AMP, there is a performance loss when the PDP is not exactly known.

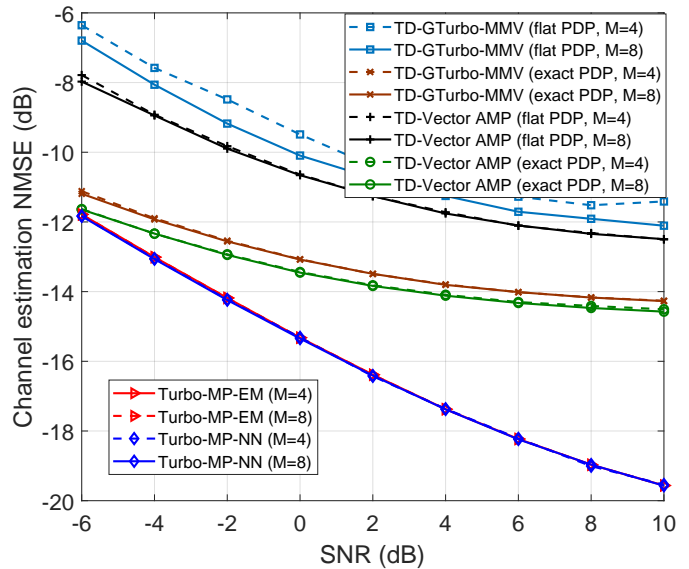


Fig. 6. Channel estimation NMSE versus SNR. The number of OFDM symbols $T = 8$ and the sub-block number $Q = 4$.

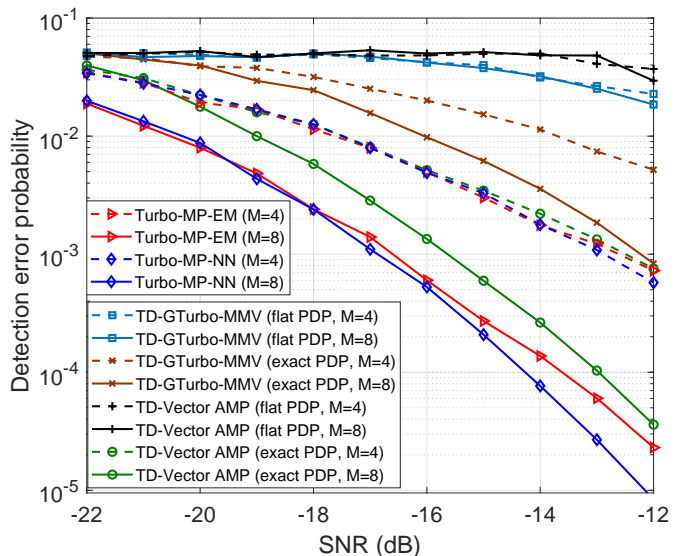


Fig. 7. Detection error probability versus SNR. The number of OFDM symbols $T = 8$ and the sub-block number $Q = 4$.

To further demonstrate the performance superiority of Turbo-MP, Fig. 6 shows the channel estimation NMSE against the SNR. With the increase of the SNR, the performance gap between Turbo-MP and the baselines becomes larger, which suggests that if the BS adopts the Turbo-MP algorithm to reach a high CE performance, the devices will consume much lower power. It is also found that the CE performance of each

algorithm at different BS antenna number is similar. Fig. 7 shows the detection error probability versus the SNR. It is seen that as the antenna number M increases, the detection performances of the tested algorithms improve more than one order of magnitude. This is because the increase of BS antennas leads to a larger dimension of the block-sparsity vector. Among the tested algorithms, Turbo-MP-NN has superiority over the other algorithms especially when antenna number $M = 8$.

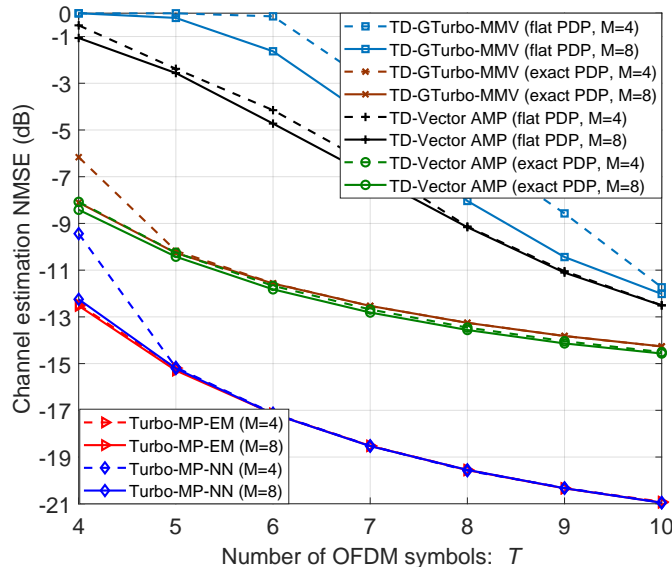


Fig. 8. Channel estimation NMSE versus the number of OFDM symbols T . SNR = 10 dB and the sub-block number $Q = 4$.

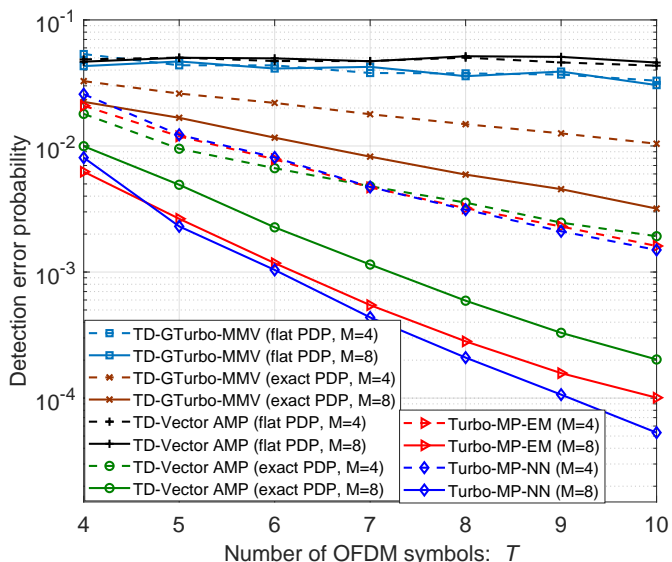


Fig. 9. Detection error probability versus the number of OFDM symbols T . SNR = -15 dB and the sub-block number $Q = 4$.

Fig. 8 shows the channel estimation NMSE against the different number of OFDM symbols. There is a clear performance gap between Turbo-MP and the baselines at different T . Moreover, to reach NMSE = -15 dB, the pilot overhead

($T = 5$) for Turbo-MP is only half of that ($T = 10$) for TD-Vector AMP and TD-Gturbo-MMV with exact PDP. It is shown that our proposed scheme can support mMTC with a dramatically reduced overhead. In Fig. 9, the detection error probability versus OFDM symbols is shown. The trend is similar to Fig. 7. Turbo-MP achieves the best ADD performance at the different number of OFDM symbols.

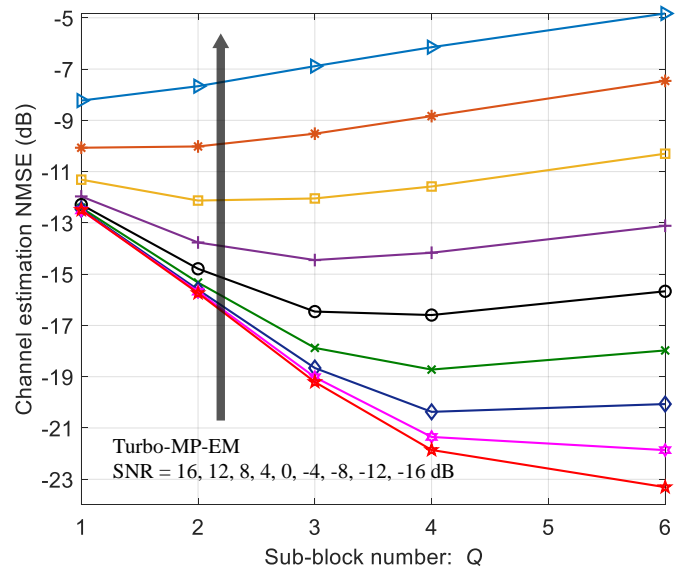


Fig. 10. Channel estimation NMSE versus sub-block number Q at different SNR. The number of BS antennas $M = 8$ and the number of OFDM symbols $T = 10$.

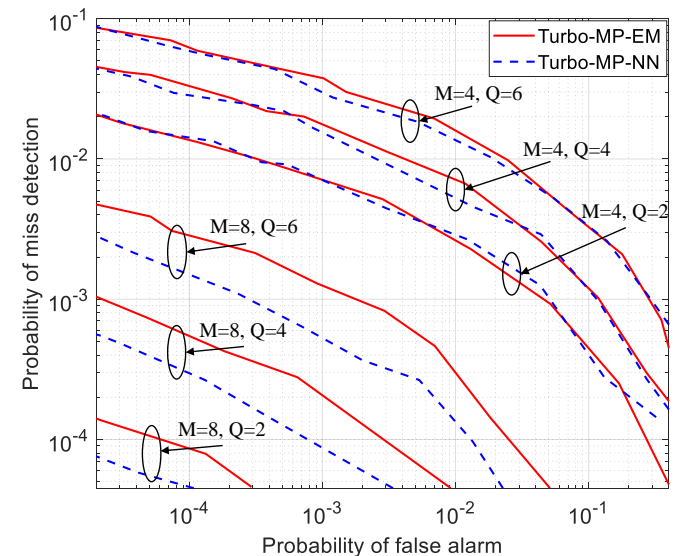


Fig. 11. Probability of miss detection versus probability of false alarm at different sub-block number and antennas' number. SNR = -15 dB and the number of OFDM symbols $T = 10$.

Fig. 10 illustrates the impact of the sub-block number on the CE performance at different SNR, which implies the trade-off between the estimation accuracy and the model accuracy. In specific, it is seen that at low SNR, a smaller sub-block number corresponds to better NMSE performance while the case is contrary at high SNR. The reason is that the NMSE

performance is mainly affected by AWGN at low SNR, in which case a small sub-block number helps to improve the estimation accuracy. However, at high SNR, the CE performance is limited by the model mismatch which can be reduced by increasing the sub-block number. In practice, a proper sub-block number needs to be chosen according to the wireless environment. Fig. 11 shows the miss detection probability versus false alarm probability, where we modify the threshold λ^{thr} to reach the trade-off between miss detection and false alarm. Clearly, Turbo-MP-NN outperforms Turbo-MP-EM at different settings of sub-block number and antennas number. Besides, we see that Turbo-MP with the smaller sub-block number corresponds to a better ADD performance at SNR = -15 dB. This result is consistent with the CE performance in Fig. 10. From Fig. 10 and Fig. 11, we find that the algorithms can achieve excellent device detection performance at relative low SNR while the accurate channel estimation requires a higher SNR. Such observation suggests that when the BS is equipped with multiple antennas, the bottleneck in mMTC is CE instead of ADD.

VII. CONCLUSION

In this paper, a frequency-domain block-wise linear channel model was established in the MIMO-OFDM-based grant-free NOMA system to effectively compensate the channel frequency-selectivity and reduce the number with a small number of variables to be determined in channel estimation. From the perspective of Bayesian inference, we designed the low-complexity Turbo-MP algorithm to solve the ADD and CE problem, where machine learning was incorporated to learn the prior parameters. We numerically show that Turbo-MP designed for the proposed block-wise linear model significantly outperforms the state-of-the-art counterpart algorithms.

REFERENCES

- [1] C. Bockelmann, N. Pratas, H. Nikopour, K. Au, T. Svensson, C. Stefanovic, P. Popovski, and A. Dekorsy, "Massive machine-type communications in 5g: physical and mac-layer solutions," *IEEE Commun. Mag.*, vol. 54, no. 9, pp. 59–65, Sep. 2016.
- [2] *Study on New Radio Access Technology*, Std. TR 38.901 version 14.2.0 Release 14, 3GPP, Sep. 2017.
- [3] B. Wang, L. Dai, Y. Zhang, T. Mir, and J. Li, "Dynamic compressive sensing-based multi-user detection for uplink grant-free noma," *IEEE Commun. Lett.*, vol. 20, no. 11, pp. 2320–2323, Nov. 2016.
- [4] Y. Du, C. Cheng, B. Dong, Z. Chen, X. Wang, J. Fang, and S. Li, "Block-sparsity-based multiuser detection for uplink grant-free NOMA," *IEEE Trans. Wireless Commun.*, vol. 17, no. 12, pp. 7894–7909, Dec. 2018.
- [5] B. K. Jeong, B. Shim, and K. B. Lee, "MAP-based active user and data detection for massive machine-type communications," *IEEE Trans. Veh. Technol.*, vol. 67, no. 9, pp. 8481–8494, Sep. 2018.
- [6] L. Liu, C. Yuen, Y. L. Guan, Y. Li, and C. Huang, "Gaussian message passing for overloaded massive MIMO-NOMA," *IEEE Trans. Wireless Commun.*, vol. 18, no. 1, pp. 210–226, Jan. 2018.
- [7] 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.
- [8] D. L. Donoho, "Compressed sensing," *IEEE Trans. Inf. Theory*, vol. 52, no. 4, pp. 1289–1306, April 2006.
- [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, June 2018.
- [10] Z. Chen, F. Sahrabi, and W. Yu, "Sparse activity detection for massive connectivity," *IEEE Transactions on Signal Processing*, vol. 66, no. 7, pp. 1890–1904, April 2018.
- [11] T. Liu, S. Jin, C. Wen, M. Matthaiou, and X. You, "Generalized channel estimation and user detection for massive connectivity with mixed-ADC massive MIMO," *IEEE Trans. Wireless Commun.*, vol. 18, no. 6, pp. 3236–3250, June 2019.
- [12] J. Ma, X. Yuan, and L. Ping, "Turbo compressed sensing with partial DFT sensing matrix," *IEEE Signal Process. Lett.*, vol. 22, no. 2, pp. 158–161, Feb. 2015.
- [13] M. Bayati and A. Montanari, "The dynamics of message passing on dense graphs, with applications to compressed sensing," *IEEE Trans. Inf. Theory*, vol. 57, no. 2, pp. 764–785, Feb. 2011.
- [14] Q. Zou, H. Zhang, D. Cai, and H. Yang, "A low-complexity joint user activity, channel and data estimation for grant-free massive MIMO systems," *IEEE Signal Process. Lett.*, vol. 27, pp. 1290–1294, July 2020.
- [15] 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, July 2019.
- [16] 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.
- [17] 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 Trans. Veh. Technol.*, vol. 68, no. 12, pp. 12000–12012, Dec. 2019.
- [18] 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, Jan 2020.
- [19] X. Kuai, L. Chen, X. Yuan, and A. Liu, "Structured turbo compressed sensing for downlink massive mimo-ofdm channel estimation," *IEEE Trans. Wireless Commun.*, vol. 18, no. 8, pp. 3813–3826, Aug. 2019.
- [20] Z. Xue, X. Yuan, and Y. Yang, "Denosing-based turbo message passing for compressed video background subtraction," *IEEE Trans. Image Process.*, vol. 30, pp. 2682–2696, Feb. 2021.
- [21] R. Ratasuk, N. Mangalvedhe, D. Bhatoolaul, and A. Ghosh, "LTE-M evolution towards 5G massive MTC," in *in proc. 2017 IEEE Globecom Workshops (GC Wkshps)*, 2017, pp. 1–6.
- [22] F. R. Kschischang, B. J. Frey, and H. J. Loeliger, "Factor graphs and the sum-product algorithm," *IEEE Trans. Inf. Theory*, Feb 2001.
- [23] T. K. Moon, "The expectation-maximization algorithm," *IEEE Signal Processing Magazine*, vol. 13, no. 6, pp. 47–60, Nov 1996.
- [24] M. Chen, Y. Miao, Y. Hao, and K. Hwang, "Narrow band internet of things," *IEEE Access*, vol. 5, pp. 20 557–20 577, Sep. 2017.
- [25] P. Popovski, K. F. Trillingsgaard, O. Simeone, and G. Durisi, "5G wireless network slicing for eMBB, URLLC, and mMTC: A communication-theoretic view," *IEEE Access*, vol. 6, pp. 55 765–55 779, Sep. 2018.
- [26] Y. Li, L. J. Cimini, and N. R. Sollenberger, "Robust channel estimation for ofdm systems with rapid dispersive fading channels," *IEEE Trans. Inf. Theory*, vol. 46, no. 7, pp. 902–915, July 1998.
- [27] *5G; Study on channel model for frequencies from 0.5 to 100 GHz*, Std. TR 38.901 version 14.0.0 Release 14, 3GPP, May 2017.
- [28] S. M. Kay, *Fundamentals of Statistical Signal Processing: Estimation Theory*. Prentice Hall International Editions, 1993.
- [29] C. Berrou and A. Glavieux, "Near optimum error correcting coding and decoding: turbo-codes," *IEEE Trans. Commun.*, Oct. 1996.
- [30] T. P. Minka, "Expectation propagation for approximate bayesian inference," *arXiv preprint arXiv:1301.2294*, Jan. 2013.
- [31] M. Borgerting, P. Schniter, and S. Rangan, "AMP-inspired deep networks for sparse linear inverse problems," *IEEE Trans. Signal Process.*, vol. 65, no. 16, pp. 4293–4308, August 2017.
- [32] X. He, Z. Xue, and X. Yuan, "Learned turbo message passing for affine rank minimization and compressed robust principal component analysis," *IEEE Access*, vol. 7, pp. 140 606–140 617, 2019.