

FDD Massive MIMO Channel Estimation with Arbitrary 2D-Array Geometry

Jisheng Dai, An Liu, and Vincent K. N. Lau

Abstract—This paper addresses the problem of downlink channel estimation in frequency-division duplexing (FDD) massive multiple-input multiple-output (MIMO) systems. The existing methods usually exploit hidden sparsity under a discrete Fourier transform (DFT) basis to estimate the downlink channel. However, there are at least two shortcomings of these DFT-based methods: 1) they are applicable to uniform linear arrays (ULAs) only, since the DFT basis requires a special structure of ULAs; and 2) they always suffer from a performance loss due to the leakage of energy over some DFT bins. To deal with the above shortcomings, we introduce an off-grid model for downlink channel sparse representation with arbitrary 2D-array antenna geometry, and propose an efficient sparse Bayesian learning (SBL) approach for the sparse channel recovery and off-grid refinement. The main idea of the proposed off-grid method is to consider the sampled grid points as adjustable parameters. Utilizing an in-exact block majorization-minimization (MM) algorithm, the grid points are refined iteratively to minimize the off-grid gap. Finally, we further extend the solution to uplink-aided channel estimation by exploiting the angular reciprocity between downlink and uplink channels, which brings enhanced recovery performance.

Index Terms—Channel estimation, massive multiple-input multiple-output (MIMO), sparse Bayesian learning (SBL), majorization-minimization (MM), off-grid refinement.

I. INTRODUCTION

Massive multiple-input multiple-output (MIMO) has attracted significant attention in wireless communications, and has been widely considered as a key candidate technology to meet the capacity demand in 5G wireless networks [1], [2]. To fully reap the benefit of excessive base station (BS) antennas, knowledge of channel state information at the transmitter (CSIT) is essentially required [3]. Many research efforts have been devoted to time-division duplexing (TDD) massive MIMO, because the CSIT in the TDD mode can be obtained by exploiting channel reciprocity, where the pilot-aided training overhead is proportional to the number of active mobile users (MUs) only [4], [5]. However, in the frequency-division duplexing (FDD) mode, the conventional training overhead for the CSIT acquisition grows proportionally with the BS antenna size [6], [7], which can be quite large in massive MIMO systems. Hence, it appears to be an extremely

challenging task to obtain accurate CSIT in FDD massive MIMO systems.

Fortunately, due to the limited local scattering effect in the propagation environment, the elements in the massive MIMO channel are highly correlated. Many works have shown that the effective dimension of a massive MIMO channel is much less than its original dimension [8]–[11]. Specifically, if the BS is equipped with a large uniform linear array (ULA), the massive MIMO channel has an *approximately* sparse representation under the discrete Fourier transform (DFT) basis [10], [12], [13]. Exploiting such hidden sparsity, many efficient downlink channel estimation and feedback algorithms have been proposed in recent years [8], [10], [11], [14]–[19]. Nevertheless, it is worth noting that the validity of the DFT basis as a sparse representation of a massive MIMO channel depends on ULAs. When the antenna geometry deviates from a ULA, the aforementioned methods will fail to work.

DFT-based channel estimation methods always have a performance loss, even for ULA systems, because of the leakage of energy in the DFT basis. As shown in [20]–[22], the DFT basis actually provides a fixed sampling grid that discretely covers the angular domain of the massive MIMO system. Since signals usually come from random directions, the leakage energy caused by direction mismatch is unavoidable. To achieve a better sparse representation, Ding and Rao [20]–[22] considered an overcomplete DFT basis, which corresponds to a denser sampling grid on the angular domain. The overcomplete DFT basis may still lead to a high direction mismatch, if the grid is not sufficiently dense. On the other hand, if a very dense sampling grid is used, the l_1 -norm-based recovery methods may not work well due to high correlation between the basis vectors. To overcome the leakage issue and to generalize for general antenna geometry, dictionary learning techniques were also proposed in [20]–[22]. However, the standard dictionary learning approach has several drawbacks: 1) its convergence is not theoretically guaranteed; and 2) learning a comprehensive dictionary requires collecting a large amount of channel measurements as training samples from all locations in a specific cell, which may pose great challenges in practical implementations.

In this paper, we consider a generic off-grid model for channel sparse representation of massive MIMO systems with an arbitrary 2D-array geometry, and we propose an efficient sparse Bayesian learning (SBL) approach [23], [24] for joint sparse channel recovery and off-grid refinement. The main idea of the proposed method is to consider the sampled grid points as adjustable parameters. Then, we utilize an in-exact block majorization-minimization (MM) algorithm [25], [26] to refine

J. Dai is with the Department of Electronic Engineering, Jiangsu University, Zhenjiang 212013, China, and also with the Department of ECE, The Hong Kong University of Science and Technology, Hong Kong (e-mail: jsdai@ujs.edu.cn).

A. Liu was with the Department of ECE, The Hong Kong University of Science and Technology. He is now with the College of Information Science and Electronic Engineering, Zhejiang University (e-mail: wendaolstr@gmail.com).

V. Lau is with the Department of ECE, The Hong Kong University of Science and Technology, Hong Kong (e-mail: eeknlau@ust.hk).

the grid points iteratively. After several iterations, the refined points will approach the actual directions of arrival/departure, so the proposed method can significantly alleviate direction mismatch in the angular domain. The following summarizes the contributions of this paper.

- **Model-based Off-Grid Sparse Basis**

We provide a novel off-grid model for massive MIMO channel sparse representation with an arbitrary 2D-array geometry. Off-grid models have been applied widely to the direction-of-arrival in array signal processing [27]–[29]. However, the commonly used linear approximation off-grid model does not work well, especially when the grid is not sufficiently fine [30]. Our proposed model avoids using any approximations, and thus can significantly alleviate the modeling error.

- **Joint Sparse Channel Recovery and Off-Grid Refinement with Autonomous Learning**

We propose an SBL-based framework based on in-exact block MM algorithm for joint sparse channel recovery and off-grid refinement. The proposed solution outperforms l_1 -norm recovery [31]–[33],¹ and has an inherent learning capability, so no prior knowledge about the sparsity level, noise variance or direction mismatch is required. We show that the solution converges to the stationary solution of the optimization problem. Simulation results reveal substantial performance gains over the existing state-of-the-art baselines.

- **Enhanced Recovery Performance with Angular Reciprocity**

We further extend the solution to uplink-aided channel estimation by exploiting angular reciprocity² between downlink and uplink channels. Characterizing the joint sparse structure with angular reciprocity was first addressed in [22]. However, it always has a performance loss due to the fact that the joint sparse structure only holds approximately. Our new extension strictly characterizes the joint sparse structure by the inherent mechanism of the off-grid model, bringing enhanced recovery performance.

The rest of the paper is organized as follows. In Section II, we present the system model and review the state-of-the-art DFT-based channel estimation for massive MIMO systems. In Section III, we provide the SBL-based off-grid channel estimation method for a linear array, and then, in Section IV, we extend it to an arbitrary 2D-array geometry. In Section V, we exploit angular reciprocity to improve channel estimation performance. Numerical experiments and discussions follow in Sections VI and VII, respectively.

Notations : \mathbb{C} denotes complex number, $\|\cdot\|_p$ denotes

¹SBL methods include l_1 -norm-based methods as a special case when a maximum *a posteriori* (MAP) optimal estimate is adopted with a fixed Laplace signal prior, and theoretical and empirical results show that SBL methods with better priors can achieve enhanced performance over the l_1 -norm-based methods [24], [34].

²We consider an FDD system, so the reciprocity of the channel realization between the uplink and downlink does not hold. However, the directions of arrival and departure of the uplink are reciprocal with those of the downlink due to the fact that both the uplink and downlink face the same scattering structure [22].

p -norm, $(\cdot)^T$ denotes transpose, $(\cdot)^H$ denotes Hermitian transpose, $(\cdot)^\dagger$ denotes pseudoinverse, \mathbf{I} denotes identity matrix, \mathbf{A}_Ω denotes the sub-matrix formed by collecting the columns from Ω , $\mathcal{CN}(\cdot, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ denotes complex Gaussian distribution with mean $\boldsymbol{\mu}$ and variance $\boldsymbol{\Sigma}$, $\text{supp}(\cdot)$ denotes the set of indices of nonzero elements, $\text{tr}(\cdot)$ denotes trace operator, $\text{diag}(\cdot)$ denotes diagonal operator, and \odot denotes Hadamard product.

II. MASSIVE MIMO CHANNEL MODEL AND EXISTING SOLUTIONS

A. Massive MIMO Channel Model

Consider a massive MIMO system operating in FDD mode. There is one BS with N ($\gg 1$) antennas and K MUs equipped with a single antenna. The array at the BS has an arbitrary geometry in the plane. Without loss of generality, we define the reference plane (the X-Y plane) to be the plane of 2D array, and set the origin of a polar coordinate system to be at the first element of the array, as illustrated in Fig. 1. We consider a flat fading channel, and the downlink channel vector from the BS to the k -th user is given by [35], [36]

$$\mathbf{h}_k = \sum_{c=1}^{N_c} \sum_{s=1}^{N_s} \xi_{c,s}^k \mathbf{a}(\theta_{c,s}^k, \varphi_{c,s}^k), \quad (1)$$

where N_c stands for the number of scattering clusters, N_s stands for the number of sub-paths per scattering cluster, $\xi_{c,s}^k$ is the complex gain of the s -th sub-path in the c -th scattering cluster for the k -th MU, $\theta_{c,s}^k$ and $\varphi_{c,s}^k$ are the corresponding azimuth and elevation angles-of-departure (AoDs), respectively. The steering vector $\mathbf{a}(\theta, \varphi) \in \mathbb{C}^{N \times 1}$ is

$$\mathbf{a}(\theta, \varphi) = [1, e^{-j2\pi \frac{d_2}{\lambda_d} \cos(\varphi) \sin(\theta - \phi_2)}, \dots, e^{-j2\pi \frac{d_N}{\lambda_d} \cos(\varphi) \sin(\theta - \phi_N)}]^T, \quad (2)$$

where (d_n, ϕ_n) is the coordinates of the n -th sensor, and λ_d is the wavelength of the downlink propagation. For a linear array, $\mathbf{a}(\theta, \varphi)$ can be simplified by

$$\mathbf{a}(\theta) = [1, e^{-j2\pi \frac{d_2}{\lambda_d} \sin(\theta)}, \dots, e^{-j2\pi \frac{d_N}{\lambda_d} \sin(\theta)}]^T. \quad (3)$$

Specifically, for a ULA, $\mathbf{a}(\theta)$ becomes

$$\mathbf{a}(\theta) = [1, e^{-j2\pi \frac{d}{\lambda_d} \sin(\theta)}, \dots, e^{-j2\pi \frac{(N-1)d}{\lambda_d} \sin(\theta)}]^T, \quad (4)$$

where d stands for the distance between adjacent sensors.

According to the geometry-based stochastic channel model (GSCM) [37], the number of scattering clusters N_c is usually small, and the sub-paths associated with each scattering cluster are likely to concentrate in a small range around the line-of-sight (LOS) direction between the BS and the scattering cluster. Therefore, only a few dimensions in the angular domain are occupied, which, in return, brings a low dimensional representation for the massive MIMO channels [8], [12], [13].

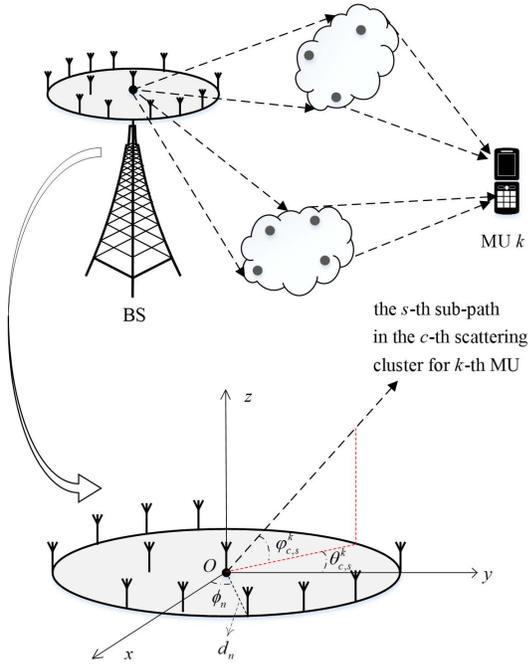


Fig. 1. Illustration of downlink channel model for a massive MIMO system. Note that we define the coordinate system with the X-Y plane being the plane of 2D array, regardless of how the plane of 2D array is placed at the BS (e.g., horizontally or perpendicularly).

B. Review of Downlink Channel Estimation

In this subsection, we review the state-of-the-art DFT-based sparse channel estimation for the downlink channel in an FDD system. Assume that the BS is equipped with a ULA, and it broadcasts a sequence of T training pilot symbols, denoted by $\mathbf{X} \in \mathbb{C}^{T \times N}$, for each MU to estimate the downlink channel. Then, the downlink received signal $\mathbf{y}_k \in \mathbb{C}^{T \times 1}$ at the k -th MU is given by

$$\mathbf{y}_k = \mathbf{X}\mathbf{h}_k + \mathbf{n}_k, \quad (5)$$

where $\mathbf{n}_k \in \mathbb{C}^{T \times 1}$ stands for the additive complex Gaussian noise with each element being zero mean and variance σ^2 in the downlink, and $\text{tr}(\mathbf{X}\mathbf{X}^H) = PTN$ with P/σ^2 measuring the training SNR. Since the number of antennas N at the BS is large, it is unlikely to obtain a robust recovery of \mathbf{h}_k by using conventional channel estimation techniques, e.g., least squares (LS) method. Recently, the emerging compressed sensing (CS) technique has given new interest in the problem of downlink channel estimation with limited training overhead. The main idea of these methods is to find a sparse representation of \mathbf{h}_k in the DFT basis [35], i.e.,

$$\mathbf{h}_k = \mathbf{F}\mathbf{t}_k \quad (6)$$

where $\mathbf{F} \in \mathbb{C}^{N \times N}$ denotes the DFT matrix and \mathbf{t}_k is the sparse representation channel vector. Then, the received signal \mathbf{y}_k in (5) can be formulated as

$$\mathbf{y}_k = \mathbf{X}\mathbf{F}\mathbf{t}_k + \mathbf{n}_k, \quad (7)$$

and the corresponding sparse signal recovery problem is given by

$$\min_{\mathbf{t}_k} \|\mathbf{t}_k\|_0, \quad \text{subject to } \|\mathbf{y}_k - \mathbf{X}\mathbf{F}\mathbf{t}_k\|_2 \leq \epsilon, \quad (8)$$

where ϵ is a constant determined by the upper bound of $\|\mathbf{n}_k\|_2$. As l_0 -norm is non-convex, it is usually relaxed by l_1 -norm, i.e.,

$$\min_{\mathbf{t}_k} \|\mathbf{t}_k\|_1, \quad \text{subject to } \|\mathbf{y}_k - \mathbf{X}\mathbf{F}\mathbf{t}_k\|_2 \leq \epsilon. \quad (9)$$

C. Challenges for the DFT-based Method

In this subsection, we discuss challenges for the DFT-based method. Firstly, this method is applicable to ULAs only, which is explained as follows. The DFT matrix can be written in the form of

$$\mathbf{F} = [\mathbf{f}(-\frac{1}{2}), \mathbf{f}(-\frac{1}{2} + \frac{1}{N}), \dots, \mathbf{f}(\frac{1}{2} - \frac{1}{N})] \quad (10)$$

with

$$\mathbf{f}(x) = \frac{1}{\sqrt{N}} [1, e^{-j2\pi x}, \dots, e^{-j2\pi x(N-1)}]^T, \quad (11)$$

which provides a fixed grid that uniformly covers the range $[-\frac{1}{2}, \frac{1}{2}]$ with $N + 1$ sampling points, i.e., $\{-\frac{1}{2}, (-\frac{1}{2} + \frac{1}{N}), \dots, (\frac{1}{2} - \frac{1}{N}), \frac{1}{2}\}$ ³. As illustrated in (4), the steering vectors of ULAs share the same structure with $\mathbf{f}(x)$. For each sampling point (e.g., the n -th point), we can always find a $\hat{\theta}_n$ in the angular domain such that $\frac{d}{\lambda_d} \sin(\hat{\theta}_n) = -\frac{1}{2} + \frac{n-1}{N}$. Hence, it is equivalent to saying the DFT basis actually provides a fixed sampling grid in the angular domain. When the true azimuth AoDs $\theta_{c,s}^k$ lie on (or, practically, close to) the sampling points $\{\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_{N+1}\}$, the channel vector \mathbf{h}_k definitely has a sparse representation in the DFT basis. Since the sparse property hinges strongly on the shared structure between the DFT basis and the ULA steering, the DFT-based method is applicable to ULAs only.

The other shortcoming of the DFT-based method is that it always has a performance loss, even for ULA systems, due to the leakage of energy. As will be illustrated shortly, the leakage of energy caused by direction mismatch is unavoidable in practice. According to (6), we have

$$\mathbf{t}_k = \mathbf{F}^H \mathbf{h}_k = \sum_{c=1}^{N_c} \sum_{s=1}^{N_s} \xi_{c,s}^k \mathbf{v}(\theta_{c,s}^k), \quad (12)$$

where $\mathbf{v}(\theta_{c,s}^k) = \mathbf{F}^H \mathbf{a}(\theta_{c,s}^k)$. Then, the n -th element of $\mathbf{v}(\theta)$ can be calculated as

$$\begin{aligned} v_n(\theta) &= \frac{1}{\sqrt{N}} \sum_{i=0}^{N-1} e^{j2\pi i(-\frac{1}{2} + \frac{n-1}{N})} e^{-j2\pi \frac{id}{\lambda_d} \sin(\theta)} \\ &= \frac{1}{\sqrt{N}} \frac{1 - e^{j2\pi(\frac{n-1}{N} - \frac{1}{2} - \frac{d}{\lambda_d} \sin(\theta))N}}{1 - e^{j2\pi(\frac{n-1}{N} - \frac{1}{2} - \frac{d}{\lambda_d} \sin(\theta))}} \\ &= \frac{1}{\sqrt{N}} \frac{\sin(\pi \varrho(\theta)N)}{\sin(\pi \varrho(\theta))} e^{j\pi \varrho(\theta)(N-1)}, \end{aligned} \quad (13)$$

where $\varrho(\theta) = \frac{n-1}{N} - \frac{1}{2} - \frac{d}{\lambda_d} \sin(\theta)$. Clearly, the modulus of $v_n(\theta)$ (denoted as $|v_n(\theta)|$) is a periodic function w.r.t. θ , and

³Only the first N points are used in the DFT matrix since $\mathbf{f}(-\frac{1}{2}) = \mathbf{f}(\frac{1}{2})$.

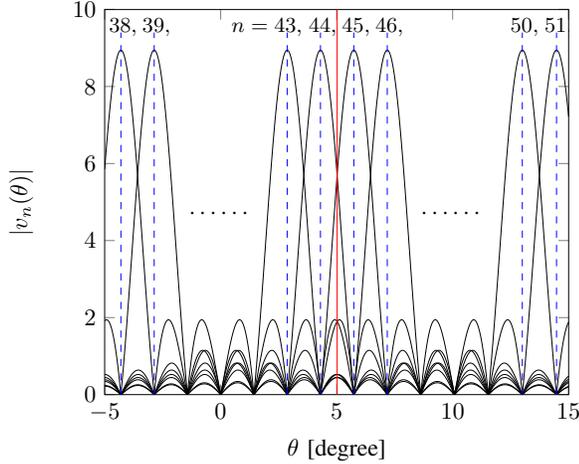


Fig. 2. Illustration of problem of energy leakage for the DFT basis, where the sampling points for the DFT basis in the angular domain are denoted by the dotted blue lines, the true azimuth AoD at $\theta = 5.0198^\circ$ is denoted by the red line, and the distance between the red line and the nearest dotted blue line is called the direction mismatch.

it achieves the maximum value at $\theta = \hat{\theta}_n$. If the true azimuth AoDs are located on the predefined points $\{\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_{N+1}\}$, there is no energy leakage. In practice, however, direction mismatch is unavoidable because signals usually come from random directions. Any direction mismatch will result in the leakage of energy in the DFT basis. Fig. 2 shows an example of the energy leakage, where the ULA is of size $N = 80$ and the inter-antenna spacing is a half wavelength. For an off-grid azimuth AoD (e.g. $\theta^* = 5.0198^\circ$), there is a very serious energy leakage, where both $|v_{44}(\theta^*)|$ and $|v_{45}(\theta^*)|$ have large values, as well as some significant values with $|v_{43}(\theta^*)|$ and $|v_{46}(\theta^*)|$.

To achieve a better sparse representation, [20]–[22] applied an overcomplete DFT basis, which corresponds to a denser sampling grid covering the angular domain with more points. However, if the grid is not sufficiently dense, the overcomplete DFT method may still lead to a high direction mismatch. In order to solve the problem of direction mismatch, as well as apply the sparse channel estimation method to more general array geometry, we propose an efficient SBL-based off-grid method for downlink channel estimation. In the following, we first focus on a simple case, where the BS is equipped with a linear array, and thus only the azimuth angle is involved in the steering vector. This simple case will help to simplify the algorithm design and link the proposed off-grid method with the state-of-the-art methods that are usually applicable to ULAs only. After that, we extend the proposed off-grid method to an arbitrary 2D-array geometry, where both the azimuth and elevation angles are presented in the steering vector. Finally, we exploit angular reciprocity to further improve channel estimation performance.

III. OFF-GRID DOWNLINK CHANNEL ESTIMATION FOR LINEAR ARRAY

In this section, we will propose an efficient SBL-based off-grid method for downlink channel estimation with a linear

array, which includes a ULA as a special case. For ease of exposition, we proceed as follows. We begin by introducing a model-based off-grid basis to handle the direction mismatch for a linear array. Then, we apply this off-grid model to the downlink channel estimation, and an in-exact MM algorithm is provided, as well as its convergence analysis.

A. Off-Grid Basis for Massive MIMO Channels

For ease of notation, we drop the MU's index k and denote the true azimuth AoDs as $\{\theta_l, l = 1, 2, \dots, L\}$, where $L = N_c N_s$. Let $\hat{\vartheta} = \{\hat{\vartheta}_l\}_{l=1}^{\hat{L}}$ be a fixed sampling grid that uniformly covers the angular domain $[-\frac{\pi}{2}, \frac{\pi}{2}]$, where \hat{L} denotes the number of grid points. If the grid is fine enough such that all the true DOAs $\theta_l, l = 1, 2, \dots, L$, lie on (or practically close to) the grid, we can use the following model for \mathbf{h} :

$$\mathbf{h} = \mathbf{A}\mathbf{w}, \quad (14)$$

where $\mathbf{A} = [\mathbf{a}(\hat{\vartheta}_1), \mathbf{a}(\hat{\vartheta}_2), \dots, \mathbf{a}(\hat{\vartheta}_{\hat{L}})] \in \mathbb{C}^{N \times \hat{L}}$, $\mathbf{a}(\theta)$ is a steering vector for a linear array [defined in (3)], and $\mathbf{w} \in \mathbb{C}^{\hat{L} \times 1}$ is a sparse vector whose non-zero elements correspond to the true directions at $\{\theta_l, l = 1, 2, \dots, L\}$. For example, if the \hat{l} -th element of \mathbf{w} is nonzero and the corresponding true direction is θ_l , then we have $\theta_l = \hat{\vartheta}_{\hat{l}}$. Note that \mathbf{A} includes the DFT basis as a special case.

As mentioned in Section II-B, the assumption of the true directions being located on the predefined spatial grid is usually invalid in practice. To handle the direction mismatch, we adopt an off-grid model. Specifically, if $\theta_l \notin \{\hat{\vartheta}_i\}_{i=1}^{\hat{L}}$ and $\hat{\vartheta}_{n_l}, n_l \in \{1, 2, \dots, \hat{L}\}$, is the nearest grid point to θ_l , we write θ_l as

$$\theta_l = \hat{\vartheta}_{n_l} + \beta_{n_l}, \quad (15)$$

where β_{n_l} corresponds to the off-grid gap. Using (15), we have $\mathbf{a}(\theta_l) = \mathbf{a}(\hat{\vartheta}_{n_l} + \beta_{n_l})$. Then, \mathbf{h} can be rewritten as

$$\mathbf{h} = \mathbf{A}(\boldsymbol{\beta})\mathbf{w}, \quad (16)$$

where $\boldsymbol{\beta} = [\beta_1, \beta_2, \dots, \beta_{\hat{L}}]^T$, $\mathbf{A}(\boldsymbol{\beta}) = [\mathbf{a}(\hat{\vartheta}_1 + \beta_1), \mathbf{a}(\hat{\vartheta}_2 + \beta_2), \dots, \mathbf{a}(\hat{\vartheta}_{\hat{L}} + \beta_{\hat{L}})]$, and

$$\beta_{n_l} = \begin{cases} \theta_l - \hat{\vartheta}_{n_l}, & l = 1, 2, \dots, L \\ 0, & \text{otherwise} \end{cases}.$$

Note that with the off-grid basis, the model can significantly alleviate the direction mismatch because there always exists some β_{n_l} making (15) hold exactly. The received signal \mathbf{y} in (5) can be rewritten by

$$\mathbf{y} = \mathbf{X}\mathbf{A}(\boldsymbol{\beta})\mathbf{w} + \mathbf{n} = \boldsymbol{\Phi}(\boldsymbol{\beta})\mathbf{w} + \mathbf{n}, \quad (17)$$

where $\boldsymbol{\Phi}(\boldsymbol{\beta}) \triangleq \mathbf{X}\mathbf{A}(\boldsymbol{\beta})$. Since the coefficient vector $\boldsymbol{\beta}$ is unknown, the current l_1 -norm minimization algorithm can not be applied to the off-grid channel model (17) directly. To jointly recover the sparse signal and refine the grid points, we adopt the SBL algorithm [23], [24], which is one of the most popular approaches for sparse recovery and perturbation calibration. Theoretical and empirical results show that SBL

methods can achieve enhanced performance over l_1 regularized optimization (please also refer to our simulations). In the following, we will discuss how to jointly recover the sparse signal and refine the grid.

B. Sparse Bayesian Learning Formulation

Under the assumption of circular symmetric complex Gaussian noises, we have

$$p(\mathbf{y}|\mathbf{w}, \alpha, \boldsymbol{\beta}) = \mathcal{CN}(\mathbf{y}|\Phi(\boldsymbol{\beta})\mathbf{w}, \alpha^{-1}\mathbf{I}), \quad (18)$$

where $\alpha = \sigma^{-2}$ stands for the noise precision. Since α is usually unknown, we model it as a Gamma hyperprior $p(\alpha) = \Gamma(\alpha; 1 + a, b)$, where we set $a, b \rightarrow 0$ as in [23], [24] so as to obtain a broad hyperprior. We assume a noninformative uniform prior for $\boldsymbol{\beta}$:

$$\boldsymbol{\beta} \sim U\left(\left[-\frac{\pi}{2}, \frac{\pi}{2}\right]^{\hat{L}}\right). \quad (19)$$

Following the commonly used sparse Bayesian model [23], we further assign a non-stationary Gaussian prior distribution with a distinct precision γ_i for each element of \mathbf{w} . Letting $\boldsymbol{\gamma} = [\gamma_1, \gamma_2, \dots, \gamma_{\hat{L}}]^T$, we have

$$p(\mathbf{w}|\boldsymbol{\gamma}) = \mathcal{CN}(\mathbf{w}|\mathbf{0}, \text{diag}(\boldsymbol{\gamma}^{-1})). \quad (20)$$

Similarly, we model γ_i s as independent Gamma distributions, i.e.,

$$p(\boldsymbol{\gamma}) = \prod_{i=1}^{\hat{L}} \Gamma(\gamma_i; 1 + a, b). \quad (21)$$

This two-stage hierarchical prior gives

$$p(\mathbf{w}) = \int_0^\infty \mathcal{CN}(\mathbf{w}|\mathbf{0}, \text{diag}(\boldsymbol{\gamma}^{-1}))p(\boldsymbol{\gamma})d\boldsymbol{\gamma} \\ \propto \prod_{i=1}^{\hat{L}} (b + |w_i|^2)^{-(a+\frac{3}{2})}, \quad (22)$$

which is recognized as encouraging sparsity due to the heavy tails and sharp peak at zero with a small b [23], [34]. In fact, it can be shown that finding a MAP estimate of \mathbf{w} with the prior (22) is equivalent to finding the minimum l_0 -norm solution using FOCUS with $p \rightarrow 0$ [38]. This explains why SBL methods can achieve enhanced performance over the l_1 -norm-based methods. Since directly finding the aforementioned MAP estimate of \mathbf{w} is difficult, SBL methods introduce a two-stage hierarchical prior to get around the problematic MAP estimate. We refer interested readers to Section V of [34] for details.

It is worth noting that the precisions γ_i s in (20) fully indicate the support of \mathbf{w} . For example, if γ_l is large, the l -th element of \mathbf{w} tends to zero; otherwise, the value of the l -th element is significant. As a consequence, once we obtain the precision vector $\boldsymbol{\gamma}$, as well as the off-grid gap $\boldsymbol{\beta}$, the estimated downlink channel \mathbf{h}^e can be obtained by

$$\mathbf{h}^e = \mathbf{A}_\Omega(\boldsymbol{\beta}) (\Phi_\Omega(\boldsymbol{\beta}))^\dagger \mathbf{y}, \quad (23)$$

where $\Omega = \text{supp}(\mathbf{w})$. Therefore, in the rest part of this section, we only need to focus on finding the optimal $\boldsymbol{\gamma}$ and $\boldsymbol{\beta}$. As the

noise precision α is still unknown, we find the most-probable values α^* , $\boldsymbol{\gamma}^*$ and $\boldsymbol{\beta}^*$ together by maximizing the posteriori $p(\alpha, \boldsymbol{\gamma}, \boldsymbol{\beta}|\mathbf{y})$, i.e.,

$$(\alpha^*, \boldsymbol{\gamma}^*, \boldsymbol{\beta}^*) = \arg \max_{\alpha, \boldsymbol{\gamma}, \boldsymbol{\beta}} p(\alpha, \boldsymbol{\gamma}, \boldsymbol{\beta}|\mathbf{y}), \quad (24)$$

or, equivalently,

$$(\alpha^*, \boldsymbol{\gamma}^*, \boldsymbol{\beta}^*) = \arg \max_{\alpha, \boldsymbol{\gamma}, \boldsymbol{\beta}} \ln p(\mathbf{y}, \alpha, \boldsymbol{\gamma}, \boldsymbol{\beta}). \quad (25)$$

The above objective is a high-dimensional non-convex function. It is difficult to directly use the gradient ascent method on the original objective function because gradient ascent is known to have a slow convergence speed, and moreover, the gradient of the original objective function has no closed-form expression. To overcome this challenge, we propose a novel in-exact block MM algorithm to find a stationary point of (25).

C. Overview of the In-exact Block MM Algorithm

The principle behind the block MM algorithm is to iteratively construct a continuous surrogate function (lower bound) for the objective function $\ln p(\mathbf{y}, \alpha, \boldsymbol{\gamma}, \boldsymbol{\beta})$, and then alternately maximize the surrogate function with respect to α , $\boldsymbol{\gamma}$ and $\boldsymbol{\beta}$. The surrogate function is chosen such that the alternating maximization w.r.t. each variable has a closed-form/simple solution.

Specifically, let $\mathcal{U}(\alpha, \boldsymbol{\gamma}, \boldsymbol{\beta}|\hat{\alpha}, \hat{\boldsymbol{\gamma}}, \hat{\boldsymbol{\beta}})$ be the surrogate function constructed at some fixed point $(\hat{\alpha}, \hat{\boldsymbol{\gamma}}, \hat{\boldsymbol{\beta}})$ which satisfies the following properties:

$$\mathcal{U}(\alpha, \boldsymbol{\gamma}, \boldsymbol{\beta}|\hat{\alpha}, \hat{\boldsymbol{\gamma}}, \hat{\boldsymbol{\beta}}) \leq \ln p(\mathbf{y}, \alpha, \boldsymbol{\gamma}, \boldsymbol{\beta}), \quad \forall \alpha, \boldsymbol{\gamma}, \boldsymbol{\beta}, \quad (26)$$

$$\mathcal{U}(\hat{\alpha}, \hat{\boldsymbol{\gamma}}, \hat{\boldsymbol{\beta}}|\hat{\alpha}, \hat{\boldsymbol{\gamma}}, \hat{\boldsymbol{\beta}}) = \ln p(\mathbf{y}, \hat{\alpha}, \hat{\boldsymbol{\gamma}}, \hat{\boldsymbol{\beta}}), \quad (27)$$

$$\left. \frac{\partial \mathcal{U}(\alpha, \hat{\boldsymbol{\gamma}}, \hat{\boldsymbol{\beta}}|\hat{\alpha}, \hat{\boldsymbol{\gamma}}, \hat{\boldsymbol{\beta}})}{\partial \alpha} \right|_{\alpha=\hat{\alpha}} = \left. \frac{\partial \ln p(\mathbf{y}, \alpha, \hat{\boldsymbol{\gamma}}, \hat{\boldsymbol{\beta}})}{\partial \alpha} \right|_{\alpha=\hat{\alpha}}, \quad (28)$$

$$\left. \frac{\partial \mathcal{U}(\hat{\alpha}, \boldsymbol{\gamma}, \hat{\boldsymbol{\beta}}|\hat{\alpha}, \hat{\boldsymbol{\gamma}}, \hat{\boldsymbol{\beta}})}{\partial \boldsymbol{\gamma}} \right|_{\boldsymbol{\gamma}=\hat{\boldsymbol{\gamma}}} = \left. \frac{\partial \ln p(\mathbf{y}, \hat{\alpha}, \boldsymbol{\gamma}, \hat{\boldsymbol{\beta}})}{\partial \boldsymbol{\gamma}} \right|_{\boldsymbol{\gamma}=\hat{\boldsymbol{\gamma}}}, \quad (29)$$

$$\left. \frac{\partial \mathcal{U}(\hat{\alpha}, \hat{\boldsymbol{\gamma}}, \boldsymbol{\beta}|\hat{\alpha}, \hat{\boldsymbol{\gamma}}, \hat{\boldsymbol{\beta}})}{\partial \boldsymbol{\beta}} \right|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}} = \left. \frac{\partial \ln p(\mathbf{y}, \hat{\alpha}, \hat{\boldsymbol{\gamma}}, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}}. \quad (30)$$

Then, we update α , $\boldsymbol{\gamma}$, and $\boldsymbol{\beta}$ as

$$\alpha^{(i+1)} = \arg \max_{\alpha} \mathcal{U}(\alpha, \boldsymbol{\gamma}^{(i)}, \boldsymbol{\beta}^{(i)}|\alpha^{(i)}, \boldsymbol{\gamma}^{(i)}, \boldsymbol{\beta}^{(i)}), \quad (31)$$

$$\boldsymbol{\gamma}^{(i+1)} = \arg \max_{\boldsymbol{\gamma}} \mathcal{U}(\alpha^{(i+1)}, \boldsymbol{\gamma}, \boldsymbol{\beta}^{(i)}|\alpha^{(i+1)}, \boldsymbol{\gamma}^{(i)}, \boldsymbol{\beta}^{(i)}), \quad (32)$$

$$\boldsymbol{\beta}^{(i+1)} = \arg \max_{\boldsymbol{\beta}} \mathcal{U}(\alpha^{(i+1)}, \boldsymbol{\gamma}^{(i+1)}, \boldsymbol{\beta}|\alpha^{(i+1)}, \boldsymbol{\gamma}^{(i+1)}, \boldsymbol{\beta}^{(i)}), \quad (33)$$

where $(\cdot)^{(i)}$ stands for the i -th iteration. The overall flow of the block MM algorithm is given in Fig. 3. The update rules (31)–(33) guarantee the convergence of the block MM algorithm as follows.

Lemma 1. The update rules (31)–(33) give a non-decreasing sequence $\ln p(\mathbf{y}, \alpha^{(i)}, \boldsymbol{\gamma}^{(i)}, \boldsymbol{\beta}^{(i)})$, $i = 1, 2, 3, \dots$

Proof: See Appendix A. ■

In the block MM algorithm, we need to obtain the optimal solutions for the maximization problems in (31)–(33). However, the maximization problem w.r.t. $\boldsymbol{\beta}$ in (33) is non-convex

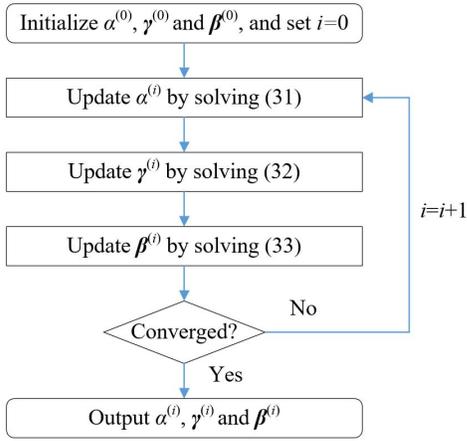


Fig. 3. The overall flow of the block MM algorithm.

and it is difficult to find its optimal solution. Therefore, in this paper, we use an in-exact MM algorithm where $\beta^{(i+1)}$ is obtained by applying a simple one-step update. In the following, we will discuss the choice of the surrogate function and the hyperparameter updates for α , γ and β , respectively. Despite the in-exact update for β , we will prove that the in-exact block MM algorithm still converges to a stationary solution of the optimization problem (25).

D. Detailed Implementations

To update α , γ and β , we first have to choose an appropriate surrogate function $\mathcal{U}(\alpha, \gamma, \beta | \cdot, \cdot, \cdot)$ that satisfies the properties mentioned in (26)–(30). Inspired by the expectation-maximization (EM) algorithm [39], we use the corresponding lower bound function as the surrogate function; i.e., for any fixed point $(\hat{\alpha}, \hat{\gamma}, \hat{\beta})$, we construct the surrogate function as

$$\begin{aligned} \mathcal{U}(\alpha, \gamma, \beta | \hat{\alpha}, \hat{\gamma}, \hat{\beta}) \\ = \int p(\mathbf{w} | \mathbf{y}, \hat{\alpha}, \hat{\gamma}, \hat{\beta}) \ln \frac{p(\mathbf{w}, \mathbf{y}, \alpha, \gamma, \beta)}{p(\mathbf{w} | \mathbf{y}, \hat{\alpha}, \hat{\gamma}, \hat{\beta})} d\mathbf{w}, \end{aligned} \quad (34)$$

and we have the following lemma.

Lemma 2. All the properties in (26)–(30) hold true with the surrogate function $\mathcal{U}(\alpha, \gamma, \beta | \cdot, \cdot, \cdot)$ given in (34).

Proof: See Appendix B. \blacksquare

Note that, from (18) and (20), $p(\mathbf{w} | \mathbf{y}, \alpha, \gamma, \beta)$ is complex Gaussian [23], [27]:

$$p(\mathbf{w} | \mathbf{y}, \alpha, \gamma, \beta) = \mathcal{CN}(\mathbf{w} | \boldsymbol{\mu}(\alpha, \gamma, \beta), \boldsymbol{\Sigma}(\alpha, \gamma, \beta)), \quad (35)$$

where

$$\begin{aligned} \boldsymbol{\mu}(\alpha, \gamma, \beta) &= \alpha \boldsymbol{\Sigma}(\alpha, \gamma, \beta) \boldsymbol{\Phi}^H(\beta) \mathbf{y}, \\ \boldsymbol{\Sigma}(\alpha, \gamma, \beta) &= (\alpha \boldsymbol{\Phi}^H(\beta) \boldsymbol{\Phi}(\beta) + \text{diag}(\gamma))^{-1}. \end{aligned}$$

With $\mathcal{U}(\alpha, \gamma, \beta | \cdot, \cdot, \cdot)$, we discuss the hyperparameter updates for α , γ and β , respectively, as follows.

1) *Update for α :* The maximization problem in (31) has a simple and closed-form solution:

Lemma 3. The optimization problem (31) has a unique solution:

$$\alpha^{(i+1)} = \frac{T + a}{b + \eta(\alpha^{(i)}, \gamma^{(i)}, \beta^{(i)})}, \quad (36)$$

where

$$\begin{aligned} \eta(\alpha, \gamma, \beta) &= \text{tr}(\boldsymbol{\Phi}(\beta) \boldsymbol{\Sigma}(\alpha, \gamma, \beta) \boldsymbol{\Phi}^H(\beta)) \\ &\quad + \|\mathbf{y} - \boldsymbol{\Phi}(\beta) \boldsymbol{\mu}(\alpha, \gamma, \beta)\|_2^2. \end{aligned}$$

Proof: See Appendix C. \blacksquare

2) *Update for γ :* The maximization problem in (32) also has a simple and closed-form solution:

Lemma 4. The optimization problem (32) has a unique solution:

$$\gamma_l^{(i+1)} = \frac{a + 1}{b + [\boldsymbol{\Xi}(\alpha^{(i+1)}, \gamma^{(i)}, \beta^{(i)})]_{ll}}, \quad \forall l, \quad (37)$$

where

$$\boldsymbol{\Xi}(\alpha, \gamma, \beta) = \boldsymbol{\Sigma}(\alpha, \gamma, \beta) + \boldsymbol{\mu}(\alpha, \gamma, \beta) \boldsymbol{\mu}^H(\alpha, \gamma, \beta).$$

Proof: See Appendix D. \blacksquare

3) *Update for β :* Since the maximization problem (33) is non-convex and it is difficult to find its optimal solution, we apply gradient update on the objective function of (33) and obtain a simple one-step update for β . We name the procedure of updating β as grid refining, because it is related with the modeling error caused by the off-grid gap. The derivative of the objective function in (33) w.r.t. β can be calculated as

$$\boldsymbol{\zeta}_\beta^{(i)} = [\zeta^{(i)}(\beta_1), \zeta^{(i)}(\beta_2), \dots, \zeta^{(i)}(\beta_{\hat{L}})]^T, \quad (38)$$

with

$$\begin{aligned} \zeta^{(i)}(\beta_l) &= 2\text{Re} \left((\mathbf{a}'(\hat{\vartheta}_l + \beta_l))^H \mathbf{X}^H \mathbf{X} (\mathbf{a}(\hat{\vartheta}_l + \beta_l)) \right) \cdot c_1^{(i)} \\ &\quad + 2\text{Re} \left((\mathbf{a}'(\hat{\vartheta}_l + \beta_l))^H \mathbf{X}^H \mathbf{c}_2^{(i)} \right), \end{aligned} \quad (39)$$

where $c_1^{(i)} = -\alpha^{(i+1)} (\chi_{ll}^{(i)} + |\mu_l^{(i)}|^2)$, $c_2^{(i)} = \alpha^{(i+1)} ((\mu_l^{(i)})^* \mathbf{y}_{-l}^{(i)} - \mathbf{X} \sum_{j \neq l} \chi_{jl}^{(i)} \mathbf{a}(\hat{\vartheta}_j + \beta_j))$, $\mathbf{y}_{-l}^{(i)} = \mathbf{y} - \mathbf{X} \cdot \sum_{j \neq l} (\mu_j^{(i)} \cdot \mathbf{a}(\hat{\vartheta}_j + \beta_j))$, $\mathbf{a}'(\hat{\vartheta}_j + \beta_l) = d\mathbf{a}(\hat{\vartheta}_j + \beta_l) / d\beta_l$, $\mu_l^{(i)}$ and $\chi_{jl}^{(i)}$ denote the l -th element and the (j, l) -th element of $\boldsymbol{\mu}(\alpha^{(i+1)}, \gamma^{(i+1)}, \beta^{(i)})$ and $\boldsymbol{\Sigma}(\alpha^{(i+1)}, \gamma^{(i+1)}, \beta^{(i)})$, respectively. The detailed derivation for (39) can be found in Appendix E. It is clear that the optimal solution for β is hard to obtain. Fortunately, due to the convergence property illustrated in (87), we just have to find a suboptimal solution that increases the value of the objective function step by step. The most popular numerical method is to update the value of β_l in the derivative direction, i.e.,

$$\beta^{(i+1)} = \beta^{(i)} + \Delta_\beta \cdot \boldsymbol{\zeta}_\beta^{(i)}, \quad (40)$$

where Δ_β is the stepsize. Here, we can resort to backtracking line search [40] to determine the maximum stepsize Δ_β , which ensures that the objective value can be strictly decreased before reaching the stationary point. The complexity of choosing the right stepsize mainly depends on calculating the cost function. If the number of calculations of the cost function in every backtracking line search is R_b , the computational complexity is $\mathcal{O}(R_b T \hat{L}^2)$ per iteration for parameter tuning. Note that the complexity in calculating $\boldsymbol{\zeta}$ is $\mathcal{O}(TN\hat{L})$ per iteration. This suggests the computational requirement of updating β is $\mathcal{O}(R_b T \hat{L}^2)$ per iteration, because \hat{L} is usually larger than

N . To reduce the computational complexity, we alternatively use a fixed stepsize to update β :

$$\beta^{(i+1)} = \beta^{(i)} + \frac{r_\theta}{100} \cdot \text{sign}(\zeta_\beta^{(i)}), \quad (41)$$

where $r_\theta = \pi/\hat{L}$ stands for the grid interval, and $\text{sign}(\cdot)$ stands for the signum function whose computational complexity is negligible. The motivation for choosing this fixed stepsize comes from the fact that a tiny difference between the obtained angles and the true angles does not affect the channel estimation performance much. The term $\frac{r_\theta}{100}$ guarantees that the final gap is smaller than 1% of r_θ , and the (approximate) true values may be attained within 100 iterations in the worst case.

Finally, following are some practical implementation tips for the proposed method. Empirical evidence shows that the proposed method usually converges within 30 iterations, and it remains very robust to the choice of initialization. We can simply set the initialization as follows: $a = b = 0.0001$, $\alpha^{(0)} = 1$, $\gamma^{(0)} = \mathbf{1}$, and $\beta^{(0)} = \mathbf{0}$. Note that MATLAB codes have been made available online at <https://sites.google.com/site/jsdaiustc/publication>.

E. Convergence Analysis and Discussion

From Lemma 1, the sequence $\ln p(\mathbf{y}, \alpha^{(i)}, \gamma^{(i)}, \beta^{(i)})$, $i = 1, 2, 3, \dots$, is non-decreasing and it converges to a limit because the evidence function has the upper bound of 1. In the following, we further prove that the sequence of iterates generated by the algorithm converges to a stationary point.

Theorem 5. For the surrogate function defined in (34), if variables are iteratively updated by (36), (37) and (40), the iterates generated by the in-exact block MM algorithm converge to a stationary solution of the optimization problem (25).

Proof: See Appendix F. ■

Next, we address the difference between the proposed in-exact block MM algorithm and the EM algorithm. The standard SBL method usually exploits the EM algorithm to perform the Bayesian inference. The EM algorithm iteratively constructs the same lower bound as in (34), and simultaneously updates α , γ and β by

$$\begin{aligned} & (\alpha^{(i+1)}, \gamma^{(i+1)}, \beta^{(i+1)}) \\ &= \arg \max_{\alpha, \gamma, \beta} \mathcal{U}(\alpha, \gamma, \beta | \mathbf{y}, \alpha^{(i)}, \gamma^{(i)}, \beta^{(i)}). \end{aligned} \quad (42)$$

The EM algorithm can find a local optimal solution and its convergence can be guaranteed, if the joint maximization problem in (42) is solvable. Unfortunately, in our problem, (42) is non-convex and is intractable in the presence of β . Hence, the EM algorithm cannot be directly applied to our problem.

One commonly used method to address this challenge is to first obtain a convex approximation of the non-convex problem (42) using the linear approximation off-grid model [27], [28], and then update α , γ and β by solving the resulting convex approximation problem in each iteration. Specifically, by replacing the steering vector $\mathbf{a}(\theta_l) = \mathbf{a}(\hat{\vartheta}_{n_l} + \beta_{n_l})$ with the linear approximation

$$\mathbf{a}(\theta_l) \approx \mathbf{a}(\hat{\vartheta}_{n_l}) + \beta_{n_l} \cdot \mathbf{a}'(\hat{\vartheta}_{n_l}), \quad (43)$$

the surrogate function in (34) becomes convex and thus can be maximized efficiently. However, we do not adopt this linear approximation method, because if the grid is not sufficiently fine, (43) may lead to a large modeling error, and the final channel estimation performance will be poor.

Finally, we discuss the computational complexity of the proposed method, whose main computational burden is given as follows.

- The complexities in calculating $\Sigma(\alpha, \gamma, \beta)$ and $\mu(\alpha, \gamma, \beta)$ in each iteration are $\mathcal{O}(T\hat{L}^2)$ and $\mathcal{O}(\hat{L}^2)$, respectively.
- The complexities in updating α and γ in each iteration are $\mathcal{O}(T\hat{L}^2)$ and $\mathcal{O}(\hat{L})$, respectively.
- The complexity in updating β is $\mathcal{O}(TN\hat{L})$ per iteration if the fixed stepsize update is used.

This suggests the total computational requirement of the proposed method is $\mathcal{O}(T\hat{L}^2)$ per iteration, because \hat{L} is usually larger than N .

IV. EXTENSION TO ARBITRARY 2D-ARRAY GEOMETRY

In this section, we extend the proposed off-grid method to an arbitrary 2D-array geometry, where the steering vector $\mathbf{a}(\theta, \varphi)$ [defined in (2)] contains both azimuth and elevation angles. Following the convention in Section III, we adopt a fixed sampling grid $\hat{\vartheta} = \{\hat{\vartheta}_l\}_{l=1}^{\hat{L}}$ to uniformly cover the azimuth domain $[-\pi, \pi]$, and define the off-grid gap β similarly to (16). Then, the received signal \mathbf{y} in (5) can be rewritten by

$$\mathbf{y} = \Phi(\beta, \hat{\varphi})\mathbf{w} + \mathbf{n}, \quad (44)$$

where $\Phi(\beta, \hat{\varphi}) = \mathbf{X}\mathbf{A}(\beta, \hat{\varphi})$, $\mathbf{A}(\beta, \hat{\varphi}) = [\mathbf{a}(\hat{\vartheta}_1 + \beta_1, \hat{\varphi}_1), \mathbf{a}(\hat{\vartheta}_2 + \beta_2, \hat{\varphi}_2), \dots, \mathbf{a}(\hat{\vartheta}_{\hat{L}} + \beta_{\hat{L}}, \hat{\varphi}_{\hat{L}})]$, and

$$\hat{\varphi}_{n_l} = \begin{cases} \varphi_l, & l = 1, 2, \dots, L \\ 0, & \text{otherwise} \end{cases}.$$

Note that the definition of n_l can be found in (15). Compared with (17), the only difference is that the measurement matrix $\Phi(\beta, \hat{\varphi})$ contains an additional unknown variable (i.e., $\hat{\varphi}$). In the following, we will show how to blend $\hat{\varphi}$ with the proposed off-grid method.

In the sparse Bayesian learning formulation for the new model (44), almost all the results in Section III-B remain unchanged, except that (18) is replaced by

$$p(\mathbf{y} | \mathbf{w}, \alpha, \beta, \hat{\varphi}) = \mathcal{CN}(\mathbf{y} | \Phi(\beta, \hat{\varphi})\mathbf{w}, \alpha^{-1}\mathbf{I}), \quad (45)$$

and the optimization problem (25) is modified by

$$(\alpha^*, \gamma^*, \beta^*, \hat{\varphi}^*) = \arg \max_{\alpha, \gamma, \beta, \hat{\varphi}} \ln p(\mathbf{y}, \alpha, \gamma, \beta, \hat{\varphi}). \quad (46)$$

For ease of notation, let $\Theta \triangleq \{\alpha, \gamma, \beta, \hat{\varphi}\}$. At some fixed point $\hat{\Theta} = \{\hat{\alpha}, \hat{\gamma}, \hat{\beta}, \hat{\varphi}\}$, we construct the surrogate function as

$$\mathcal{U}(\Theta | \hat{\Theta}) = \int p(\mathbf{w} | \mathbf{y}, \hat{\Theta}) \ln \frac{p(\mathbf{w}, \mathbf{y}, \Theta)}{p(\mathbf{w} | \mathbf{y}, \hat{\Theta})} d\mathbf{w}. \quad (47)$$

Then, in the maximization step of the $(i+1)$ -th iteration, we update α, γ, β and $\hat{\varphi}$ as

$$\alpha^{(i+1)} = \arg \max_{\alpha} \mathcal{U}(\alpha, \gamma^{(i)}, \beta^{(i)}, \hat{\varphi}^{(i)} | \Theta^{(i)}), \quad (48)$$

$$\gamma^{(i+1)} = \arg \max_{\gamma} \mathcal{U}(\alpha^{(i+1)}, \gamma, \beta^{(i)}, \hat{\varphi}^{(i)} | \Theta_1^{(i)}), \quad (49)$$

$$\beta^{(i+1)} = \arg \max_{\beta} \mathcal{U}(\alpha^{(i+1)}, \gamma^{(i+1)}, \beta, \hat{\varphi}^{(i)} | \Theta_2^{(i)}), \quad (50)$$

$$\hat{\varphi}^{(i+1)} = \arg \max_{\hat{\varphi}} \mathcal{U}(\alpha^{(i+1)}, \gamma^{(i+1)}, \beta^{(i+1)}, \hat{\varphi} | \Theta_3^{(i)}), \quad (51)$$

where $\Theta_j^{(i)}$ denotes the first j elements of $\Theta^{(i)}$ coming from the $(i+1)$ -th iteration. Applying the results in Section III-D directly, we can obtain the solutions to (48)–(50):

$$\alpha^{(i+1)} = \frac{T + a}{b + \eta(\alpha^{(i)}, \gamma^{(i)}, \beta^{(i)}, \hat{\varphi}^{(i)})}, \quad (52)$$

$$\gamma_l^{(i+1)} = \frac{a + 1}{b + [\Xi(\alpha^{(i+1)}, \gamma^{(i)}, \beta^{(i)}, \hat{\varphi}^{(i)})]_l}, \quad \forall l, \quad (53)$$

$$\beta^{(i+1)} = \beta^{(i)} + \frac{r\theta}{100} \cdot \text{sign}(\zeta_{\beta}^{(i)}), \quad (54)$$

where $\eta(\alpha, \gamma, \beta, \hat{\varphi})$, $\Xi(\alpha, \gamma, \beta, \hat{\varphi})$, and ζ_{β} follow the same definitions as in Section III-D, except for the tiny modification of adding the new variable $\hat{\varphi}$ after β for all the equalities.

What remains is to obtain the update for $\hat{\varphi}$. The last maximization problem (51) is similar to (50), where the objective function w.r.t $\hat{\varphi}$ is also non-convex. Hence, we apply the same one-step update for $\hat{\varphi}$ as in (54). Following similar procedures to those in Appendix-E, we can obtain the derivative of the objective function w.r.t $\hat{\varphi}_l$ as

$$\zeta_{\varphi}^{(i)} = [\zeta^{(i)}(\hat{\varphi}_1), \zeta^{(i)}(\hat{\varphi}_2), \dots, \zeta^{(i)}(\hat{\varphi}_L)]^T, \quad (55)$$

with

$$\begin{aligned} \zeta^{(i)}(\hat{\varphi}_l) = & 2\text{Re} \left((\mathbf{a}'_{\varphi}(\hat{\vartheta}_l + \beta_l^{(i+1)}, \hat{\varphi}_l))^H \mathbf{X}^H \mathbf{X} (\mathbf{a}(\hat{\vartheta}_l + \beta_l^{(i+1)}, \hat{\varphi}_l)) \right) \cdot c_{\varphi 1}^{(i)} \\ & + 2\text{Re} \left((\mathbf{a}'_{\varphi}(\hat{\vartheta}_l + \beta_l^{(i+1)}, \hat{\varphi}_l))^H \mathbf{X}^H \mathbf{c}_{\varphi 2}^{(i)} \right), \end{aligned} \quad (56)$$

where $c_{\varphi 1}^{(i)} = -\alpha^{(i+1)}(\chi_{\varphi, ll}^{(i)} + |\mu_{\varphi, ll}^{(i)}|^2)$, $\mathbf{c}_{\varphi 2}^{(i)} = \alpha^{(i+1)}((\mu_{\varphi, ll}^{(i)})^* \mathbf{y}_{\varphi-l} - \mathbf{X} \sum_{j \neq l} \chi_{\varphi, jl}^{(i)} \mathbf{a}(\hat{\vartheta}_j + \beta_j^{(i+1)}, \hat{\varphi}_j))$, $\mathbf{y}_{\varphi-l} = \mathbf{y} - \mathbf{X} \cdot \sum_{j \neq l} (\mu_{\varphi, j}^{(i)} \cdot \mathbf{a}(\hat{\vartheta}_j + \beta_j^{(i+1)}, \hat{\varphi}_j))$, $\mathbf{a}'_{\varphi}(\hat{\vartheta}_j + \beta_j^{(i+1)}, \hat{\varphi}_j) = d\mathbf{a}(\hat{\vartheta}_j + \beta_j^{(i+1)}, \hat{\varphi}_j)/d\hat{\varphi}_j$, $\mu_{\varphi, l}^{(i)}$ and $\chi_{\varphi, jl}^{(i)}$ denote the l -th element and the (j, l) -th element of $\boldsymbol{\mu}(\alpha^{(i+1)}, \gamma^{(i+1)}, \beta^{(i+1)}, \hat{\varphi}^{(i)})$ and $\boldsymbol{\Sigma}(\alpha^{(i+1)}, \gamma^{(i+1)}, \beta^{(i+1)}, \hat{\varphi}^{(i)})$, respectively. With (55), we are able to update $\hat{\varphi}$ similarly to (40).

There are some important tips for practical implementations. We note that the elevation angle φ ranges from $-\pi/2$ to $\pi/2$, but it is sufficient to assume that φ ranges from 0 to $\pi/2$, because the steering vector contains $\cos \varphi$ only. Hence, we initialize each $\hat{\varphi}_l$ uniformly from $[0, \pi/2]$. To reduce the computational complexity, we use a fixed stepsize to update $\hat{\varphi}$ [similarly to (41)]:

$$\hat{\varphi}^{(i+1)} = \hat{\varphi}^{(i)} + \frac{\pi}{36} \cdot \max\{(\rho)^i, 0.001\} \cdot \text{sign}(\zeta_{\beta}^{(i)}), \quad (57)$$

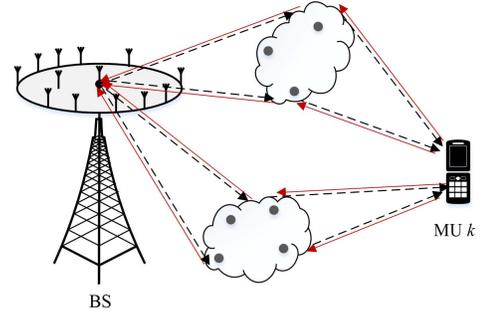


Fig. 4. Illustration of angular reciprocity for a massive MIMO system, where the downlink transmission and the uplink transmission are denoted by the dotted black lines and red lines, respectively.

where $0 < \rho < 1$ is a constant.⁴ Here, we use a different stepsize, $\frac{\pi}{36} \cdot \max\{(\rho)^i, 0.001\}$ instead of the $r_{\theta}/100$ in (41), because there is no grid covering the elevation angle domain. The motivation for choosing such a stepsize comes from the fact that 1) the term $\frac{\pi}{36}$ guarantees that the true elevation angles can be approximately approached within 20 iterations; and 2) the term $(\rho)^i$ keeps the stepsize decreasing, so as to attain a sufficiently small value [which is no less than $\frac{\pi}{36} \cdot 0.001$ due to the constant term 0.001 in (57)].

V. CHANNEL ESTIMATION WITH ANGULAR RECIPROCITY

For the downlink channel estimation, the training period T could become a bottleneck of the recovery performance, because the dimension of the measurement vector \mathbf{y} in (17) is determined by the training period T , while T is usually much less than N . The performance of the downlink channel estimation can be improved if we collect more useful information. Inspired by the angular reciprocity used in [22], we present an off-grid uplink-AoA-aided channel estimation method in this section. Here we only take the linear array as an example, but its extension to an arbitrary 2D-array antenna geometry is straightforward. Note that angular reciprocity is quite different from the commonly used channel reciprocity in TDD systems. In the first subsection, we will explain angular reciprocity in detail.

A. Angular Reciprocity

Following the downlink channel model in Section II-A, the uplink channel vector from the k -th user to the BS is given by

$$\bar{\mathbf{h}}_k = \sum_{c=1}^{N_c} \sum_{s=1}^{N_s} \bar{\xi}_{c,s}^k \bar{\mathbf{a}}(\bar{\theta}_{c,s}^k), \quad (58)$$

where $\bar{\xi}_{c,s}^k$ is similarly defined as $\xi_{c,s}^k$, $\bar{\theta}_{c,s}^k$ is the corresponding azimuth angle-of-arrival (AoA), as illustrated in Fig. 4, and $\bar{\mathbf{a}}(\theta) \in \mathbb{C}^{N \times 1}$ is the steering vector for a linear array:

$$\bar{\mathbf{a}}(\theta) = [1, e^{-j2\pi \frac{d_2}{\lambda_u} \sin(\theta)}, \dots, e^{-j2\pi \frac{d_N}{\lambda_u} \sin(\theta)}]^T,$$

⁴The maximum movement is about $\frac{\pi}{36} \sum_{i=1}^{\infty} (\rho)^i$. In order to cover the whole angle domain $[0, \pi/2]$, ρ should be chosen to be from 0.9474 to 1.

where λ_u is the wavelength of uplink propagation. Usually, channel reciprocity does not hold in FDD systems because different frequency bands are used in the downlink and uplink transmission. However, if the downlink and uplink transmissions operate closely in time, it is reasonable to have the following assumption:

Assumption 6 (Angular Reciprocity [22]). The azimuth AoAs of signals for the k -th MU in the uplink transmission almost coincide with the azimuth AoDs of signals in the downlink transmission, i.e.,

$$\theta_{c,s}^k = \bar{\theta}_{c,s}^k, \quad \forall k, c, s, \quad (59)$$

as illustrated in Fig. 4.

To exploit the angular reciprocity, Ding and Rao [22] collected the downlink and uplink channel vectors for the k -th MU in pair

$$\mathbf{h}_k = \mathbf{F}\mathbf{t}_k, \quad (60)$$

$$\bar{\mathbf{h}}_k = \mathbf{F}\bar{\mathbf{t}}_k, \quad (61)$$

where $\bar{\mathbf{t}}_k$ is the sparse representation of $\bar{\mathbf{h}}_k$ under the DFT basis for ULAs. It is worth noting that the steering vectors $\mathbf{a}(\theta_{c,s}^k)$ and $\bar{\mathbf{a}}(\bar{\theta}_{c,s}^k)$ are distinct if different frequency bands are used. Hence, the angular reciprocity between the downlink and uplink transmissions does not bring a joint sparse structure for \mathbf{t}_k and $\bar{\mathbf{t}}_k$. To get around this problem, they assume that the frequency duplex distance is not large (i.e., $\lambda_d \approx \lambda_u$). In this case, they approximately have

$$\mathbf{a}(\theta_{c,s}^k) \approx \bar{\mathbf{a}}(\bar{\theta}_{c,s}^k), \quad (62)$$

and then

$$\text{supp}(\mathbf{t}_k) \approx \text{supp}(\bar{\mathbf{t}}_k). \quad (63)$$

As the joint sparse structure only holds approximately, it may result in performance loss. To handle this drawback, we will propose a joint off-grid model in the next subsection.

B. Joint Off-Grid Model

For the uplink channel estimation, assume that each MU broadcasts a sequence of \bar{T} training pilot symbols, denoted by $\mathbf{s}_k \in \mathbb{C}^{\bar{T} \times 1}, k = 1, 2, \dots, K$. Then, the received signal $\bar{\mathbf{Y}} \in \mathbb{C}^{N \times \bar{T}}$ at the BS is given by

$$\bar{\mathbf{Y}} = \bar{\mathbf{H}}\mathbf{S} + \bar{\mathbf{N}}, \quad (64)$$

where $\bar{\mathbf{H}} = [\bar{\mathbf{h}}_1, \bar{\mathbf{h}}_2, \dots, \bar{\mathbf{h}}_K] \in \mathbb{C}^{N \times K}$ with $\bar{\mathbf{h}}_k$ being the channel vector for the k -th MU, $\mathbf{S} = [\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_K]^T \in \mathbb{C}^{K \times \bar{T}}$, and $\bar{\mathbf{N}} \in \mathbb{C}^{N \times \bar{T}}$ stands for the additive complex Gaussian noise with each element being zero mean and variance $\bar{\sigma}^2$ in the uplink. If the number of MUs is small, i.e., $K \leq \bar{T}$, the uplink channel matrix $\bar{\mathbf{H}}_u$ can be easily obtained by the conventional LS estimate, i.e.,

$$[\bar{\mathbf{h}}_1^{ls}, \bar{\mathbf{h}}_2^{ls}, \dots, \bar{\mathbf{h}}_K^{ls}] \triangleq \bar{\mathbf{Y}}\mathbf{S}^\dagger = \bar{\mathbf{H}} + \mathbf{E}, \quad (65)$$

or, equivalently,

$$\bar{\mathbf{h}}_k^{ls} = \bar{\mathbf{h}}_k + \mathbf{e}_k, \quad k = 1, 2, \dots, K, \quad (66)$$

where $\bar{\mathbf{h}}_k^{ls}$ stands for the LS estimate of $\bar{\mathbf{h}}_k$ and $\mathbf{E} \triangleq [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_K]$ stands for the estimation error. If \mathbf{S} consists of

an orthogonal pilot sequence, \mathbf{E} is i.i.d. Gaussian. Compared with the requirement $T \geq N$ for the downlink channel estimation, it is much easier to meet the requirement $\bar{T} \geq K$ for the uplink channel estimation.

With (66) and (17), we are able to exploit the sparse property of each MU independently. We drop the MU's index k for ease of notation, and then the paired sparse representation equalities can be rewritten as

$$\mathbf{y} = \Phi(\beta)\mathbf{w} + \mathbf{n}, \quad (67)$$

$$\bar{\mathbf{h}}^{ls} = \bar{\Phi}(\beta)\bar{\mathbf{w}} + \mathbf{e}, \quad (68)$$

where (67) coincides with (17), $\bar{\Phi}(\beta) = [\bar{\mathbf{a}}(\hat{\vartheta}_1 + \beta_1), \bar{\mathbf{a}}(\hat{\vartheta}_2 + \beta_2), \dots, \bar{\mathbf{a}}(\hat{\vartheta}_{\bar{L}} + \beta_{\bar{L}})]$, and $\bar{\mathbf{w}}$ is the sparse representation of $\bar{\mathbf{h}}^{ls}$. If the angular reciprocity holds, it is easy to check that $\text{supp}(\mathbf{w}) = \text{supp}(\bar{\mathbf{w}})$. Different from the approximation method [22] that hinges on the condition of (62), our off-grid model guarantees a jointly sparse structure from the angular domain directly, where neither approximation of $\lambda_d \approx \lambda_u$ nor the assumption of ULA at the BS is required. In the following, we will show how to jointly recover the sparse vectors \mathbf{w} and $\bar{\mathbf{w}}$ in the framework of SBL with the in-exact MM algorithm. Since the results can be similarly derived by following the procedures in Section III, detailed derivations are omitted for brevity.

C. Sparse Bayesian Learning Formulation

Under the assumption of circular symmetric complex Gaussian noises, we have

$$p(\mathbf{y}|\mathbf{w}, \alpha, \beta) = \mathcal{CN}(\mathbf{y}|\Phi(\beta)\mathbf{w}, \alpha^{-1}\mathbf{I}), \quad (69)$$

$$p(\bar{\mathbf{h}}^{ls}|\bar{\mathbf{w}}, \bar{\alpha}, \beta) = \mathcal{CN}(\bar{\mathbf{h}}^{ls}|\bar{\Phi}(\beta)\bar{\mathbf{w}}, \bar{\alpha}^{-1}\mathbf{I}), \quad (70)$$

where $\bar{\alpha}$ stands for the noise precision of \mathbf{e} , which is further modeled as a Gamma hyperprior $p(\bar{\alpha}) = \Gamma(\bar{\alpha}; a, b)$. Recall that, in Section III-B, we have used γ to control the sparsity of \mathbf{w} as follows:

$$p(\mathbf{w}|\gamma) = \mathcal{CN}(\mathbf{w}|\mathbf{0}, \text{diag}(\gamma^{-1})). \quad (71)$$

If we let $\boldsymbol{\tau} = [\tau_1, \tau_2, \dots, \tau_{\bar{L}}]^T$ be a nonnegative vector and

$$p(\bar{\mathbf{w}}|\gamma) = \mathcal{CN}(\bar{\mathbf{w}}|\mathbf{0}, \text{diag}((\gamma \odot \boldsymbol{\tau})^{-1})), \quad (72)$$

\mathbf{w} and $\bar{\mathbf{w}}$ will share a joint sparse structure. For example, γ_l^{-1} tends to zero, so is $(\gamma_l \tau_l)^{-1}$. Therefore, (71) and (72) provide a mathematic representation of the angular reciprocity. The estimated uplink channel $\bar{\mathbf{h}}^{ls}$ contains the AoA information in the uplink. The proposed method only exploits the azimuth AoA information in $\bar{\mathbf{h}}^{ls}$ to help in identifying the azimuth AoD in the downlink (downlink channel support) via the angular reciprocity in (71) and (72). The small scale fading information contained in the estimated uplink channel $\bar{\mathbf{h}}^{ls}$ cannot be exploited since the channel reciprocity does not hold for FDD systems.

D. Bayesian Inference

For ease of notation, let $\mathbf{w}^a \triangleq \{\mathbf{w}, \bar{\mathbf{w}}\}$, $\mathbf{y}^a \triangleq \{\mathbf{y}, \bar{\mathbf{h}}^{ls}\}$, $\Theta \triangleq \{\alpha, \bar{\alpha}, \gamma, \tau, \beta\}$. We assume a noninformative uniform prior for τ . At some fixed point $\Theta = \{\hat{\alpha}, \hat{\bar{\alpha}}, \hat{\gamma}, \hat{\tau}, \hat{\beta}\}$, we construct the surrogate function as

$$\mathcal{U}(\Theta|\hat{\Theta}) = \int p(\mathbf{w}^a|\mathbf{y}^a, \hat{\Theta}) \ln \frac{p(\mathbf{w}^a, \mathbf{y}^a, \Theta)}{p(\mathbf{w}^a|\mathbf{y}^a, \hat{\Theta})} d\mathbf{w}^a, \quad (73)$$

where

$$p(\mathbf{w}^a, \mathbf{y}^a, \Theta) = p(\mathbf{y}|\mathbf{w}, \alpha, \beta) p(\mathbf{w}|\gamma) p(\bar{\mathbf{h}}^{ls}|\bar{\mathbf{w}}, \bar{\alpha}, \beta) \\ \cdot p(\bar{\mathbf{w}}|\gamma, \tau) p(\alpha) p(\bar{\alpha}) p(\gamma) p(\tau) p(\beta)$$

and

$$p(\mathbf{w}^a|\mathbf{y}^a, \Theta) = \mathcal{CN}(\mathbf{w}|\boldsymbol{\mu}(\alpha, \gamma, \beta), \boldsymbol{\Sigma}(\alpha, \gamma, \beta)) \\ \cdot \mathcal{CN}(\bar{\mathbf{w}}|\bar{\boldsymbol{\mu}}(\bar{\alpha}, \gamma, \tau, \beta), \bar{\boldsymbol{\Sigma}}(\bar{\alpha}, \gamma, \tau, \beta)),$$

with

$$\bar{\boldsymbol{\mu}}(\bar{\alpha}, \gamma, \tau, \beta) = \bar{\alpha} \bar{\boldsymbol{\Sigma}}(\bar{\alpha}, \gamma, \tau, \beta) \bar{\boldsymbol{\Phi}}^H(\beta) \bar{\mathbf{h}}^{ls}, \\ \bar{\boldsymbol{\Sigma}}(\bar{\alpha}, \gamma, \tau, \beta) = (\bar{\alpha} \bar{\boldsymbol{\Phi}}^H(\beta) \bar{\boldsymbol{\Phi}}(\beta) + \text{diag}(\gamma \odot \tau))^{-1}.$$

Note that $\boldsymbol{\mu}(\alpha, \gamma, \beta)$ and $\boldsymbol{\Sigma}(\alpha, \gamma, \beta)$ have been defined in (35).

In the maximization step of the $(i+1)$ -th iteration, we update $\alpha, \bar{\alpha}, \gamma, \tau, \beta$ as

$$\alpha^{(i+1)} = \arg \max_{\alpha} \mathcal{U}(\alpha, \bar{\alpha}^{(i)}, \gamma^{(i)}, \tau^{(i)}, \beta^{(i)} | \Theta^{(i)}), \quad (74)$$

$$\bar{\alpha}^{(i+1)} = \arg \max_{\bar{\alpha}} \mathcal{U}(\alpha^{(i+1)}, \bar{\alpha}, \gamma^{(i)}, \tau^{(i)}, \beta^{(i)} | \Theta_1^{(i)}), \quad (75)$$

$$\gamma^{(i+1)} = \arg \max_{\gamma} \mathcal{U}(\alpha^{(i+1)}, \bar{\alpha}^{(i+1)}, \gamma, \tau^{(i)}, \beta^{(i)} | \Theta_2^{(i)}), \quad (76)$$

$$\tau^{(i+1)} = \arg \max_{\tau} \mathcal{U}(\alpha^{(i+1)}, \bar{\alpha}^{(i+1)}, \gamma^{(i+1)}, \tau, \beta^{(i)} | \Theta_3^{(i)}), \quad (77)$$

$$\beta^{(i+1)} = \arg \max_{\beta} \mathcal{U}(\alpha^{(i+1)}, \bar{\alpha}^{(i+1)}, \gamma^{(i+1)}, \tau^{(i+1)}, \beta | \Theta_4^{(i)}). \quad (78)$$

Extending Lemmas 3 and 4, the updates for $\alpha, \bar{\alpha}, \gamma$ and τ can be obtained as follows:

$$\alpha^{(i+1)} = \frac{T + a}{b + \eta_d(\Theta^{(i)})}, \quad (79)$$

$$\bar{\alpha}^{(i+1)} = \frac{N + a}{b + \eta_u(\Theta_1^{(i)})}, \quad (80)$$

$$\gamma_l^{(i+1)} = \frac{a + 2}{b + \left[\Xi_d(\Theta_2^{(i)}) + \tau_l \Xi_u(\Theta_2^{(i)}) \right]_l}, \quad \forall l, \quad (81)$$

$$\tau_l^{(i+1)} = \frac{1}{\left[\gamma_l^{(i+1)} \Xi_u(\Theta_3^{(i)}) \right]_l}, \quad \forall l, \quad (82)$$

where

$$\eta_d(\Theta) = \text{tr}(\boldsymbol{\Phi}(\beta) \boldsymbol{\Sigma}(\alpha, \gamma, \beta) \boldsymbol{\Phi}^H(\beta)) \\ + \|\mathbf{y} - \boldsymbol{\Phi}(\beta) \boldsymbol{\mu}(\alpha, \gamma, \beta)\|_2^2, \\ \eta_u(\Theta) = \text{tr}(\bar{\boldsymbol{\Phi}}(\beta) \bar{\boldsymbol{\Sigma}}(\bar{\alpha}, \gamma, \tau, \beta) \bar{\boldsymbol{\Phi}}^H(\beta)) \\ + \|\bar{\mathbf{h}}^{ls} - \bar{\boldsymbol{\Phi}}(\beta) \bar{\boldsymbol{\mu}}(\bar{\alpha}, \gamma, \tau, \beta)\|_2^2, \\ \Xi_d(\Theta) = \boldsymbol{\Sigma}(\alpha, \gamma, \beta) + \boldsymbol{\mu}(\alpha, \gamma, \beta) \boldsymbol{\mu}^H(\alpha, \gamma, \beta), \\ \Xi_u(\Theta) = \bar{\boldsymbol{\Sigma}}(\bar{\alpha}, \gamma, \tau, \beta) + \bar{\boldsymbol{\mu}}(\bar{\alpha}, \gamma, \tau, \beta) \bar{\boldsymbol{\mu}}^H(\bar{\alpha}, \gamma, \tau, \beta).$$

Finally, we discuss how to refine the grid for the uplink-AoA-aided channel estimation. Ignoring the terms independent of β , the objective function in (78) becomes

$$\mathcal{U}(\alpha^{(i+1)}, \bar{\alpha}^{(i+1)}, \gamma^{(i+1)}, \tau^{(i+1)}, \beta | \Theta_4^{(i)}) \\ = -\alpha^{(i+1)} \left\| \mathbf{y} - \boldsymbol{\Phi}(\beta) \boldsymbol{\mu}^{(i)} \right\|_2^2 \\ - \alpha^{(i+1)} \text{tr} \left(\boldsymbol{\Phi}(\beta) \boldsymbol{\Sigma}^{(i)} \boldsymbol{\Phi}^H(\beta) \right) \\ - \bar{\alpha}^{(i+1)} \left\| \bar{\mathbf{h}}^{ls} - \bar{\boldsymbol{\Phi}}(\beta) \bar{\boldsymbol{\mu}}^{(i)} \right\|_2^2 \\ - \bar{\alpha}^{(i+1)} \text{tr} \left(\bar{\boldsymbol{\Phi}}(\beta) \bar{\boldsymbol{\Sigma}}^{(i)} \bar{\boldsymbol{\Phi}}^H(\beta) \right), \quad (83)$$

where $\boldsymbol{\mu}^{(i)} \triangleq \boldsymbol{\mu}(\alpha^{(i+1)}, \gamma^{(i+1)}, \beta^{(i)})$, $\boldsymbol{\Sigma}^{(i)} \triangleq \boldsymbol{\Sigma}(\alpha^{(i+1)}, \gamma^{(i+1)}, \beta^{(i)})$, $\bar{\boldsymbol{\mu}}^{(i)} \triangleq \bar{\boldsymbol{\mu}}(\bar{\alpha}^{(i+1)}, \gamma^{(i+1)}, \tau^{(i+1)}, \beta^{(i)})$, and $\bar{\boldsymbol{\Sigma}}^{(i)} \triangleq \bar{\boldsymbol{\Sigma}}(\bar{\alpha}^{(i+1)}, \gamma^{(i+1)}, \tau^{(i+1)}, \beta^{(i)})$. Calculating the derivative of (83) w.r.t. β_l leads to

$$\zeta^{(i)}(\beta_l) = 2\text{Re} \left((\mathbf{a}'(\hat{\vartheta}_l + \beta_l))^H \mathbf{X}^H \mathbf{X} (\mathbf{a}(\hat{\vartheta}_l + \beta_l)) \right) \cdot c_{d1}^{(i)} \\ + 2\text{Re} \left((\mathbf{a}'(\hat{\vartheta}_l + \beta_l))^H \mathbf{X}^H \mathbf{c}_{d2}^{(i)} \right) \\ + 2\text{Re} \left((\bar{\mathbf{a}}'(\hat{\vartheta}_l + \beta_l))^H (\bar{\mathbf{a}}(\hat{\vartheta}_l + \beta_l)) \right) \cdot c_{u1}^{(i)} \\ + 2\text{Re} \left((\bar{\mathbf{a}}'(\hat{\vartheta}_l + \beta_l))^H \mathbf{c}_{u2}^{(i)} \right), \quad (84)$$

where $\bar{\mathbf{a}}'(\hat{\vartheta}_j + \beta_l) = d\bar{\mathbf{a}}(\hat{\vartheta}_j + \beta_l)/d\beta_l$, $c_{d1}^{(i)} = -\alpha^{(i+1)} (\chi_{d,u}^{(i)} + |\mu_{d,l}^{(i)}|^2)$, $\mathbf{c}_{d2}^{(i)} = \alpha^{(i+1)} ((\mu_{d,l}^{(i)})^* \mathbf{y}_{d-l}^{(i)} - \mathbf{X} \sum_{j \neq l} \chi_{d,jl}^{(i)} (\mathbf{a}(\hat{\vartheta}_j + \beta_j)))$, $c_{u1}^{(i)} = -\bar{\alpha}^{(i+1)} (\chi_{u,u}^{(i)} + |\mu_{u,l}^{(i)}|^2)$, $\mathbf{c}_{u2}^{(i)} = \bar{\alpha}^{(i+1)} ((\mu_{u,l}^{(i)})^* \bar{\mathbf{h}}_{-l}^{(i)} - \sum_{j \neq l} \chi_{u,jl}^{(i)} (\bar{\mathbf{a}}(\hat{\vartheta}_j + \beta_j)))$, $\mathbf{y}_{d-l}^{(i)} = \mathbf{y} - \mathbf{X} \cdot \sum_{j \neq l} (\mu_{d,j}^{(i)} \cdot \mathbf{a}(\hat{\vartheta}_j + \beta_j))$, $\bar{\mathbf{h}}_{-l}^{(i)} = \bar{\mathbf{h}}^{ls} - \sum_{j \neq l} (\mu_{u,j}^{(i)} \cdot \bar{\mathbf{a}}(\hat{\vartheta}_j + \beta_j))$, $\mu_{d,l}^{(i)}$ ($\mu_{u,l}^{(i)}$) and $\chi_{d,jl}^{(i)}$ ($\chi_{u,jl}^{(i)}$) denote the l -th element and the (j, l) -th element of $\boldsymbol{\mu}^{(i)}$ ($\bar{\boldsymbol{\mu}}^{(i)}$) and $\boldsymbol{\Sigma}^{(i)}$ ($\bar{\boldsymbol{\Sigma}}^{(i)}$), respectively. With (84), we are able to update β similarly as in (40) or (41). Noth that following the same initializations mentioned in Section III-D, we set $\alpha^{(0)} = \bar{\alpha}^{(0)} = 1$, $\gamma^{(0)} = \tau^{(0)} = 1$ and $\beta^{(0)} = \mathbf{0}$ for the off-grid uplink-AoA-aided method.

VI. SIMULATION RESULTS

In this section, we conduct simulations to investigate the performance of our proposed methods. The proposed methods are compared with the following methods:

- **Baseline 1** (SBL): \mathbf{h}_k is recovered using the standard SBL method [23] with the dictionary \mathbf{A} defined in (14).
- **Baseline 2** (DFT): \mathbf{h}_k is recovered using the l_1 -norm minimization algorithm [31]–[33] with a DFT basis.
- **Baseline 3** (Overcomplete DFT): \mathbf{h}_k is recovered using the l_1 -norm minimization algorithm [31]–[33] with the dictionary \mathbf{A} defined in (14).
- **Baseline 4** (Dictionary Learning): \mathbf{h}_k is recovered using the method proposed in [22] with the dictionary \mathbf{A} defined in (14).

Since the state-of-the-art DFT methods work for ULAs only, we first focus on simulations for ULAs, where we use the 3GPP spatial channel model (SCM) [36] to generate the channel coefficients for an urban microcell. The uplink frequency

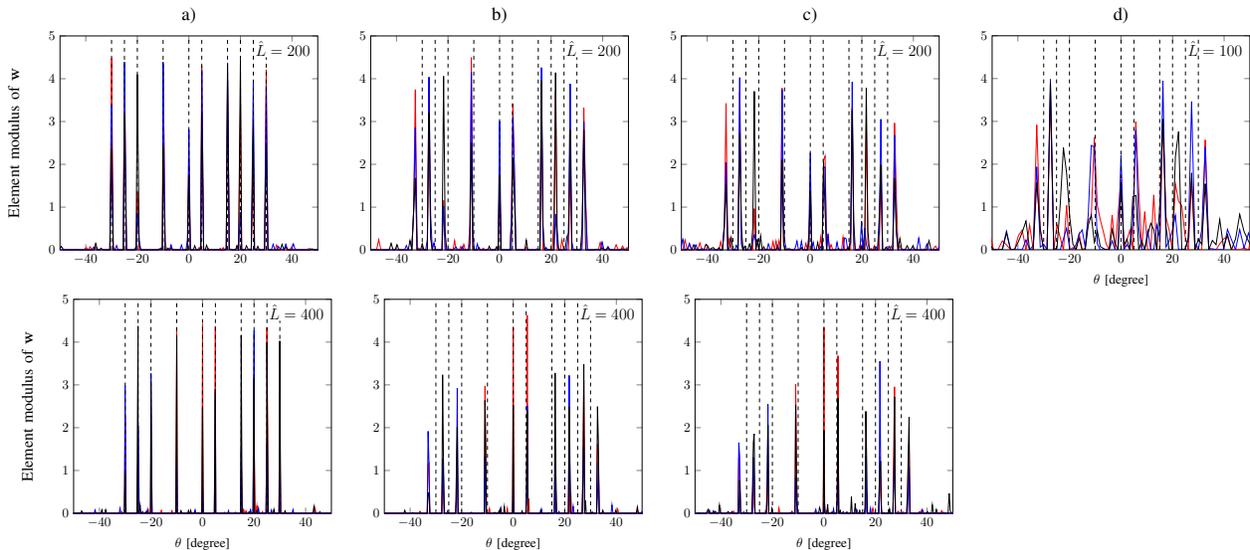


Fig. 5. Element modulus of \mathbf{w} for three independent trials with $N = 100$, $T = 40$ and $\text{SNR} = 10$ dB. The true azimuth AoDs are denoted by dotted lines. a) Off-Grid; b) SBL; c) Overcomplete DFT; d) DFT.

is 1980 MHz, the downlink frequency is 2170 MHz, and the inter-antenna spacing is $d = c/(2f_0)$, with c being the light speed and $f_0 = 2000$ MHz. Then, we run simulations with the 3GPP 3D channel model [41], which provides a 2D array model involving both azimuth and elevation angles. All the parameters of the 3D channel model follow 3D-UMa-NOLS (see Table 7.3-6 in [41]) and the downlink frequency is 2170 MHz. The normalized mean square error (NMSE) is defined as

$$\frac{1}{M_c} \sum_{m=1}^{M_c} \frac{\|\mathbf{h}_m^e - \mathbf{h}_m\|_2^2}{\|\mathbf{h}_m\|_2^2}, \quad (85)$$

where \mathbf{h}_m^e is the estimate of \mathbf{h}_m at the n -th Monte Carlo trial and $M_c = 200$ is the number of Monte Carlo trials.

A. Recovered Channel Sparsity in the Angular Domain for ULA

In Fig. 5, we illustrate the effect of direction mismatch on the channel sparse representation performance for different channel estimation strategies. Consider a simple scenario where a ULA with 100 antennas at the BS is used to send the training pilot symbols with ten azimuth AoDs in total, which are simply denoted as $\theta_1 = -30^\circ$, $\theta_2 = -25^\circ$, $\theta_3 = -20^\circ$, $\theta_4 = -10^\circ$, $\theta_5 = 0^\circ$, $\theta_6 = 5^\circ$, $\theta_7 = 15^\circ$, $\theta_8 = 20^\circ$, $\theta_9 = 25^\circ$, and $\theta_{10} = 30^\circ$. The training pilots are randomly generated with $T = 40$ and the SNR is set to 10 dB. Fig. 5 shows the element modulus of the recovered channel sparse representation \mathbf{w} , where the number of grid points \hat{L} is fixed to 200 or 400 for all the methods, except for the classical DFT method. It is observed that 1) the solution of the classical DFT method is not exactly sparse, and it has a significant performance loss due to the leakage of energy over many bins; 2) the standard SBL method and the overcomplete DFT method can achieve better sparse representations, especially for a dense grid ($\hat{L} = 400$), but direction mismatch always exists; and 3) our proposed off-grid method can greatly improve the sparsity and accuracy of

the channel representation, and the direction mismatch can be almost eliminated.

B. Channel Estimation Performance Versus T for ULA

In Fig. 6, Monte Carlo trials are carried out to investigate the impact of the number of pilot symbols on the downlink channel estimation performance for ULA. Assume that the ULA at the BS is equipped with 150 antennas, and the system supports ten MUs, where each MU has a single antenna. All the results are obtained by averaging over 200 Monte Carlo channel realizations. Every channel realization consists of $N_c = 3$ random scattering clusters ranging from -40° to 40° , and each cluster contains $N_s = 10$ sub-paths concentrated in a 20° angular spread. The training pilots are randomly generated, the SNR is chosen as 0 dB or 10 dB, and the number of grid points is fixed to 200 for all but the DFT method. Fig. 6 shows the NMSE performance of the downlink channel estimate achieved by the different channel estimation strategies versus the number of training pilot symbols T . It can be seen that 1) the NMSEs of all the methods decrease as the number of training pilot symbols increases, and the DFT method gives the worst performance; 2) compared with the DFT method, the state-of-the-art methods (the overcomplete DFT and dictionary learning method) can improve the NMSE performance, but the improvement is not significant (they are all worse than the standard SBL method); and 3) our proposed off-grid method always outperforms the state-of-the-art methods, and the uplink-AoA-aided method can further improve the performance of the downlink channel estimation, as it collects more useful information than the off-grid method.

C. Channel Estimation Performance Versus \hat{L} for ULA

In Fig. 7, we study the impact of the number of grid points on the downlink channel estimation performance for ULA. We consider the same scenario as in Section VI-B, except that

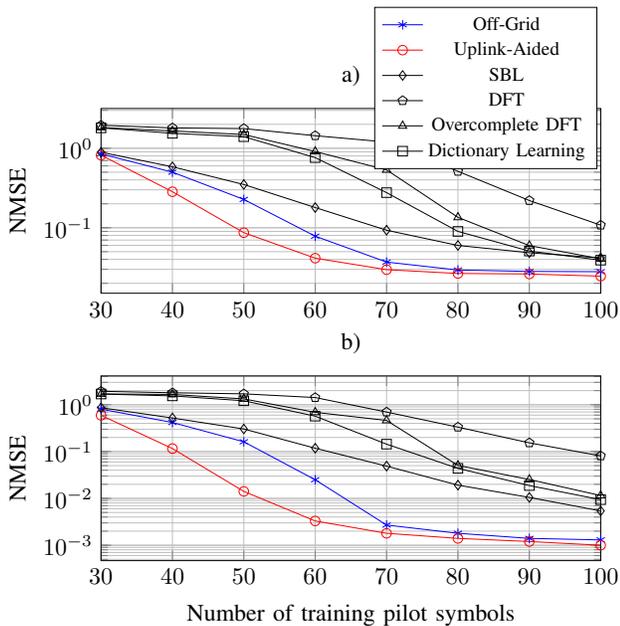


Fig. 6. NMSE of downlink channel estimate versus the number of training pilot symbols for ULA. a) SNR = 0 dB; b) SNR = 10 dB.

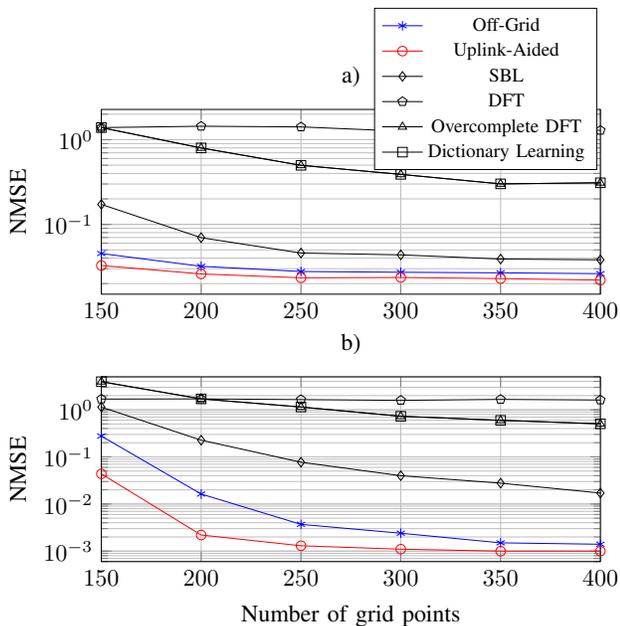


Fig. 7. NMSE of downlink channel estimate versus the number of grid points for ULA. a) $N = 150$, $N_c = 2$ and SNR = 0 dB; b) $N = 200$, $N_c = 3$ and SNR = 10 dB.

the number of training pilot symbols is fixed to 70, and the scattering clusters range from -90° to 90° . All the results are obtained by averaging over 200 Monte Carlo channel realizations. Fig. 7 shows the NMSE performance of the downlink channel estimate achieved by the different channel estimation strategies versus the number of grid points \hat{L} . It is shown that the overcomplete DFT method and dictionary learning method achieve the same performance, because there is no benefit in learning the true AoDs which range from -90° to 90° . The NMSEs of the DFT method, overcomplete DFT

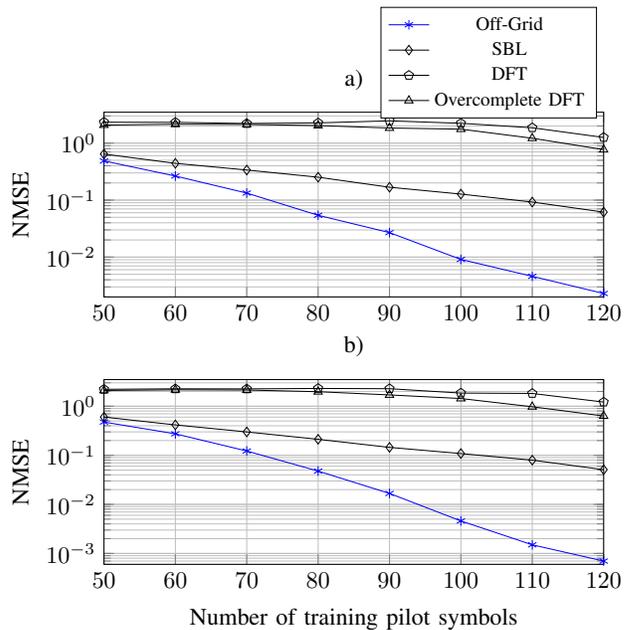


Fig. 8. NMSE of downlink channel estimate versus the number of training pilot symbols for 2D array. a) $\hat{L} = 250$ and SNR = 0 dB; b) $\hat{L} = 300$ and SNR = 10 dB.

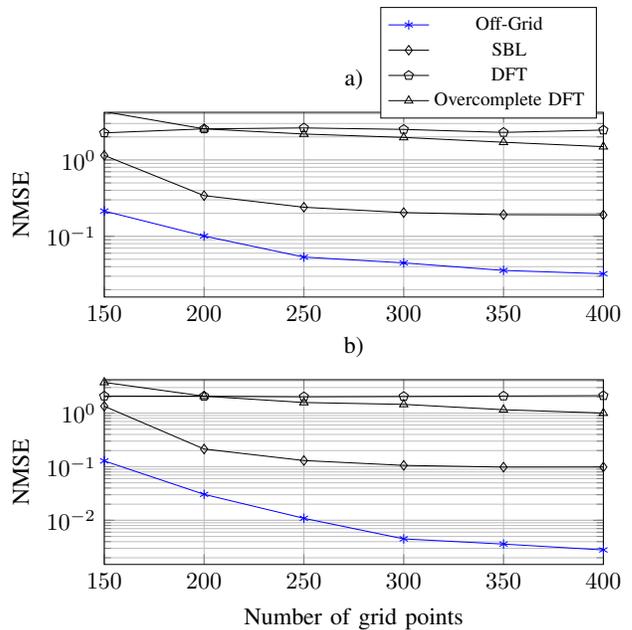


Fig. 9. NMSE of downlink channel estimate versus the number of grid points for 2D array with SNR = 0 dB. a) $T = 80$; b) $T = 100$.

method and SBL method coincide with each other at $\hat{L} = 150$ in Fig. 7-a and $\hat{L} = 200$ in Fig. 7-b, respectively, because they use the same grid in the case of $N = \hat{L}$. The NMSEs of our methods decrease as the number of grid points increases, and they always outperform the others, no matter what number of grid points is used.

D. Channel Estimation Performance with 2D Array

In Figs. 8 and 9, Monte Carlo trials are carried out to investigate the channel estimation performance with the 2D

array. Assume that the 2D planar array at the BS is equipped with 20×10 antennas, where both the horizontal and vertical inter-antenna spacings are a half wavelength. Every channel realization consists of $N_c = 20$ random scattering clusters, and each cluster contains $N_s = 20$ subpaths. The AoDs are randomly generated by the 3GPP 3D channel model, where the azimuth AoDs range from -180° to 180° and the elevation AoDs range from -90° to 90° . The training pilots are randomly generated, the SNR is chosen as 0 dB or 10 dB, and ρ in (57) is set to 0.95. All the results are obtained by averaging over 200 Monte Carlo channel realizations. Fig. 8 shows the NMSE performance of the downlink channel estimate achieved by the different channel estimation strategies versus the number of training pilot symbols T , and Fig. 9 shows the NMSE performance of the downlink channel estimate achieved by the different channel estimation strategies versus the number of grid points \hat{L} . It can be seen that 1) the DFT method and overcomplete DFT method give very poor performance, because they can not work for non-ULAs; 2) the standard SBL method outperforms the DFT-based methods, but the performance improvement is not significant; and 3) our proposed off-grid method indeed works for the 2D array, and it can substantially improve the channel estimation performance.

VII. CONCLUSION

The problem of downlink channel estimation in FDD massive MIMO systems is addressed in this paper. We provide a novel off-grid model for massive MIMO channel sparse representation, which can greatly improve the sparsity and accuracy of the channel representation. To the best of our knowledge, our work is the first to utilize an off-grid channel model to combat modeling error for channel estimation. The proposed off-grid model and the SBL-based framework have wide applicability. They do not require any prior knowledge about the sparsity of channels, nor the variance of noises, and all the parameters are automatically tuned by the in-exact MM algorithm. Extending the results to MUs with multiple antennas is straightforward in the framework of SBL.

APPENDIX

A. Proof of Lemma 1

The non-decreasing property can be achieved as

$$\begin{aligned} & \ln p(\mathbf{y}, \alpha^{(i+1)}, \gamma^{(i+1)}, \beta^{(i+1)}) \\ & \geq \mathcal{U}(\alpha^{(i+1)}, \gamma^{(i+1)}, \beta^{(i+1)} | \alpha^{(i+1)}, \gamma^{(i+1)}, \beta^{(i)}) \end{aligned} \quad (86)$$

$$\geq \mathcal{U}(\alpha^{(i+1)}, \gamma^{(i+1)}, \beta^{(i)} | \alpha^{(i+1)}, \gamma^{(i+1)}, \beta^{(i)}) \quad (87)$$

$$= \ln p(\mathbf{y}, \alpha^{(i+1)}, \gamma^{(i+1)}, \beta^{(i)}) \quad (88)$$

$$\geq \mathcal{U}(\alpha^{(i+1)}, \gamma^{(i+1)}, \beta^{(i)} | \alpha^{(i+1)}, \gamma^{(i)}, \beta^{(i)}) \quad (89)$$

$$\geq \mathcal{U}(\alpha^{(i+1)}, \gamma^{(i)}, \beta^{(i)} | \alpha^{(i+1)}, \gamma^{(i)}, \beta^{(i)}) \quad (90)$$

$$= \ln p(\mathbf{y}, \alpha^{(i+1)}, \gamma^{(i)}, \beta^{(i)}) \quad (91)$$

$$\geq \mathcal{U}(\alpha^{(i+1)}, \gamma^{(i)}, \beta^{(i)} | \alpha^{(i)}, \gamma^{(i)}, \beta^{(i)}) \quad (92)$$

$$\geq \mathcal{U}(\alpha^{(i)}, \gamma^{(i)}, \beta^{(i)} | \alpha^{(i)}, \gamma^{(i)}, \beta^{(i)}) \quad (93)$$

$$= \ln p(\mathbf{y}, \alpha^{(i)}, \gamma^{(i)}, \beta^{(i)}), \quad (94)$$

where (86), (89) and (92) follow (26); (88), (91) and (94) follow (27); and (87), (90) and (93) follow (33), (32) and (31), respectively.

B. Proof of Lemma 2

Letting $q(\mathbf{w})$ be an arbitrary distribution, the lower bound of $\ln p(\mathbf{y}, \alpha, \gamma, \beta)$ can be written as

$$\begin{aligned} \ln p(\mathbf{y}, \alpha, \gamma, \beta) &= \ln \int p(\mathbf{w}, \mathbf{y}, \alpha, \gamma, \beta) d\mathbf{w} \\ &= \ln \int q(\mathbf{w}) \frac{p(\mathbf{w}, \mathbf{y}, \alpha, \gamma, \beta)}{q(\mathbf{w})} d\mathbf{w} \\ &\geq \int q(\mathbf{w}) \ln \frac{p(\mathbf{w}, \mathbf{y}, \alpha, \gamma, \beta)}{q(\mathbf{w})} d\mathbf{w}, \end{aligned} \quad (95)$$

where Jensen's inequality is applied in the last step. The equality holds when $\frac{p(\mathbf{w}, \mathbf{y}, \alpha, \gamma, \beta)}{q(\mathbf{w})} = c$ for a constant c that does not depend on \mathbf{w} . As $q(\mathbf{w})$ is a distribution, we have $\int q(\mathbf{w}) d\mathbf{w} = 1$. This further indicates that

$$c = \int p(\mathbf{w}, \mathbf{y}, \alpha, \gamma, \beta) d\mathbf{w} = p(\mathbf{y}, \alpha, \gamma, \beta) \quad (96)$$

and

$$q(\mathbf{w}) = p(\mathbf{w} | \mathbf{y}, \alpha, \gamma, \beta). \quad (97)$$

With (95) and (97), it is easy to check that the constructed surrogate function $\mathcal{U}(\alpha, \gamma, \beta | \hat{\alpha}, \hat{\gamma}, \hat{\beta})$ always satisfies (26) and (27) for any fixed $(\hat{\alpha}, \hat{\gamma}, \hat{\beta})$.

To prove (28), we first rewrite the left side of (28) as

$$\begin{aligned} & \left. \frac{\partial \mathcal{U}(\alpha, \hat{\gamma}, \hat{\beta} | \hat{\alpha}, \hat{\gamma}, \hat{\beta})}{\partial \alpha} \right|_{\alpha=\hat{\alpha}} \\ &= \int p(\mathbf{w} | \mathbf{y}, \hat{\alpha}, \hat{\gamma}, \hat{\beta}) \left. \frac{\partial \ln p(\mathbf{w}, \mathbf{y}, \alpha, \hat{\gamma}, \hat{\beta})}{\partial \alpha} d\mathbf{w} \right|_{\alpha=\hat{\alpha}} \\ &= \int \frac{p(\mathbf{w} | \mathbf{y}, \hat{\alpha}, \hat{\gamma}, \hat{\beta})}{p(\mathbf{w}, \mathbf{y}, \alpha, \hat{\gamma}, \hat{\beta})} \left. \frac{\partial p(\mathbf{w}, \mathbf{y}, \alpha, \hat{\gamma}, \hat{\beta})}{\partial \alpha} d\mathbf{w} \right|_{\alpha=\hat{\alpha}} \\ &= \frac{1}{p(\mathbf{y}, \hat{\alpha}, \hat{\gamma}, \hat{\beta})} \int \left. \frac{\partial p(\mathbf{w}, \mathbf{y}, \alpha, \hat{\gamma}, \hat{\beta})}{\partial \alpha} d\mathbf{w} \right|_{\alpha=\hat{\alpha}} \\ &= \frac{1}{p(\mathbf{y}, \hat{\alpha}, \hat{\gamma}, \hat{\beta})} \cdot \left. \frac{\partial p(\mathbf{y}, \alpha, \hat{\gamma}, \hat{\beta})}{\partial \alpha} \right|_{\alpha=\hat{\alpha}}. \end{aligned} \quad (98)$$

On the other hand, the right side of (28) is

$$\frac{\partial \ln p(\mathbf{y}, \alpha, \hat{\gamma}, \hat{\beta})}{\partial \alpha} = \frac{1}{p(\mathbf{y}, \alpha, \hat{\gamma}, \hat{\beta})} \frac{\partial p(\mathbf{y}, \alpha, \hat{\gamma}, \hat{\beta})}{\partial \alpha}. \quad (99)$$

Combining (98) and (99), we achieve the equality in (28). Since the proofs for (29)–(30) can be similarly achieved, they are omitted for brevity.

C. Proof of Lemma 3

For α , ignoring the independent terms, the objective function in (31) can be rewritten as

$$\begin{aligned}
& \mathcal{U}(\alpha, \boldsymbol{\gamma}^{(i)}, \boldsymbol{\beta}^{(i)} | \alpha^{(i)}, \boldsymbol{\gamma}^{(i)}, \boldsymbol{\beta}^{(i)}) \\
&= \int p(\mathbf{w} | \mathbf{y}, \alpha^{(i)}, \boldsymbol{\gamma}^{(i)}, \boldsymbol{\beta}^{(i)}) \ln p(\mathbf{y} | \mathbf{w}, \alpha, \boldsymbol{\beta}^{(i)}) d\mathbf{w} \\
&\quad + \int p(\mathbf{w} | \mathbf{y}, \alpha^{(i)}, \boldsymbol{\gamma}^{(i)}, \boldsymbol{\beta}^{(i)}) \ln p(\alpha) d\mathbf{w} \\
&= -\alpha \int p(\mathbf{w} | \mathbf{y}, \alpha^{(i)}, \boldsymbol{\gamma}^{(i)}, \boldsymbol{\beta}^{(i)}) \left\| \mathbf{y} - \boldsymbol{\Phi}(\boldsymbol{\beta}^{(i)}) \mathbf{w} \right\|_2^2 d\mathbf{w} \\
&\quad + T \ln \alpha + (a) \ln \alpha - b\alpha \\
&= (T + a) \ln \alpha - \alpha (b + \eta(\alpha^{(i)}, \boldsymbol{\gamma}^{(i)}, \boldsymbol{\beta}^{(i)})). \tag{100}
\end{aligned}$$

Since (100) is a strict concave function related to α , setting its derivative to zero gives the unique optimal solution

$$\alpha^{(i+1)} = \frac{T + a}{b + \eta(\alpha^{(i)}, \boldsymbol{\gamma}^{(i)}, \boldsymbol{\beta}^{(i)})}.$$

D. Proof of Lemma 4

For $\boldsymbol{\gamma}$, ignoring the independent terms of the objective function in (32), we obtain

$$\begin{aligned}
& \mathcal{U}(\alpha^{(i+1)}, \boldsymbol{\gamma}, \boldsymbol{\beta}^{(i)} | \alpha^{(i+1)}, \boldsymbol{\gamma}^{(i)}, \boldsymbol{\beta}^{(i)}) \\
&= \int p(\mathbf{w} | \mathbf{y}, \alpha^{(i+1)}, \boldsymbol{\gamma}^{(i)}, \boldsymbol{\beta}^{(i)}) \ln p(\mathbf{w} | \boldsymbol{\gamma}) d\mathbf{w} \\
&\quad + \int p(\mathbf{w} | \mathbf{y}, \alpha^{(i+1)}, \boldsymbol{\gamma}^{(i)}, \boldsymbol{\beta}^{(i)}) \ln p(\boldsymbol{\gamma}) d\mathbf{w} \\
&= -\ln |\text{diag}(\boldsymbol{\gamma}^{-1})| + (a) \sum_{l=1}^{\hat{L}} \ln \gamma_l - b \sum_{l=1}^{\hat{L}} \gamma_l \\
&\quad - \int p(\mathbf{w} | \mathbf{y}, \alpha^{(i+1)}, \boldsymbol{\gamma}^{(i)}, \boldsymbol{\beta}^{(i)}) ((\mathbf{w})^H \text{diag}(\boldsymbol{\gamma}) \mathbf{w}) d\mathbf{w} \\
&= \sum_{l=1}^{\hat{L}} \ln \gamma_l + (a) \sum_{l=1}^{\hat{L}} \ln \gamma_l - b \sum_{l=1}^{\hat{L}} \gamma_l \\
&\quad - \text{tr} \left(\boldsymbol{\Xi}(\alpha^{(i+1)}, \boldsymbol{\gamma}^{(i)}, \boldsymbol{\beta}^{(i)}) \cdot \text{diag}(\boldsymbol{\gamma}) \right).
\end{aligned}$$

Differentiating w.r.t. each γ_l yields

$$\begin{aligned}
& \frac{\partial \mathcal{U}(\alpha^{(i+1)}, \boldsymbol{\gamma}, \boldsymbol{\beta}^{(i)} | \alpha^{(i+1)}, \boldsymbol{\gamma}^{(i)}, \boldsymbol{\beta}^{(i)})}{\partial \gamma_l} \\
&= \frac{a + 1}{\gamma_l} - b - \left[\boldsymbol{\Xi}(\alpha^{(i+1)}, \boldsymbol{\gamma}^{(i)}, \boldsymbol{\beta}^{(i)}) \right]_{ll}.
\end{aligned}$$

Then, setting the derivative to zero and solving for γ_l give the unique optimal solution

$$\gamma_l^{(i+1)} = \frac{a + 1}{b + \left[\boldsymbol{\Xi}(\alpha^{(i+1)}, \boldsymbol{\gamma}^{(i)}, \boldsymbol{\beta}^{(i)}) \right]_{ll}}.$$

E. Derivation for Eq. (38)

Ignoring the independent terms, the objective function in (33) becomes

$$\begin{aligned}
& \mathcal{U}(\alpha^{(i+1)}, \boldsymbol{\gamma}^{(i+1)}, \boldsymbol{\beta} | \alpha^{(i+1)}, \boldsymbol{\gamma}^{(i+1)}, \boldsymbol{\beta}^{(i)}) \\
&= \int p(\mathbf{w} | \mathbf{y}, \alpha^{(i+1)}, \boldsymbol{\gamma}^{(i+1)}, \boldsymbol{\beta}^{(i)}) \ln p(\mathbf{y} | \mathbf{w}, \alpha^{(i+1)}, \boldsymbol{\beta}) d\mathbf{w} \\
&= -\alpha^{(i+1)} \int p(\mathbf{w} | \mathbf{y}, \alpha^{(i+1)}, \boldsymbol{\gamma}^{(i+1)}, \boldsymbol{\beta}^{(i)}) \\
&\quad \cdot \left\| \mathbf{y} - \boldsymbol{\Phi}(\boldsymbol{\beta}) \mathbf{w} \right\|_2^2 d\mathbf{w} \\
&= -\alpha^{(i+1)} \left\| \mathbf{y} - \boldsymbol{\Phi}(\boldsymbol{\beta}) \boldsymbol{\mu}(\alpha^{(i+1)}, \boldsymbol{\gamma}^{(i+1)}, \boldsymbol{\beta}^{(i)}) \right\|_2^2 \\
&\quad - \alpha^{(i+1)} \text{tr} \left(\boldsymbol{\Phi}(\boldsymbol{\beta}) \boldsymbol{\Sigma}(\alpha^{(i+1)}, \boldsymbol{\gamma}^{(i+1)}, \boldsymbol{\beta}^{(i)}) \boldsymbol{\Phi}^H(\boldsymbol{\beta}) \right).
\end{aligned}$$

For ease of notation, we simply denote $\boldsymbol{\mu}(\alpha^{(i+1)}, \boldsymbol{\gamma}^{(i+1)}, \boldsymbol{\beta}^{(i)})$ and $\boldsymbol{\Sigma}(\alpha^{(i+1)}, \boldsymbol{\gamma}^{(i+1)}, \boldsymbol{\beta}^{(i)})$ by $\boldsymbol{\mu}^{(i)}$ and $\boldsymbol{\Sigma}^{(i)}$, respectively. Calculating the derivative of each term in the above equality w.r.t. β_l , we obtain

$$\begin{aligned}
& \frac{\partial \left\| \mathbf{y} - \boldsymbol{\Phi}(\boldsymbol{\beta}) \boldsymbol{\mu}^{(i)} \right\|_2^2}{\partial \beta_l} \\
&= \frac{\partial \left\| \mathbf{y}_{-l}^{(i)} - \mu_l^{(i)} \cdot \mathbf{X}(\mathbf{a}(\hat{\vartheta}_l + \beta_l)) \right\|_2^2}{\partial \beta_l} \\
&= 2\text{Re} \left((\mathbf{a}'(\hat{\vartheta}_l + \beta_l))^H \mathbf{X}^H \mathbf{X} \mathbf{a}(\hat{\vartheta}_l + \beta_l) \right) \cdot |\mu_l^{(i)}|^2 \\
&\quad - 2\text{Re} \left((\mathbf{a}'(\hat{\vartheta}_l + \beta_l))^H \mathbf{X}^H \cdot (\mu_l^{(i)})^* \mathbf{y}_{-l}^{(i)} \right)
\end{aligned}$$

and

$$\begin{aligned}
& \frac{\partial \text{tr} \left(\boldsymbol{\Phi}(\boldsymbol{\beta}) \boldsymbol{\Sigma}^{(i)} \boldsymbol{\Phi}^H(\boldsymbol{\beta}) \right)}{\partial \beta_l} \\
&= 2\text{Re} \left((\mathbf{a}'(\hat{\vartheta}_l + \beta_l))^H \mathbf{X}^H \mathbf{X} \mathbf{a}(\hat{\vartheta}_l + \beta_l) \right) \cdot \chi_{ll}^{(i)} \\
&\quad + 2\text{Re} \left((\mathbf{a}'(\hat{\vartheta}_l + \beta_l))^H \mathbf{X}^H \mathbf{X} \cdot \sum_{j \neq l} \chi_{jl}^{(i)} \mathbf{a}(\hat{\vartheta}_j + \beta_j) \right).
\end{aligned}$$

Hence, the derivative of $\mathcal{U}(\alpha^{(i+1)}, \boldsymbol{\gamma}^{(i+1)}, \boldsymbol{\beta} | \alpha^{(i+1)}, \boldsymbol{\gamma}^{(i+1)}, \boldsymbol{\beta}^{(i)})$ w.r.t β_l is same as (39).

F. Proof of Theorem 5

According to Theorem 2-b in [25], the block MM algorithm will converge to a stationary solution if the following additional conditions are satisfied:

- All the properties in (26)–(30) hold true with the surrogate function.
- At least two of the problems (31)–(33) have a unique solution.

Lemmas 2–4 guarantee that the above two conditions hold true, respectively

REFERENCES

- [1] T. L. Marzetta, "Noncooperative cellular wireless with unlimited numbers of base station antennas," *IEEE Trans. Wireless Commun.*, vol. 9, no. 11, pp. 3590–3600, 2010.

- [2] F. Rusek, D. Persson, B. K. Lau, E. G. Larsson, T. L. Marzetta, O. Edfors, and F. Tufvesson, "Scaling up MIMO: Opportunities and challenges with very large arrays," *IEEE Signal Process. Mag.*, vol. 30, no. 1, pp. 40–60, 2013.
- [3] J.-C. Shen, J. Zhang, K.-C. Chen, and K. B. Letaief, "High-dimensional CSI acquisition in massive MIMO: Sparsity-inspired approaches," *IEEE Systems Journal*, vol. 11, no. 1, pp. 32–40, 2017.
- [4] E. G. Larsson, O. Edfors, F. Tufvesson, and T. L. Marzetta, "Massive MIMO for next generation wireless systems," *IEEE Commun. Mag.*, vol. 52, no. 2, pp. 186–195, 2014.
- [5] J. Hoydis, S. Ten Brink, and M. Debbah, "Massive MIMO in the UL/DL of cellular networks: How many antennas do we need?" *IEEE J. Sel. Areas Commun.*, vol. 31, no. 2, pp. 160–171, 2013.
- [6] B. Hassibi and B. M. Hochwald, "How much training is needed in multiple-antenna wireless links?" *IEEE Trans. Infor. Theory*, vol. 49, no. 4, pp. 951–963, 2003.
- [7] L. Lu, G. Y. Li, A. L. Swindlehurst, A. Ashikhmin, and R. Zhang, "An overview of massive MIMO: Benefits and challenges," *IEEE J. Sel. Topics in Signal Process.*, vol. 8, no. 5, pp. 742–758, 2014.
- [8] Z. Gao, L. Dai, Z. Wang, and S. Chen, "Spatially common sparsity based adaptive channel estimation and feedback for FDD massive MIMO," *IEEE Trans. Signal Process.*, vol. 63, no. 23, pp. 6169–6183, 2015.
- [9] J. Hoydis, C. Hoek, T. Wild, and S. ten Brink, "Channel measurements for large antenna arrays," in *International Symposium on ISWCS 2012*. IEEE, 2012, pp. 811–815.
- [10] X. Rao and V. K. Lau, "Distributed compressive CSIT estimation and feedback for FDD multi-user massive MIMO systems," *IEEE Trans. Signal Process.*, vol. 62, no. 12, pp. 3261–3271, 2014.
- [11] Z. Gao, C. Zhang, Z. Wang, and S. Chen, "Prior-Information aided iterative hard threshold: A low-complexity high-accuracy compressive sensing based channel estimation for TDS-OFDM," *IEEE Trans. Wireless Commun.*, vol. 14, no. 1, pp. 242–251, 2015.
- [12] Z. Chen and C. Yang, "Pilot decontamination in wideband massive MIMO systems by exploiting channel sparsity," *IEEE Trans. Wireless Commun.*, vol. 15, no. 7, pp. 5087–5100, 2016.
- [13] J.-C. Shen, J. Zhang, E. Alsusa, and K. B. Letaief, "Compressed CSI acquisition in FDD massive MIMO: How much training is needed?" *IEEE Trans. Wireless Commun.*, vol. 15, no. 6, pp. 4145–4156, 2016.
- [14] X. Rao and V. K. Lau, "Compressive sensing with prior support quality information and application to massive MIMO channel estimation with temporal correlation," *IEEE Trans. Signal Process.*, vol. 63, no. 18, pp. 4914–4924, 2015.
- [15] J. Choi, D. J. Love, and P. Bidigare, "Downlink training techniques for FDD massive MIMO systems: Open-loop and closed-loop training with memory," *IEEE J. Sel. Topics in Signal Process.*, vol. 8, no. 5, pp. 802–814, 2014.
- [16] L. You, X. Gao, A. L. Swindlehurst, and W. Zhong, "Channel acquisition for massive MIMO-OFDM with adjustable phase shift pilots," *IEEE Trans. Signal Process.*, vol. 64, no. 6, pp. 1461–1476, 2016.
- [17] Z. Gao, L. Dai, W. Dai, B. Shim, and Z. Wang, "Structured compressive sensing-based spatio-temporal joint channel estimation for FDD massive MIMO," *IEEE Trans. Commun.*, vol. 64, no. 2, pp. 601–617, 2016.
- [18] A. Liu, V. K. Lau, and W. Dai, "Exploiting burst-sparsity in massive MIMO with partial channel support information," *IEEE Trans. Wireless Commun.*, vol. 15, no. 11, pp. 7820–7830, 2016.
- [19] A. Liu, F. Zhu, and V. K. Lau, "Closed-loop autonomous pilot and compressive CSIT feedback resource adaptation in multi-user FDD massive MIMO systems," *IEEE Trans. Signal Process.*, vol. 65, no. 1, pp. 173–183, 2017.
- [20] Y. Ding and B. D. Rao, "Channel estimation using joint dictionary learning in FDD massive MIMO systems," in *IEEE GlobSIP 2015*. IEEE, 2015, pp. 185–189.
- [21] —, "Compressed downlink channel estimation based on dictionary learning in FDD massive MIMO systems," in *IEEE GLOBECOM 2015*. IEEE, 2015, pp. 1–6.
- [22] —, "Dictionary learning based sparse channel representation and estimation for FDD massive MIMO systems," *arXiv preprint arXiv:1612.06553*, 2016.
- [23] M. E. Tipping, "Sparse Bayesian learning and the relevance vector machine," *Journal of Machine Learning Research*, vol. 1, no. Jun, pp. 211–244, 2001.
- [24] S. Ji, Y. Xue, and L. Carin, "Bayesian compressive sensing," *IEEE Trans. Signal Process.*, vol. 56, no. 6, pp. 2346–2356, 2008.
- [25] M. Razaviyayn, "Successive convex approximation: Analysis and applications," Ph.D. dissertation, University of Minnesota, 2014.
- [26] Y. Sun, P. Babu, and D. P. Palomar, "Majorization-minimization algorithms in signal processing, communications, and machine learning," *IEEE Trans. Signal Process.*, vol. 65, no. 3, pp. 794–816, 2017.
- [27] Z. Yang, L. Xie, and C. Zhang, "Off-grid direction of arrival estimation using sparse Bayesian inference," *IEEE Trans. Signal Process.*, vol. 61, no. 1, pp. 38–43, 2013.
- [28] Q. Liu, H. C. So, and Y. Gu, "Off-grid DOA estimation with nonconvex regularization via joint sparse representation," *Signal Process.*, 2017.
- [29] J. Dai and H. C. So, "Sparse Bayesian learning approach for outlier-resistant direction-of-arrival estimation," *IEEE Trans. Signal Process.*, vol. 66, no. 3, pp. 744–756, 2018.
- [30] J. Dai, X. Bao, W. Xu, and C. Chang, "Root sparse Bayesian learning for off-grid DOA estimation," *IEEE Signal Process. Letters*, vol. 24, no. 1, pp. 46–50, 2017.
- [31] D. Donoho and Y. Tsaig, "Fast solution of l_1 -norm minimization problems when the solution may be sparse," *IEEE Trans. Information Theory*, vol. 54, no. 11, pp. 4789–4812, 2008.
- [32] D. L. Donoho, "Compressed sensing," *IEEE Trans. Infor. theory*, vol. 52, no. 4, pp. 1289–1306, 2006.
- [33] E. J. Candès, J. Romberg, and T. Tao, "Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information," *IEEE Trans. Infor. Theory*, vol. 52, no. 2, pp. 489–509, 2006.
- [34] D. P. Wipf and B. D. Rao, "Sparse Bayesian learning for basis selection," *IEEE Trans. Signal Process.*, vol. 52, no. 8, pp. 2153–2164, 2004.
- [35] D. Tse and P. Viswanath, *Fundamentals of Wireless Communication*. Cambridge university press, 2005.
- [36] 3GPP, "Universal mobile telecommunications system (UMTS); Spatial channel model for multiple input multiple output (MIMO) simulations," *3GPP TR 25.996 version 11.0.0 Release 11*, 2012.
- [37] A. F. Molisch, A. Kuchar, J. Laurila, K. Hugl, and R. Schmalenberger, "Geometry-based directional model for mobile radio channels: principles and implementation," *Trans. Emerging TeleComm. Technologies*, vol. 14, no. 4, pp. 351–359, 2003.
- [38] B. D. Rao and K. Kreutz-Delgado, "An affine scaling methodology for best basis selection," *IEEE Trans. Signal Process.*, vol. 47, no. 1, pp. 187–200, 1999.
- [39] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *Journal of the Royal Statistical Society. Series B (Methodological)*, pp. 1–38, 1977.
- [40] S. J. Wright and J. Nocedal, "Numerical optimization," *Springer Science*, vol. 35, no. 67–68, p. 7, 1999.
- [41] 3GPP, "3rd generation partnership project; Technical specification group radio access network; Study on 3D channel model for LTE," *3GPP TR 36.873 version 12.2.0 Release 12*, 2015.