

# Multi-Channel Hankel Matrix Completion through Nonconvex Optimization

Shuai Zhang, Yingshuai Hao, Meng Wang, *Member, IEEE*, Joe H. Chow, *Fellow, IEEE*

**Abstract**—This paper studies the multi-channel missing data recovery problem when the measurements are generated by a dynamical system. A new model, termed multi-channel low-rank Hankel matrices, is proposed to characterize the intrinsic low-dimensional structures in multi-channel time series. The data recovery problem is formulated as a nonconvex optimization problem, and two fast algorithms (AM-FIHT and RAM-FIHT), both with linear convergence rates, are developed to recover the missing points with provable performance guarantees. The required number of observations is significantly reduced compared with conventional low-rank completion methods. Our methods are verified through numerical experiments on synthetic data and recorded synchrophasor data in power systems.

**Index Terms**—low-rank matrix completion, nonconvex optimization, Hankel matrix, linear dynamic systems, synchrophasor data

## I. INTRODUCTION

Missing data recovery is an important task in various applications such as covariance estimation from partially observed correlations in remote sensing [10], multi-class learning in machine learning [3], [8], the Netflix Prize [1] problem and other similar questions in collaborative filtering [19]. Moreover, the recent framework of super-resolution enables accurate signal recovery from sparsely sampled measurements [9]. Example applications include magnetic resonance imaging (MRI) [21], [25], [43] and target localization in radar imaging [12], [47]. In power system monitoring, Phasor Measurement Units (PMU) [39] can measure voltage and current phasors directly at various locations and transmit the measurements to the operator for state estimation [2], [14] or disturbance identification [31]. Some PMU data points, however, do not reach the operator due to PMU malfunction or communication congestions. These missing data points should be recovered for the subsequent applications on PMU data [18].

Since practical datasets often have intrinsic low-dimensional structures, the missing data recovery problem can be formulated as a low-rank matrix completion problem, which is nonconvex due to the rank constraint. Its convex relaxation, termed Nuclear Norm Minimization (NNM) problem, has been extensively investigated [8], [11], [16], [20]. Given an  $n_c \times n$  ( $n_c \leq n$ ) matrix with rank  $r$  ( $r \ll n$ ), as long as  $O(rn \log^2 n)$ <sup>1</sup>

The first two authors contributed equally to this paper. The authors are with the Dept. of Electrical, Computer, and Systems Engineering, Rensselaer Polytechnic Institute, 110 8th Street, Troy, NY, 12180. Email: {zhangs21, haoy2, wangm7, chowj}@rpi.edu. Phone: 518-276-3842. Fax: 518-276-6261. Partial and preliminary results have appeared in [49].

<sup>1</sup> $f(n) = O(g(n))$  means that if for some constant  $C > 0$ ,  $f(n) \leq Cg(n)$  holds when  $n$  is sufficiently large.  $f(n) = \Theta(g(n))$  means that for some constants  $C_1 > 0$  and  $C_2 > 0$ ,  $C_1g(n) \leq f(n) \leq C_2g(n)$  holds when  $n$  is sufficiently large.

entries are observed, one can recover the remaining entries accurately by solving NNM [8], [11], [20].

Although elegant theoretical analyses exist, convex approaches like NNM have high computational complexity and poor convergence rate. For example, to decompose an  $n_c \times n$  matrix, the per-iteration complexity of the best specialized implementation is  $O(n_c^2 n)$  [36]. To reduce the computational complexity, first-order algorithms like [24] have been developed to solve the non-convex problem directly. Despite the numerical superiority, the theoretical analyses of the convergence and recovery performance of these nonconvex methods are still open problems. Only a few recent work such as [7], [24] provided such analyses on a case-by-base basis.

The low-rank matrix model, however, does not capture the temporal correlations in time series. A permutation of measurements at different time steps would result in different time series, but the rank of the data matrix remains the same. As a result, low-rank matrix completion methods require at least  $r$  entries in each column/row to recover the missing points and would fail if a complete column/row was lost. They cannot recover simultaneous data losses among all channels. Simultaneous data losses are not uncommon in power systems due to communication congestions.

There is limited study of the coupling of low-dimensional models and temporal correlations. Parametric models like hidden Markov models [33], [35] and autoregression (AR) models [22], [34] are employed to model temporal correlations. The accuracy of the algorithms depends on the correct estimation of model parameters, and no theoretical analysis is reported.

This paper develops a new model to characterize the intrinsic structures of multiple time series that are generated by a linear dynamical system. Our model of *multi-channel low-rank Hankel matrix* characterizes the temporal correlations in time series like PMU data without directly modeling the dynamical systems and estimating the system parameters. Our model can also be viewed as an extension of the single-channel low-rank Hankel matrix model with a  $\Theta(r)$  degree of freedom in [7], [12] to an  $n_c$ -channel matrix with a  $\Theta(n_c r)$  degree of freedom. It can also characterize spectrally sparse signals in applications like radar imaging [41] and magnetic resonance imaging [30].

Building upon the FIHT algorithm [7], this paper proposes two fast algorithms, termed accelerated multi-channel fast iterative hard thresholding (AM-FIHT) and robust AM-FIHT (RAM-FIHT) for multi-channel low-rank Hankel matrix completion. They can recover missing points for simultaneous data losses. The heavy ball method [29], [40] is employed to accelerate the convergence rate, and the acceleration is evaluated theoretically and numerically. Our algo-

gorithms converge linearly with a low per iteration complexity  $O(r^2 n_c n + r n_c n \log n + r^3)$  to the original matrix (noiseless measurements) or a sufficiently close matrix depending on the noise level (noisy measurements). Theoretical analyses of FIHT with only noiseless measurements are reported [7]. Moreover, the recovery is successful as long as the number of observed measurements is  $O(r^2 \log^2 n)$ , significantly lower than  $O(r n \log^2 n)$  for general rank- $r$  matrices. This number is also a constant fraction of the required number of measurements by applying the single-channel Hankel matrix completion methods like FIHT [7] to each channel separately.

The rest of the paper is organized as follows. Sections II-III describe the problem formulation and the connection with the existing work. Sections IV-V present our algorithms and the theoretical analyses. Section VI records the numerical results on synthetic data and recorded PMU data. Section VII concludes the paper. All the proofs are summarized in Appendix.

*Notation:* Vectors are bold lowercase, matrices are bold uppercase, and scalars are in normal font. For example,  $\mathbf{Z}$  is a matrix and  $\mathbf{z}$  is vector.  $\mathbf{Z}_{i*}$  denotes the  $i$ th row of  $\mathbf{Z}$ , and  $Z_{ij}$  denotes the  $(i, j)$ -th entry of  $\mathbf{Z}$ .  $\mathbf{I}$  and  $\mathbf{e}_i$  denote the identity matrix and the  $i$ th standard basis vector.  $\mathbf{Z}^T$  and  $\mathbf{Z}^*$  denote the transpose and conjugate transpose of  $\mathbf{Z}$ , so as  $\mathbf{z}^T$  and  $\mathbf{z}^*$ . The inner product between two vectors is  $\langle \mathbf{z}_1, \mathbf{z}_2 \rangle = \mathbf{z}_2^* \mathbf{z}_1$ , and corresponding  $l_2$  norm is  $\|\mathbf{z}\| = \langle \mathbf{z}, \mathbf{z} \rangle^{1/2}$ . For matrices, the inner product is defined as  $\langle \mathbf{Z}_1, \mathbf{Z}_2 \rangle = \text{Tr}(\mathbf{Z}_2^* \mathbf{Z}_1)$ .  $\|\mathbf{Z}\|_F$  stands for the Frobenius norm with  $\|\mathbf{Z}\|_F = \langle \mathbf{Z}, \mathbf{Z} \rangle^{1/2}$ . The spectral norm of matrix  $\mathbf{Z}$  is denoted by  $\|\mathbf{Z}\|$ . The maximum entry (in absolute value) of  $\mathbf{Z}$  is denoted as  $\|\mathbf{Z}\|_\infty$ . Linear operators on matrix spaces will be denoted by calligraphic letters. In particular,  $\mathcal{I}$  is the identity operator. The spectral norm of a linear operator  $\mathcal{A}$  on matrix spaces is denoted as  $\|\mathcal{A}\| = \sup_{\langle \mathbf{Z}, \mathbf{Z} \rangle \leq 1} \|\mathcal{A}\mathbf{Z}\|_F$ . The adjoint operator of  $\mathcal{A}$  is denoted as  $\mathcal{A}^*$ , which satisfies  $\langle \mathcal{A}\mathbf{Z}_1, \mathbf{Z}_2 \rangle = \langle \mathbf{Z}_1, \mathcal{A}^* \mathbf{Z}_2 \rangle$ .

## II. PROBLEM FORMULATION

Consider an  $n_p$ -th order linear dynamical system after an impulse response. Let  $\mathbf{s}_t \in \mathbb{C}^{n_p}$  and  $\mathbf{x}_t \in \mathbb{C}^{n_c}$  denote deviations of state variables and observations at time  $t$  from the equilibrium point. Then we have

$$\mathbf{s}_{t+1} = \mathbf{A}\mathbf{s}_t, \quad \mathbf{x}_t = \mathbf{C}\mathbf{s}_t, \quad t = 1, 2, \dots, n, \quad (1)$$

where  $\mathbf{A} \in \mathbb{C}^{n_p \times n_p}$ , and  $\mathbf{C} \in \mathbb{C}^{n_c \times n_p}$ . Let  $\mathbf{X}$  contain the measurements from time 1 to  $n$ ,

$$\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \in \mathbb{C}^{n_c \times n}. \quad (2)$$

Further, the Hankel matrix of  $\mathbf{X}$  is defined as

$$\mathcal{H}\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_{n_2} \\ \mathbf{x}_2 & \mathbf{x}_3 & \dots & \mathbf{x}_{n_2+1} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_{n_1} & \mathbf{x}_{n_1+1} & \dots & \mathbf{x}_n \end{bmatrix} \in \mathbb{C}^{n_c n_1 \times n_2}, \quad (3)$$

where  $n_1 + n_2 = n + 1$ .

Suppose  $\mathbf{A}$  could be diagonalized, denoted by  $\mathbf{A} = \mathbf{P}\mathbf{\Lambda}\mathbf{P}^{-1}$ , where  $\mathbf{P} = [\mathbf{l}_1, \mathbf{l}_2, \dots, \mathbf{l}_{n_p}]$ ,  $\mathbf{P}^{-1} = [\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_{n_p}]^*$ ,

and  $(\cdot)^*$  stands for the conjugate transpose.  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_{n_p})$  contains the eigenvalues of  $\mathbf{A}$ . Then

$$\begin{aligned} \mathbf{x}_{t+1} &= \mathbf{C} \underbrace{\mathbf{A} \cdots \mathbf{A}}_{t \text{ times}} \mathbf{s}_1 = \mathbf{C}\mathbf{A}^t \mathbf{s}_1 = \mathbf{C}\mathbf{P}\mathbf{\Lambda}^t \mathbf{P}^{-1} \mathbf{s}_1 \\ &= \sum_{i=1}^{n_p} \lambda_i^t \mathbf{r}_i^* \mathbf{s}_1 \mathbf{C}\mathbf{l}_i. \end{aligned} \quad (4)$$

All  $n_p$  modes of the system are considered in (4). In practice, a mode might be highly damped ( $|\lambda_i| \approx 0$ ), or not excited by the input ( $|\mathbf{r}_i^* \mathbf{s}_1| \approx 0$ ), or not directly measured ( $\|\mathbf{C}\mathbf{l}_i\| \approx 0$ ). If only  $r$  ( $r \ll n$ ) out of  $n$  modes are significant, assuming these modes to be  $\lambda_1, \dots, \lambda_r$  for simplicity, we have

$$\mathbf{x}_{t+1} \simeq \sum_{i=1}^r \lambda_i^t \mathbf{r}_i^* \mathbf{s}_1 \mathbf{C}\mathbf{l}_i. \quad (5)$$

Then the corresponding Hankel matrix can be written as

$$\mathcal{H}\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_{n_2} \\ \mathbf{x}_2 & \mathbf{x}_3 & \dots & \mathbf{x}_{n_2+1} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_{n_1} & \mathbf{x}_{n_1+1} & \dots & \mathbf{x}_n \end{bmatrix} = \mathbf{P}_L \mathbf{\Gamma} \mathbf{P}_R^T, \quad (6)$$

where

$$\mathbf{P}_L = \begin{bmatrix} \mathbf{I}_{n_c} & \mathbf{I}_{n_c} & \dots & \mathbf{I}_{n_c} \\ \lambda_1 \mathbf{I}_{n_c} & \lambda_2 \mathbf{I}_{n_c} & \dots & \lambda_r \mathbf{I}_{n_c} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_1^{n_1-1} \mathbf{I}_{n_c} & \lambda_2^{n_1-1} \mathbf{I}_{n_c} & \dots & \lambda_r^{n_1-1} \mathbf{I}_{n_c} \end{bmatrix} \in \mathbb{C}^{n_c n_1 \times n_c r}, \quad (7)$$

$$\mathbf{\Gamma} = \begin{bmatrix} \mathbf{r}_1^* \mathbf{x}_1 \mathbf{C}\mathbf{l}_1 & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{r}_2^* \mathbf{x}_1 \mathbf{C}\mathbf{l}_2 & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{r}_r^* \mathbf{x}_1 \mathbf{C}\mathbf{l}_r \end{bmatrix} \in \mathbb{C}^{n_c r \times r}, \quad (8)$$

and

$$\mathbf{P}_R = \begin{bmatrix} 1 & 1 & \dots & 1 \\ \lambda_1 & \lambda_2 & \dots & \lambda_r \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_1^{n_2-1} & \lambda_2^{n_2-1} & \dots & \lambda_r^{n_2-1} \end{bmatrix} \in \mathbb{C}^{n_2 \times r}, \quad (9)$$

where  $\mathbf{I}_{n_c} \in \mathbb{C}^{n_c \times n_c}$  is the identity matrix. One can check that both  $\mathbf{X}$  and  $\mathcal{H}\mathbf{X}$  are rank  $r$  matrices<sup>2</sup>.

Let  $\mathbf{N} \in \mathbb{C}^{n_c \times n}$  denote the measurement noise.  $\mathbf{M} = \mathbf{X} + \mathbf{N}$  denotes the noisy measurements. Some entries of  $\mathbf{M}$  are not observed due to data losses. Let  $\Omega$  denote the index set of observed entries. The objective of missing data recovery is to reconstruct the missing data based on the observed entries  $\mathcal{P}_{\Omega}(\mathbf{M})$ . Since the rank of  $\mathcal{H}\mathbf{X}$  is  $r$ , the data recovery problem can be formulated as

$$\min_{\mathbf{Z} \in \mathbb{C}^{n_c \times n}} \|\mathcal{P}_{\Omega}(\mathbf{Z} - \mathbf{M})\|_F^2 \quad \text{subject to } \text{rank}(\mathcal{H}\mathbf{Z}) = r, \quad (10)$$

where  $\mathcal{P}_{\Omega}(\cdot)$  is the sampling operator with  $(\mathcal{P}_{\Omega}(\mathbf{Z}))_{ij} = Z_{ij}$  if  $(i, j) \in \Omega$  and 0 otherwise. (10) is a nonconvex problem due

<sup>2</sup>We assume  $\mathcal{H}\mathbf{X}$  is exactly rank  $r$  throughout the paper. The methods analyses can be extended to approximately low-rank matrices with minor modifications. If  $\mathcal{H}\mathbf{X}$  is approximately low-rank, i.e., its rank- $r$  approximation error is very small, we seek to find the best rank- $r$  approximation to  $\mathcal{H}\mathbf{X}$ . Then the recovery error is at least the approximation error.

to the rank constraint. It reduces to the conventional matrix completion problem when  $n_1 = 1$ .

Clearly, the recovery is impossible if  $\mathbf{X}$  is in the null space of  $\mathcal{P}_{\hat{\Omega}}(\cdot)$ . Here we follow the standard incoherence assumption in low-rank matrix completion [11].

**Definition 1.** A matrix  $\mathbf{Z} \in \mathbb{C}^{l_1 \times l_2}$  with singular value decomposition (SVD) as  $\mathbf{Z} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$ , is said to be incoherent with parameter  $\mu$  if

$$\max_{1 \leq k_1 \leq l_1} \|\mathbf{e}_{k_1}^* \mathbf{U}\|^2 \leq \frac{\mu r}{l_1}, \quad \max_{1 \leq k_2 \leq l_2} \|\mathbf{e}_{k_2}^* \mathbf{V}\|^2 \leq \frac{\mu r}{l_2}, \quad (11)$$

where  $\mathbf{e}_{k_1}, \mathbf{e}_{k_2}$  are the coordinate unit vectors.

The incoherence definition guarantees that the singular vectors of the matrix are sufficiently spread, and  $\mathcal{P}_{\hat{\Omega}}(\cdot)$  samples enough information about the matrix. We focus on recovering  $\mu$ -incoherence matrices in this paper.

### III. BACKGROUND AND RELATED WORK

The low-rank property of a Hankel matrix is also recently exploited in the direction of arrival (DOA) problem in array signal processing [12], [47], MRI image recovery from under-sampled measurements [21], [37], [48], video inpainting [15] and system identification [17]. To see the connection with our model, the  $k$ th row of  $\mathbf{X}$  in (2), denoted by  $\mathbf{X}_{k*}$ , can be equivalently viewed as the discrete samples of a spectrally sparse signal  $g_k(t)$ , which is a weighted sum of  $r$  damped or undamped sinusoids at  $t = \{0, \dots, n-1\}$ , where

$$g_k(t) = \sum_{i=1}^r d_{k,i} e^{(2\pi i f_i - \tau_i)t}, \quad k = 1, \dots, n_c, \quad (12)$$

and  $f_i$  and  $d_{k,i}$  are the frequency and the normalized complex amplitude of the  $i$ th sinusoid, respectively.  $\iota$  is the imaginary unit. The connection between (12) and (2) is that  $\lambda_i = e^{2\pi i f_i - \tau_i}$  and  $d_{k,i} = \mathbf{r}_i^* \mathbf{s}_1 \mathbf{C}_{k*} \mathbf{l}_i$ .

The signal of interest itself in array signal processing is spectrally sparse. In MRI imaging, if a signal reduces to a sparse linear combination of Dirac delta functions under some transformations, then its Fourier transform is a sum of a few sinusoids [25], [38], [48]. Most existing work on low-rank Hankel matrices studied single-channel signals, i.e.,  $n_c = 1$  in our setup. References [12], [26], [37], [38] considered 2-dimensional (2-D) and higher-dimensional signals, while a 2-D signal is still a sum of  $r$  2-D sinusoids, and the degree of freedom is still  $\Theta(r)$ . The focus of this paper is multi-channel signals with  $n_c > 1$ . Each signal is a weighted sum of the same set of  $r$  sinusoids, while the weights  $d_{k,i}$  are different for each channel  $k = 1, \dots, n_c$ . The degree of the freedom of (12) is  $\Theta(n_c r)$ .

The multi-channel signal in (12) is related to the multiple measurement vector (MMV) problem [13]. References [27], [47] considered data recovery of MMV when the signals are linear combinations of undamped sinusoids, i.e.,  $\tau_i = 0$  for all  $i$  in (12). The data recovery is achieved in [27], [47] through atomic norm minimization, which requires solving large-scale semidefinite programs. Besides the high computational complexity, it is not clear how the atomic norm can be extended to handle damped sinusoids, i.e.,  $\tau_i \neq 0$ . References [4], [25] studied multi-channel signal recovery using Hankel structures

and can thus handle damped sinusoids. Despite the numerical evaluations, there is no theoretical analysis of the recovery guarantee in [4], [25]. This paper provides analytical recovery guarantees for multi-channel damped and undamped sinusoids.

The recovery of a low-rank Hankel matrix can be formulated as a convex optimization, for example, nuclear norm minimization for missing data recovery [15], [17], [25], [37], [45], [48] and minimizing a weighted sum of the nuclear norm and the  $\ell_1$  norm for bad data correction [26]. Since it is computationally challenging to solve these convex problems for high-dimensional Hankel matrices, fast algorithms to recover missing points in single-channel [7] and multi-channel Hankel matrices [5], [15] are proposed recently. Although numerical results are reported in [5], [15], only [7] provides the theoretical performance analysis of the proposed fast iterative hard thresholding (FIHT) algorithm for single-channel Hankel matrix recovery. FIHT is a projected gradient descent method. In each iteration, the algorithm updates the estimate along the gradient descent direction and then projects it to a rank- $r$  matrix. To reduce the computational complexity, instead of solving singular value decomposition (SVD) directly, FIHT first projects a matrix onto a  $2r$ -dimensional subspace and then computes the SVD of the rank- $2r$  matrix. The per-iteration complexity of FIHT is  $O(r^2 n + rn \log n + r^3)$ .

Motivated by PMU data analysis in power systems, this paper connects dynamical systems with low-rank Hankel matrices. It develops fast data recovery algorithms for multi-channel Hankel matrices with provable performance guarantees.

### IV. DATA RECOVERY ALGORITHMS

Here we describe two algorithms to solve (10) and defer the theoretical analyses to Section V. One is accelerated multi-channel fast iterative hard thresholding algorithm (AM-FIHT), and the other one is robust AM-FIHT (RAM-FIHT). Both algorithms are built upon the FIHT [7] with some major differences. First, FIHT recovers the missing points of one spectrally sparse signal, while (R)AM-FIHT recovers the missing points of  $n_c$  signals simultaneously. The simultaneous recovery can reduce the required number of measurements, as quantified in Theorem 5. Second, (R)AM-FIHT has a heavy-ball step [29], [40], e.g., term  $\beta(\mathbf{W}_{l-1} - \mathbf{W}_{l-2})$  in line 5 of Algorithm 1 and line 14 of Algorithm 2, while FIHT does not. The basic idea of the heavy ball method is to compute the search direction using a linear combination of the gradient at the current iterate and the update direction in the previous step, rather than being memoryless of the past iterates' trajectory [29]. We will show analytically that with the heavy-ball step, AM-FIHT converges faster while maintaining the recovery accuracy (Theorem 2). Third, we provide the theoretical guarantee of data recovery when the measurements are noisy (Theorem 4), while [7] only has the performance guarantee of FIHT using noiseless measurements.

In both algorithms,  $\mathbf{M}, \mathbf{X}_l, \mathbf{G}_l \in \mathbb{C}^{n_c \times n}$ , and  $\mathbf{W}_l, \Delta \mathbf{W}_l, \mathbf{L}_l \in \mathbb{C}^{n_c n_1 \times n_2}$ .  $\mathbf{L}_l$  is a rank- $r$  matrix and its SVD is denoted as  $\mathbf{L}_l = \mathbf{U}_l \mathbf{\Sigma}_l \mathbf{V}_l^*$ , where  $\mathbf{U}_l \in \mathbb{C}^{n_c n_1 \times r}$ ,  $\mathbf{V}_l \in \mathbb{C}^{n_2 \times r}$  and  $\mathbf{\Sigma}_l \in \mathbb{C}^{r \times r}$ .  $\mathcal{S}_l$  is the tangent subspace of the rank- $r$  Riemannian

nian manifold at  $L_l$ , and for any matrix  $Z \in \mathbb{C}^{n_c n_1 \times n_2}$ , the projection of  $Z$  onto  $S_l$  is defined as

$$\mathcal{P}_{S_l}(Z) = U_l U_l^* Z + Z V_l V_l^* - U_l U_l^* Z V_l V_l^*. \quad (13)$$

$\mathcal{Q}_r$  finds the best rank- $r$  approximation as

$$\mathcal{Q}_r(Z) = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^*, \quad (14)$$

if  $Z = \sum_i \sigma_i \mathbf{u}_i \mathbf{v}_i^*$  is the SVD of  $Z$  with  $\sigma_1 \geq \sigma_2 \geq \dots$ .  $\mathcal{H}^\dagger$  is the Moore-Penrose pseudoinverse of  $\mathcal{H}$ . For any matrix  $Z \in \mathbb{C}^{n_c n_1 \times n_2}$ ,  $(\mathcal{H}^\dagger Z) \in \mathbb{C}^{n_c \times n}$  satisfies

$$\langle \mathcal{H}^\dagger Z, \mathbf{e}_k \mathbf{e}_t^* \rangle = \frac{1}{w_t} \sum_{k_1+k_2=t+1} Z_{(k_1-1)n_c+k, k_2}, \quad (15)$$

where  $w_t = \#\{(k_1, k_2) | k_1 + k_2 = t + 1, 1 \leq k_1 \leq n_1, 1 \leq k_2 \leq n_2\}$  as the number of elements in the  $t$ -th anti-diagonal of an  $n_1 \times n_2$  matrix.

---

**Algorithm 1** AM-FIHT for Data Recovery from Noiseless Measurements

---

**Require:**  $\mathcal{P}_{\hat{\Omega}}(M)$ ,  $n_1$ ,  $n_2$ ,  $r$

- 1: Set  $\mathbf{W}_{-2} = \mathbf{0}$ ,  $\mathbf{W}_{-1} = p^{-1} \mathcal{H} \mathcal{P}_{\hat{\Omega}}(M)$ ,  $\mathbf{L}_0 = \mathcal{Q}_r(\mathbf{W}_{-1})$ ;
  - 2: Initialize  $\mathbf{X}_0 = \mathcal{H}^\dagger \mathbf{L}_0$ ;
  - 3: **for**  $l = 0, 1, \dots$  **do**
  - 4:    $\mathbf{G}_l = \mathcal{P}_{\hat{\Omega}}(M - \mathbf{X}_l)$ ;
  - 5:    $\mathbf{W}_l = \mathcal{P}_{S_l}(\mathcal{H}(\mathbf{X}_l + p^{-1} \mathbf{G}_l) + \beta(\mathbf{W}_{l-1} - \mathbf{W}_{l-2}))$ ;
  - 6:    $\mathbf{L}_{l+1} = \mathcal{Q}_r(\mathbf{W}_l)$ ;
  - 7:    $\mathbf{X}_{l+1} = \mathcal{H}^\dagger \mathbf{L}_{l+1}$ ;
  - 8: **end for**
  - 9: **return**  $\mathbf{X}_l$
- 

The key steps in AM-FIHT are as follows. Here the measurements are noiseless, thus  $M = X$ . In each iteration, we first update current  $X_l$  along the gradient descent direction  $G_l$ , with a step size  $p^{-1} = \frac{n_c n}{m}$ , where  $m$  is the number of observed entries. To improve the convergence rate, the update is further combined with additional heavy-ball term  $\beta(\mathbf{W}_{l-1} - \mathbf{W}_{l-2})$ , which represents the update direction in the previous iteration. Next,  $\mathcal{H}(\mathbf{X}_l + p^{-1} \mathbf{G}_l) + \beta(\mathbf{W}_{l-1} - \mathbf{W}_{l-2})$  is projected to a rank- $r$  matrix. To reduce the computational complexity, we first project it to the  $2r$ -dimensional space  $S_l$  and then apply SVD on the rank- $2r$  matrix [7], instead of directly computing its SVD. The rank- $r$  matrix  $L_{l+1}$  is obtained in line 7 by thresholding the singular values of the rank- $2r$  matrix  $W_l$ . Finally,  $X_{l+1}$  is updated by  $\mathcal{H}^\dagger L_{l+1}$ .

The analysis of the computational complexity of AM-FIHT is similar to that of FIHT [7] with some modifications for the  $n_c$ -channel signal and the heavy-ball step. Details are in the supplementary materials. The computational complexity of solving SVD of a matrix in  $\mathbb{C}^{n_c n_1 \times n_2}$  is generally  $O(n_c n^2 r)$ . Due to the low rank structure of the matrices in  $S_l$ , the SVD of  $W_l \in \mathbb{C}^{n_c n_1 \times n_2}$  can be computed in  $O(n_c n r^2 + r^3)$  via QR decompositions and SVD on a  $2r \times 2r$  matrix [7]. Moreover, it is not necessary to construct Hankel matrices following (3) explicitly. The matrix multiplication of  $U_l^* \mathcal{H} X_l \in \mathbb{C}^{r \times n_2}$  and  $(\mathcal{H} X_l) V_l \in \mathbb{C}^{n_c n_1 \times r}$  in line 5 can be completed via fast convolution algorithms with  $O(n_c n r \log(n))$  flops, instead of

the conventional complexity of  $O(n_c n^2 r)$ . Similar analysis can be applied to line 7, which costs  $O(n_c n r \log(n))$  flops to compute  $X_{l+1}$  from the SVD of  $L_{l+1}$  directly.

With the heavy ball term, since the SVDs of  $W_{l-1}$  and  $W_{l-2}$  have been obtained in the last two steps, we compute  $\mathcal{P}_{S_l}(W_{l-1}) - \mathcal{P}_{S_l}(W_{l-2})$  in line 5. From (13), the computation of  $U_l U_l^* Z V_l V_l^*$  plays the dominant part in computing  $\mathcal{P}_{S_l}(Z)$ . Let  $W_l = U_{W_l} \Sigma_{W_l} V_{W_l}^*$  denote the SVD of  $W_l$ , where  $U_{W_l} \in \mathbb{C}^{n_c n_1 \times 2r}$ ,  $V_{W_l} \in \mathbb{C}^{n_2 \times 2r}$ . Then computing  $U_l^* U_{W_{l-1}}$  and  $V_{W_{l-1}}^* V_l$  requires  $O(n_c n r^2)$  and  $O(n r^2)$  flops, respectively. Computing  $U_l^* U_{W_{l-1}} \Sigma_{W_{l-1}} V_{W_{l-1}}^* V_l$  further requires  $O(r^3)$  flops.

From the above analysis, line 4 requires  $O(n_c n)$  flops. The complexity of line 5 is  $O(n_c n r \log(n) + n_c n r^2 + r^3)$ . Line 6 requires  $O(n_c n r^2 + r^3)$  flops, and line 7 requires  $O(n_c n r \log(n))$  flops. Thus, the total per-iteration complexity of AM-FIHT is  $O(r^2 n_c n + r n_c n \log n + r^3)$ .

---

**Algorithm 2** RAM-FIHT

---

**Require:**  $\mathcal{P}_{\hat{\Omega}}(M)$ ,  $n_1$ ,  $n_2$ ,  $r$ ,  $\mu$ , and  $\beta$ .

- 1: Partition  $\hat{\Omega}$  into  $L + 1$  disjoint sets  $\hat{\Omega}_0, \hat{\Omega}_1, \dots, \hat{\Omega}_L$  of equal size  $\hat{m}$ , let  $\hat{p} = \frac{\hat{m}}{n}$ .
  - 2: Set  $\mathbf{W}_{-2} = \mathbf{0}$ ,  $\mathbf{W}_{-1} = \hat{p}^{-1} \mathcal{H} \mathcal{P}_{\hat{\Omega}_0}(M)$ ,  $\mathbf{L}_0 = \mathcal{Q}_r(\mathbf{W}_{-1})$ ;
  - 3: **for**  $l = 0, 1, \dots, L - 1$  **do**;
  - 4:    $[\mathbf{U}_l, \Sigma_l, \mathbf{V}_l] = \text{SVD}(\mathbf{L}_l)$ ;
  - 5:   **for**  $i = 1, 2, \dots, n_c n_1$  **do**
  - 6:      $(\mathbf{A}_l)_{i*} = \frac{(\mathbf{U}_l)_{i*}}{\|(\mathbf{U}_l)_{i*}\|} \min \left\{ \|(\mathbf{U}_l)_{i*}\|, \sqrt{\frac{\mu r}{n_c n_1}} \right\}$ ;
  - 7:   **end for**
  - 8:   **for**  $i = 1, 2, \dots, n_2$  **do**
  - 9:      $(\mathbf{B}_l)_{i*} = \frac{(\mathbf{V}_l)_{i*}}{\|(\mathbf{V}_l)_{i*}\|} \min \left\{ \|(\mathbf{V}_l)_{i*}\|, \sqrt{\frac{\mu r}{n_2}} \right\}$ ;
  - 10:   **end for**
  - 11:    $\mathbf{L}'_l = \mathbf{A}_l \Sigma_l \mathbf{B}_l^*$ ;
  - 12:    $\hat{\mathbf{X}}_l = \mathcal{H}^\dagger \mathbf{L}'_l$ ;
  - 13:    $\mathbf{G}_l = \mathcal{P}_{\hat{\Omega}_{l+1}}(M - \hat{\mathbf{X}}_l)$ ;
  - 14:    $\mathbf{W}_l = \mathcal{P}_{S_l}(\mathcal{H}(\hat{\mathbf{X}}_l + \hat{p}^{-1} \mathbf{G}_l) + \beta(\mathbf{W}_{l-1} - \mathbf{W}_{l-2}))$ ;
  - 15:    $\mathbf{L}_{l+1} = \mathcal{Q}_r(\mathbf{W}_l)$ ;
  - 16: **end for**
  - 17: **return**  $\mathbf{X}_L = \mathcal{H}^\dagger \mathbf{L}_L$ ;
- 

RAM-FIHT differs from AM-FIHT mainly in resampling (line 1) and trimming (lines 5-10). The resampling and trimming are used in [7] to improve the initialization of FIHT. Here we apply these ideas in the data recovery algorithm and prove in Theorem 4 that the resulting RAM-FIHT can recover the matrix even when the observed measurements are noisy. There is no analytical analysis of FIHT on noisy measurements in [7]. Moreover, compared with AM-FIHT, we provide tighter bounds of the required number of observations for RAM-FIHT (comparing Theorems 1 and 3).

In RAM-FIHT, the sampling set  $\hat{\Omega}$  is divided into  $L$  disjoint subsets  $\hat{\Omega}_l$ 's. During the  $l$ -th iteration,  $L_l$  is updated using the observed entries in  $\hat{\Omega}_l$ , instead of using all the entries in  $\hat{\Omega}$  as in AM-FIHT. The partition of the sampling set is a standard technique in analyzing matrix completion (MC) problems [42]. The disjointness of  $L_l$ 's in different iterations simplifies

the theoretical analyses<sup>3</sup>, since it ensures the independence between  $\mathbf{X}_l$  and  $\mathbf{X}_{l+1}$ . The trimming procedure ensures that the estimate in each iteration remains close to  $\mu$ -incoherent, which in turn helps to obtain tighter bounds of the recovery performance in Theorem 3. We remark that the resampling and trimming steps in RAM-FIHT are introduced mainly to simplify the theoretical analyses and obtain tighter bounds, while we observe numerically that AM-FIHT and RAM-FIHT perform similarly in Section VI. The per iteration computational complexity of RAM-FIHT is  $O(r^2 n_c n + r n_c n \log n + r^3)$ .

## V. THEORETICAL ANALYSES

The theoretical analyses of the convergence rates and recovery accuracy of AM-FIHT and RAM-FIHT are summarized in the following four theorems. All the proofs are deferred to the Appendix. Theorem 1 records the recovery performance of AM-FIHT using noiseless measurements with  $\beta = 0$ . Theorem 2 shows that the convergence rate of AM-FIHT can be further improved by using a small positive  $\beta$ . Theorems 3 and 4 discuss the recovery performance of RAM-FIHT from noiseless and noisy measurements, respectively. We also compare the recovery performance with recovering missing points on each individual row of  $\mathbf{X}$  separately and quantify the performance gain of our algorithms in Theorem 5.

**Theorem 1.** (*AM-FIHT with noiseless measurements.*) Assume  $\mathcal{H}\mathbf{X}$  is  $\mu$ -incoherent. Let  $0 < \varepsilon_0 < \frac{1}{10}$  be a numerical constant and  $\nu = 6\varepsilon_0 < 1$ . Then with probability at least  $1 - 3n_c n^{-2}$ , the iterates  $\mathbf{X}_l$ 's generated by AM-FIHT with  $\beta = 0$  satisfy

$$\|\mathbf{X}_l - \mathbf{X}\|_F \leq \nu^{l-1} \|\mathbf{L}_0 - \mathcal{H}\mathbf{X}\|_F, \quad (16)$$

provided that

$$m \geq C_1 \max \left\{ \frac{\mu c_s r \log(n)}{\varepsilon_0^2}, \frac{1 + \varepsilon_0}{\varepsilon_0} (n_c \mu c_s n)^{\frac{1}{2}} \kappa r \log^{\frac{3}{2}}(n) \right\} \quad (17)$$

for some constant  $C_1 > 0$ , where  $\kappa = \frac{\sigma_{\max}(\mathcal{H}\mathbf{X})}{\sigma_{\min}(\mathcal{H}\mathbf{X})}$  denotes the condition number of  $\mathcal{H}\mathbf{X}$  and  $c_s = \max\{\frac{n}{n_1}, \frac{n}{n_2}\}$ .

Theorem 1 indicates that if the number of noiseless observations is  $O(r n_c^{1/2} n^{1/2} \log^{3/2}(n))$ , then AM-FIHT is guaranteed to recover  $\mathbf{X}$  exactly. Moreover, from (16), the iterates generated by AM-FIHT converge linearly to the groundtruth  $\mathbf{X}$ , and the rate of convergence is  $\nu$ . Since  $\mathbf{X}$  is rank  $r$ , if one directly applies a conventional low-rank matrix completion method such as NNM ([8], [11], [20]), the required number of observations is  $O(r n \log^2(n))$ . Thus, when  $n$  is large, by exploiting the low-rank Hankel structure of correlated time series, the required number of measurements is significantly reduced. Note that the degree of freedom of  $\mathbf{X}$  is  $\Theta(n_c r)$ , as one can see from (12), the required number of observations by Theorem 1 is suboptimal due to the dependence upon  $n$ . This results from the artefacts in our proof techniques. We will provide a tighter bound for RAM-FIHT in Theorem 3.

The required number of measurements depends on  $c_s$ , which is minimized when  $n_1 = n_2 = \frac{n+1}{2}$ . In practice, the selection of  $n_1$  and  $n_2$  of the Hankel matrix is also affected by the accuracy of the low-rank approximation.

<sup>3</sup>In fact, we only need the mutual independence among the subsets  $\widehat{\Omega}_i$ 's, and the disjoint partition is a sufficient condition.

We set  $\beta$  as 0 in Theorem 1 to simplify the analyses. The improvement of the convergence rate by using a positive  $\beta$  is quantified in the following theorem.

**Theorem 2.** (*Faster convergence with a heavy-ball step*) Given any  $\beta \in [0, \tau)$  for some  $\tau > 0$ , let  $\mathbf{X}_l$ 's denote the convergent iterates returned by AM-FIHT. There exists an integer  $s_0$ , a constant  $q \in (0, 1)$  that depends on  $\beta$  such that

$$\|\mathbf{X}_{s_0+k} - \mathbf{X}\|_F \leq c(\delta)(q(\beta) + \delta)^k, \quad \forall k \geq 0 \quad (18)$$

holds for any  $\delta \in (0, 1 - q(\beta))$  and a positive  $c(\delta)$  that depends on  $\delta$ . Moreover,

$$q(0) > q(\beta), \quad \forall \beta \in (0, \tau). \quad (19)$$

The exact expressions of  $q$  and  $\tau$  are deferred to the proofs in Appendix (equation (53)). Theorem 2 indicates that by adding a heavy-ball term, when close enough to the groundtruth  $\mathbf{X}$ , the iterates converge linearly to  $\mathbf{X}$ , and the rate of convergence is  $q(\beta) + \delta$ . Moreover, from (19), with a small positive  $\beta$ , the iterates converge faster than those without the heavy-ball step. Such improvement is numerically evaluated in Section VI.

**Theorem 3.** (*RAM-FIHT with noiseless measurements*) Assume  $\mathcal{H}\mathbf{X}$  is  $\mu$ -incoherent. Let  $0 < \varepsilon_0 < \frac{1}{2}$  and

$$L = \left\lceil \varepsilon_0^{-1} \log \left( \frac{\sigma_{\max}(\mathcal{H}\mathbf{X})}{128\kappa^3\varepsilon} \right) \right\rceil. \quad (20)$$

Define  $\nu = 2\varepsilon_0 < 1$ . Then with probability at least  $1 - (2L + 3)n_c n^{-2}$ , for any arbitrarily small constant  $\varepsilon > 0$ , the iterates  $\mathbf{L}_l$ 's and  $\mathbf{X}_L$  generated by RAM-FIHT with  $\beta = 0$  satisfy

$$\begin{aligned} \|\mathbf{L}_l - \mathbf{X}\|_F &\leq \nu^l \|\mathbf{L}_0 - \mathcal{H}\mathbf{X}\|_F, \quad 1 \leq l \leq L, \\ \text{and } \|\mathbf{X}_L - \mathbf{X}\|_F &\leq \nu^L \|\mathbf{L}_0 - \mathcal{H}\mathbf{X}\|_F \leq \varepsilon, \end{aligned}$$

provided that

$$m \geq C_2 \varepsilon_0^{-3} \mu c_s \kappa^6 r^2 \log(n) \log \left( \frac{\sigma_{\max}(\mathcal{H}\mathbf{X})}{\kappa^3 \varepsilon} \right) \quad (21)$$

for some constant  $C_2 > 0$ .

Theorem 3 shows that the iterates of RAM-FIHT converge to the groundtruth  $\mathbf{X}$  with a linear convergence rate, and the number of required measurements is further reduced from that needed by AM-FIHT. To see this, note that  $\sigma_{\max}(\mathcal{H}\mathbf{X}) \leq \sqrt{n_c n} \|\mathbf{X}\|_\infty$ . If  $\|\mathbf{X}\|_\infty$  is a constant, and select  $\varepsilon = O(n^{-\alpha})$  with a positive constant  $\alpha$ , we have  $L = O(\log(n))$  from (20) and  $m \geq O(r^2 \log^2 n)$  from (21). Compared with the bound of  $O(r n_c^{1/2} n^{1/2} \log^{3/2}(n))$  in Theorem 1, the dependence on  $n$  is significantly reduced to  $\log^2 n$ , while the dependence on  $r$  is worse, from  $r$  to  $r^2$ . Since  $r$  is usually very small, and  $n$  is much larger,  $O(r^2 \log^2 n)$  by Theorem 3 is tighter than  $O(r n_c^{1/2} n^{1/2} \log^{3/2}(n))$  by Theorem 1. Since the degree of freedom of  $\Theta(n_c r)$ , we suspect that the bound could be improved further using better proof techniques than ours.

**Theorem 4.** (*RAM-FIHT with noisy measurements*) Assume  $\mathcal{H}\mathbf{X}$  is  $\mu$ -incoherent and

$$\|\mathbf{N}\|_\infty \leq \frac{\varepsilon_0 \|\mathcal{H}\mathbf{X}\|}{2048\kappa^3 r^{1/2} n_c^{1/2} n}. \quad (22)$$

Let  $L = \left\lceil \varepsilon_0^{-1} \log \left( \frac{\sigma_{\max}(\mathcal{H}\mathbf{X})}{128\kappa^3\varepsilon} \right) \right\rceil$  and  $0 < \varepsilon_0 < \frac{1}{4}$ . Define  $\nu = 2\varepsilon_0 < \frac{1}{2}$ . Then with probability at least  $1 - (3L+3)n_c n^{-2}$  and for any arbitrarily small constant  $\varepsilon > 0$ , the iterates  $\mathbf{L}_l$ 's ( $l = 1, \dots, L$ ) generated by RAM-FIHT with  $\beta = 0$  satisfies

$$\|\mathbf{L}_l - \mathcal{H}\mathbf{X}\|_F \leq \nu^l \|\mathbf{L}_0 - \mathcal{H}\mathbf{X}\|_F + 128n_c^{1/2}n \|\mathbf{N}\|_\infty + 8r^{1/2} \|\mathcal{H}\mathbf{N}\|,$$

$$\text{and } \|\mathbf{L}_L - \mathcal{H}\mathbf{X}\|_F \leq \varepsilon + 128n_c^{1/2}n \|\mathbf{N}\|_\infty + 8r^{1/2} \|\mathcal{H}\mathbf{N}\|, \quad (23)$$

provided that

$$m \geq C_3 \varepsilon_0^{-3} \mu c_s \kappa^6 r^3 \log(n) \log \left( \frac{\sigma_{\max}(\mathcal{H}\mathbf{X})}{\kappa^3 \varepsilon} \right) \quad (24)$$

for some constant  $C_3 > 0$ .

Theorem 4 explores the performance of RAM-FIHT in the noisy case. Note that  $\|\mathcal{H}\mathbf{X}\|_\infty = \|\mathbf{X}\|_\infty$ , and

$$\frac{n_c^{1/2}n}{\mu c_s r} \|\mathcal{H}\mathbf{X}\|_\infty \leq \|\mathcal{H}\mathbf{X}\| \leq n_c^{1/2}n \|\mathcal{H}\mathbf{X}\|_\infty.$$

If  $\mu$  and  $r$  are both constants, (22) implies that  $\|\mathbf{N}\|_\infty$  can be as large as a constant fraction of  $\|\mathbf{X}\|_\infty$ . When the number of observations is at least  $O(r^3 \log^2(n))$ , the error between the ground truth and the iterates returned by RAM-FIHT is controlled by the noise level. To evaluate the optimality of this error bound, consider a special case that  $\mathbf{X}$  is a constant matrix with each entry being  $c$ , and  $\mathbf{N}$  is a constant matrix with each entry being  $-c$ . Then every observation is zero, and the estimated matrix from partial observations by any recovery method would be a zero matrix. Then the recovery error is  $\|\mathcal{H}\mathbf{N}\|_F = \sqrt{n_c n_1 n_2} |c| = \sqrt{n_c n_1 n_2} \|\mathbf{N}\|_\infty$ . The sum of the second and the third term in the right hand side of (23) is bounded by  $(128c_s + 8r^{1/2}) \sqrt{n_c n_1 n_2} \|\mathbf{N}\|_\infty$ . Thus, the error bound of RAM-FIHT is in the same order of the minimum error by any method.

**Comparison with single-channel missing data recovery.** FIHT [7] is a single-channel Hankel matrix completion method. When  $n_c = 1$ , Theorems 1 and 3 reduce to the results in [7]. One can apply FIHT to recover the missing points of each row of  $\mathbf{X}$  and solve  $n_c$  data recovery problems separately. Let  $\mathcal{H}\mathbf{X}_{k*}$  denote the single-channel Hankel matrix constructed from the  $k$ th row of  $\mathbf{X}$ . Suppose  $\mathcal{H}\mathbf{X}_{k*}$  is  $\mu_0$ -incoherent for every  $1 \leq k \leq n_c$ . Then setting  $n_c = 1$  in Theorems 1 and 3 (or using Theorems 1 and 2 in [7]), we know that if each  $\mathcal{H}\mathbf{X}_{k*}$  is recovered separately, the required number of measurements is proportional to  $\sqrt{\mu_0}$  (AM-FIHT) or  $\mu_0$  (RAM-FIHT). Then the total number of observations to recover  $\mathbf{X}$  is proportional to  $n_c \sqrt{\mu_0}$  or  $n_c \mu_0$ . In contrast, the required number of observations by our methods is proportional to  $\sqrt{n_c \mu}$  (AM-FIHT) or  $\mu$  (RAM-FIHT). Thus, the ratio of the number of measurements by our method to FIHT is  $\sqrt{\frac{\mu}{n_c \mu_0}}$  (or  $\frac{\mu}{n_c \mu_0}$ ). We can show that our method only requires a constant fraction of the measurements by using FIHT through the following theorem.

**Theorem 5.**

$$\frac{\mu}{n_c \mu_0} < 1. \quad (25)$$

If it further holds that  $(1 - \delta)|\hat{d}| \leq |d_{k,i}| \leq (1 + \delta)|\hat{d}|, \forall k \in \{1, \dots, n_c\}, i \in \{1, \dots, r\}$  for some  $\delta \in (0, 1)$  and  $\hat{d} \in \mathbb{C}$ , where  $d_{k,i} = \mathbf{r}_i^* \mathbf{s}_1 \mathbf{C}_{k*} \mathbf{l}_i$ , we have

$$\frac{\mu}{n_c \mu_0} \leq \frac{1}{1 + (n_c - 1) \frac{(1 - \delta)^2}{\kappa_L^2 (1 + \delta)^2}}, \quad (26)$$

where  $\kappa_L$  is the conditional number of  $\mathbf{P}_L$  when  $n_c = 1$ .

Theorem 5 indicates that the required number of measurements is reduced when collectively processing  $\mathbf{X}$ . Note that  $\mu_0$  is independent of the amplitude parameters  $d_{k,i}$ 's and depends only on the separations of the frequencies  $f_i$ 's in (12). As a direct corollary of Theorem 2 in [28], if the separation among frequencies  $f_i$ 's is at least  $1/c_s n$ , then  $\mu_0$  is a constant. In contrast,  $\mu$  depends on both  $d_{k,i}$ 's and  $f_i$ 's. (25) shows that  $\mu$  is always less than  $n_c \mu_0$ . Moreover, in the special case that  $d_{k,i}$ 's are all in a small range,  $\mu/(n_c \mu_0)$  can be reduced to approximately  $\kappa_L^2/n_c$  from (26). With well separated frequencies, the maximum and minimum singular values of  $\mathbf{P}_L$  when  $n_c = 1$  are both proportional to  $\sqrt{n_1}$  [28]. That implies  $\kappa_L$  is a constant. Then,  $\kappa_L^2/n_c$  is in the order of  $1/n_c$  for large  $n_c$ , and we have  $\mu/\mu_0 = O(1/n_c)$ . Combining these results with the arguments before Theorem 3, one can see that the required number of measurements is significantly reduced by collective processing.

## VI. NUMERICAL RESULTS

We test the numerical performance of AM-FIHT and RAM-FIHT. The simulations are implemented in MATLAB on a desktop with 3.4 GHz Intel Core i7 and 16 GB memory. In all the experiments, we delete some data points in the datasets and test the recovery performance. We consider three modes of missing data patterns, as illustrated in Fig. 1. Given a data loss percentage,

- Mode 1: Data losses occur randomly and independently across time and channels.
- Mode 2: At randomly selected time instants, the data points in all channels are lost simultaneously.
- Mode 3: Starting from a randomly selected time instant, in half of the channels that are randomly selected, the data points are lost simultaneously and consecutively lost.

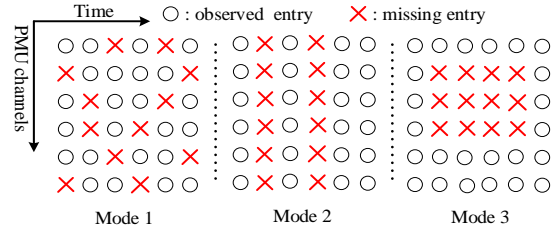


Fig. 1: Three modes of missing data

### A. Numerical experiments on synthetic data

We test the recovery performance on synthetic spectrally sparse signals. Each row of matrix  $\mathbf{X} \in \mathbb{C}^{n_c \times n}$  is a weighted sum of  $r$  sinusoids as shown in (12). Each  $f_i$  is randomly selected from  $(0, 1)$ .  $\tau_i$  is 0 for all  $i$ . The complex coefficient  $d_{k,i}$  has its angle randomly selected from  $(0, 2\pi)$  and its magnitude chosen as  $1 + 10^{a_{k,i}}$ , where  $a_{k,i}$  is randomly selected from  $(0, 1)$ .

1) *(R)AM-FIHT with noiseless measurements*: We first compare the performance of AM-FIHT and RAM-FIHT with noiseless measurements. For RAM-FIHT, instead of dividing the observation set into disjoint subsets, we use the entire observation set in every iteration. Hence, RAM-FIHT differs from AM-FIHT in the trimming step, and the thresholding is set as the ground truth  $\mu$  throughout this section. AM-FIHT is tested with both  $\beta = 0$  and  $\beta = (1-p)^2/5$ , while only  $\beta = (1-p)^2/5$  is tested on RAM-FIHT. An algorithm terminates if

$$\|\mathcal{P}_{\hat{\Omega}}(\mathbf{X}_l - \mathbf{X}_{l-1})\|_F / \|\mathcal{P}_{\hat{\Omega}}(\mathbf{X}_{l-1})\|_F \leq 10^{-6} \quad (27)$$

is satisfied before reaching the maximum iteration number, which is set as 300 here.

Figs. 2, 3, and 4 show the recovery phase transitions of AM-FIHT and RAM-FIHT with missing data patterns following different modes.  $n = 300$ ,  $n_1 = 150$ , and  $n_c = 30$ . The  $x$ -axis is the fraction of observations  $p = \frac{m}{n_c n}$ . The  $y$ -axis is the rank  $r$ . For each fixed  $p$  and  $r$ , we generate 100 independent realizations of synthetic data matrices and data erasures. We say the recovery is successful in a test case if

$$\|\mathcal{P}_{\hat{\Omega}^c}(\mathbf{X}_l - \mathbf{X})\|_F / \|\mathcal{P}_{\hat{\Omega}^c}(\mathbf{X})\|_F < 10^{-3} \quad (28)$$

holds when the algorithm terminates after  $l$ -th iteration, and  $\hat{\Omega}^c$  is the complement of  $\hat{\Omega}$ . A white block corresponds to 100% success, and a black one means failures in all 100 tests.

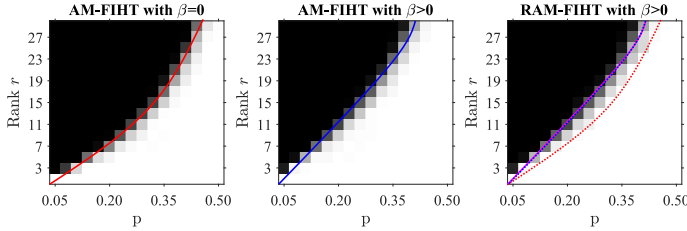


Fig. 2: Phase transition under Mode 1

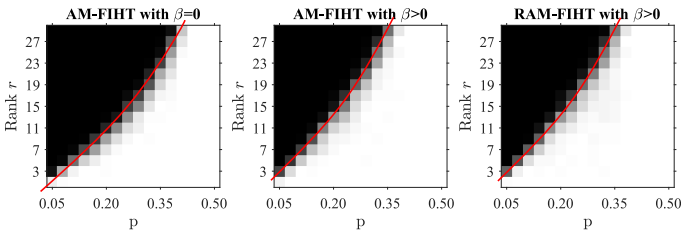


Fig. 3: Phase transition under Mode 2

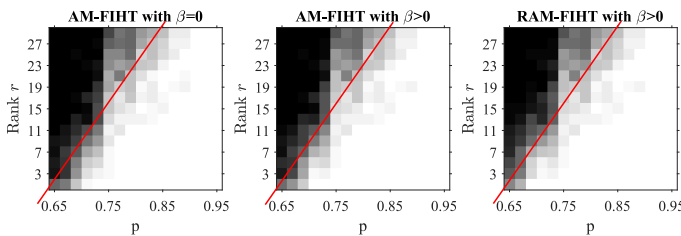


Fig. 4: Phase transition under Mode 3

Auxiliary solid lines (red for AM-FIHT with  $\beta = 0$ , blue for AM-FIHT with  $\beta > 0$ , and magenta for RAM-FIHT) are added in Figs. 2, 3, and 4 to highlight the phase transition. In the subfigures for RAM-FIHT, the phase transition curves

for AM-FIHT are repeated in dotted curves to compare. Both AM-FIHT and RAM-FIHT with  $\beta > 0$  perform very similarly, as the blue dotted line and the magenta solid line coincide in all three modes.

The phase transition threshold of  $\beta > 0$  is higher than that of  $\beta = 0$  for all the modes. The recovery improvement by the heavy-ball step can be intuitively explained as follows. Note that Theorem 2 shows that the heavy ball can speed up the convergence by reducing  $q(0)$  to  $q(\beta)$ . With a certain percentage of data losses, it might hold that  $q(\beta) < 1 < q(0)$ , which indicates that the iterates with  $\beta > 0$  are still convergent, while those with  $\beta = 0$  may be divergent.

One can see from the phase transition lines that the required ratio of observations is approximately linear in  $r$  when other parameters are fixed. Note that the degree of freedom of the signal is  $\Theta(n_c r)$ . Although our bound of the required number of measurements  $O(r^2 \log n)$  in Theorem 3 is not order-wise optimal due to the artefacts of the proof, the required number of measurements in numerical experiments is approximately linear in the degree of freedom.

2) *Comparison with existing algorithms*: We next study the recovery performance with both noiseless and noisy measurements. Here, rank  $r$  is fixed as 15, and  $n = 600$ ,  $n_1 = 300$ ,  $n_c = 20$ . All the other setups remain the same. We compare our methods with FIHT [7] and Singular Value Thresholding (SVT) [6]. Since FIHT recovers the missing points in a single channel, we convert each row of  $\mathbf{X}$  to a Hankel matrix with the size of  $300 \times 301$  and apply FIHT separately. SVT solves the convex NNM problem approximately, and the algorithm is applied on the original observation matrix and the constructed Hankel matrix, respectively. The relative recovery error is calculated as  $\|\mathcal{P}_{\hat{\Omega}^c}(\mathbf{X}_l - \mathbf{X})\|_F / \|\mathcal{P}_{\hat{\Omega}^c}(\mathbf{X})\|_F$ .

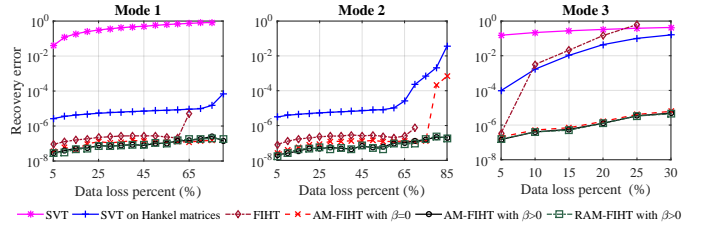


Fig. 5: Performance comparison of recovery methods in noiseless setting

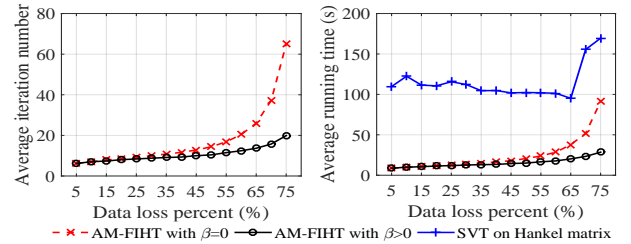


Fig. 6: Running time comparison in Mode 1

Fig. 5 compares the relative recovery error of convergent tests by different methods with noiseless measurements and different data loss patterns. (R)AM-FIHT with a nonzero  $\beta$  performs the best among all the methods. As the original data matrix is not low-rank, SVT fails in all cases. When applied to the constructed Hankel matrix, SVT exhibits better

performance, however, the recovery errors are still much larger compared with (R)AM-FIHT. SVT also needs the much longer running time, as shown in Fig. 6. To achieve the error bound of  $10^{-5}$ , SVT requires around 100 iterations at a time cost of 100 seconds, while AM-FIHT with a nonzero  $\beta$  only takes less than 12 seconds to obtain an error bound of  $10^{-7}$ . With 65% data loss in Mode 1, 13.3% tests of FIHT diverge. In contrast, all tests of AM-FIHT are convergent, even in the case with 75% data loss. A nonzero  $\beta$  also increases the percentage of convergent tests. With 80% of data loss, only 76.7% tests of AM-FIHT with  $\beta = 0$  converge, while all the tests of AM-FIHT with  $\beta > 0$  are convergent. Moreover, AM-FIHT performs much better than FIHT and SVT in Mode 3.

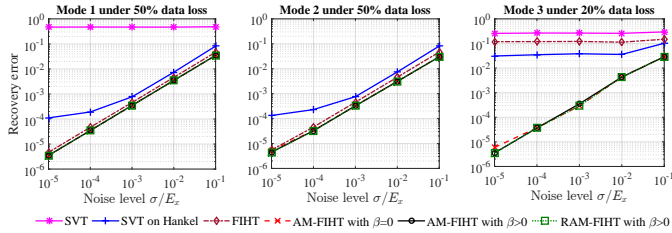


Fig. 7: Performance comparison of recovery methods in noisy setting under 50% data loss in modes 1, 2 and 20% data loss in mode 3

When measurements are noisy, every entry of  $\mathbf{N}$  is independently drawn from Gaussian  $\mathcal{N}(0, \sigma^2)$ , where  $\sigma$  is the standard deviation. Fig. 7 shows the relative recovery error of convergent tests against the relative noise level  $\sigma/E_x$ , where  $E_x$  is the average energy of  $\mathbf{X}$  calculated as  $E_x = \|\mathbf{X}\|_F / \sqrt{n_c \bar{n}}$ . The data loss percentage is fixed as 50% in modes 1, 2 and 20% for mode 3, respectively. In all three modes, AM-FIHT and RAM-FIHT perform similarly and achieve the smallest error among all the methods. The relative recovery error is proportional to the relative noise level, with a ratio between 0.3 to 0.4. FIHT is slightly worse than these two methods in modes 1 and 2, but its performance degrades significantly in mode 3. SVT has a better performance when applied on the Hankel matrix instead of the data matrix, but it is still worse than (R)AM-FIHT.

### B. Numerical experiments on actual PMU data

The low-rank property of the Hankel PMU data matrix is verified on a recorded PMU dataset in Central New York [18]. 11 voltage phasors are measured at a rate of 30 samples per second. Fig. 8 shows the recorded voltage magnitudes and angles of a 10-second dataset that contains a disturbance at 2.3s. Fig. 9 shows the approximation errors of  $\mathbf{X}$  and  $\mathcal{H}\mathbf{X}$  by rank- $r$  matrices with varying  $n_1$  and  $r$ . The approximation error of  $\mathbf{X}$  with the rank- $r$  matrix  $\mathcal{Q}_r(\mathbf{X})$  is defined as  $\|\mathbf{X} - \mathcal{Q}_r(\mathbf{X})\|_F / \|\mathbf{X}\|_F$ , and likewise for  $\mathcal{H}\mathbf{X}$ . One can see from Fig. 9 that all these data matrices can be approximated by rank-8 matrices with negligible errors.

We set  $n_1 = 8$ ,  $r = 8$  and test both  $\beta = 0$  and  $\beta = (1-p)/5$  in the simulation. Fig. 10 shows the percentage of convergent runs out of 100 runs for different algorithms. Fig. 11 compares the average recovery error of convergent runs. Overall, AM-FIHT with  $\beta > 0$  achieves a small recovery error, tolerates a

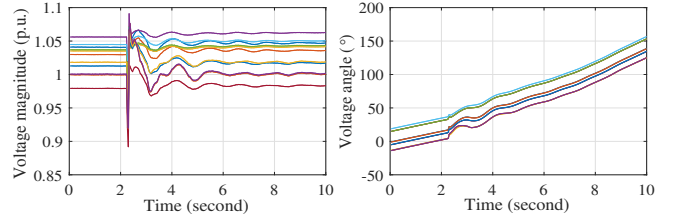


Fig. 8: The measured voltage phasors of 11 channels

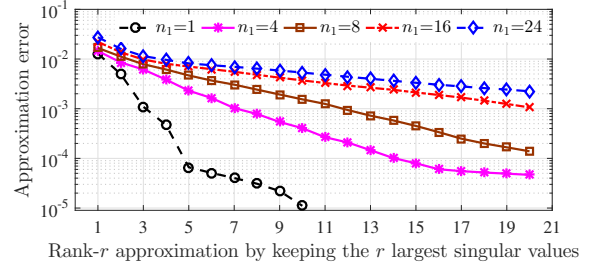


Fig. 9: The approximation errors of the data block and the Hankel matrices

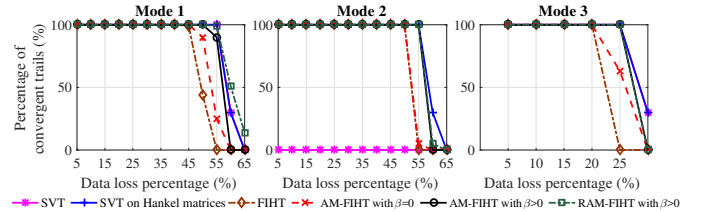


Fig. 10: Percentage of convergent trials of recovery algorithms

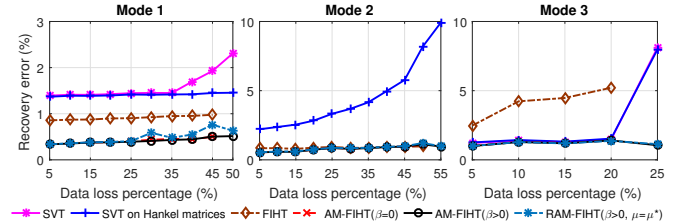


Fig. 11: Performance comparison of recovery algorithms

high data loss rate, and does not require much computation. For example, when the data loss rate is 55% in Mode 2, AM-FIHT with  $\beta = (1-p)/5$  converges every time. The number of iterations is 47.2 on average, and the running time is 0.62 seconds. It takes 4.34 seconds to run 400 iterations of SVT on Hankel matrices. AM-FIHT with  $\beta = 0$  diverges for 95% of the runs. FIHT diverges completely. Similar result to Fig. 6 about the average iteration numbers of AM-FIHT with  $\beta = 0$  and  $\beta > 0$  from respective successful trials is observed as well, thus a small positive  $\beta$  helps improve the convergence rate.

There are minor differences between AM-FIHT ( $\beta > 0$ ) and RAM-FIHT ( $\beta > 0$ ) in mode 1. RAM-FIHT tolerates a slightly higher data loss percentage, while its average recovery error of convergent runs is slightly larger than that of AM-FIHT. AM-FIHT and RAM-FIHT perform similarly in modes 2 and 3.

## VII. CONCLUSION AND DISCUSSIONS

This paper characterizes the intrinsic low-dimensional structures of correlated time series through multi-channel low-rank Hankel matrices. Two iterative hard thresholding algorithms



with linear convergence rates are proposed to solve the non-convex missing data recovery problem. Our bound of the required number of observed entries for successful recovery is  $O(r^2 \log^2 n)$ , significantly smaller than  $O(rn \log^2 n)$  by conventional low-rank matrix completion methods. Our bound is slightly larger than the degree of the freedom  $\Theta(ncr)$ , and we suspect that the bound can be improved with better proof techniques. The convergence rate is proved to be accelerated further by adding a heavy-ball step, which also increases the tolerable missing data percentage numerically.

One motivating application of our methods is power system synchrophasor data recovery. Other applications include array signal processing and MRI image recovery. This paper provides the first analytical characterization of multi-channel Hankel matrix completion methods, while existing works mostly focused on single-channel Hankel matrix recovery. One future direction is to study data recovery from both data losses and corruptions, where partial measurements contain significant errors. The bad data should be first located and removed before recovering the missing points.

#### ACKNOWLEDGEMENT

This research is supported in part by ARO W911NF-17-1-0407, NSF 1508875, EPRI 1007316, IBM Corporation, and the ERC Program of NSF and DoE under the supplement to NSF EEC-1041877 and the CURENT Industry Partnership Program. We thank New York Power Authority for providing recorded PMU datasets.

#### REFERENCES

- [1] "The Netflix Prize," <http://netflixprize.com/>, June 2006.
- [2] N. H. Abbasy and H. M. Ismail, "A unified approach for the optimal pmu location for power system state estimation," *IEEE Trans. Power Syst.*, vol. 24, no. 2, pp. 806–813, 2009.
- [3] Y. Amit, M. Fink, N. Srebro, and S. Ullman, "Uncovering shared structures in multiclass classification," in *Proc. the 24th International Conference on Machine Learning*. ACM, 2007, pp. 17–24.
- [4] A. Balachandrasekaran, V. Magnotta, and M. Jacob, "Recovery of damped exponentials using structured low rank matrix completion," *IEEE Trans. Med. Imag.*, vol. 36, no. 10, pp. 2087–2098, 2017.
- [5] A. Balachandrasekaran, G. Ongie, and M. Jacob, "Accelerated dynamic mri using structured low rank matrix completion," in *2016 IEEE International Conference on Image Processing (ICIP)*, 2016, pp. 1858–1862.
- [6] J.-F. Cai, E. J. Candès, and Z. Shen, "A singular value thresholding algorithm for matrix completion," *SIAM Journal on Optimization*, vol. 20, no. 4, pp. 1956–1982, 2010.
- [7] J.-F. Cai, T. Wang, and K. Wei, "Fast and provable algorithms for spectrally sparse signal reconstruction via low-rank hankel matrix completion," *Applied and Computational Harmonic Analysis*, 2017.
- [8] E. Candès and T. Tao, "The power of convex relaxation: Near-optimal matrix completion," *IEEE Trans. Inf. Theory*, vol. 56, no. 5, pp. 2053–2080, 2010.
- [9] E. J. Candès and C. Fernandez-Granda, "Towards a mathematical theory of super-resolution," *Communications on Pure and Applied Mathematics*, vol. 67, no. 6, pp. 906–956, 2014.
- [10] E. J. Candès and Y. Plan, "Matrix completion with noise," *Proc. IEEE*, vol. 98, no. 6, pp. 925–936, 2010.
- [11] E. J. Candès and B. Recht, "Exact matrix completion via convex optimization," *Foundations of Computational mathematics*, vol. 9, no. 6, pp. 717–772, 2009.
- [12] Y. Chen and Y. Chi, "Robust spectral compressed sensing via structured matrix completion," *IEEE Trans. Inf. Theory*, vol. 60, no. 10, pp. 6576–6601, 2014.
- [13] S. F. Cotter, B. D. Rao, K. Engan, and K. Kreutz-Delgado, "Sparse solutions to linear inverse problems with multiple measurement vectors," *IEEE Trans. Signal Process.*, vol. 53, no. 7, pp. 2477–2488, 2005.
- [14] J. De La Ree, V. Centeno, J. S. Thorp, and A. G. Phadke, "Synchronized phasor measurement applications in power systems," *IEEE Trans. Smart Grid*, vol. 1, no. 1, pp. 20–27, 2010.
- [15] T. Ding, M. Sznajder, and O. I. Camps, "A rank minimization approach to video inpainting," in *IEEE 11th International Conference on Computer Vision*. IEEE, 2007, pp. 1–8.
- [16] M. Fazel, "Matrix rank minimization with applications," Ph.D. dissertation, 2002.
- [17] M. Fazel, T. K. Pong, D. Sun, and P. Tseng, "Hankel matrix rank minimization with applications to system identification and realization," *SIAM Journal on Matrix Analysis and Applications*, vol. 34, no. 3, pp. 946–977, 2013.
- [18] P. Gao, M. Wang, S. G. Ghiocel, J. H. Chow, B. Fardanesh, and G. Stefopoulos, "Missing data recovery by exploiting low-dimensionality in power system synchrophasor measurements," *IEEE Trans. Power Syst.*, vol. 31, no. 2, pp. 1006–1013, 2016.
- [19] D. Goldberg, D. Nichols, B. M. Oki, and D. Terry, "Using collaborative filtering to weave an information tapestry," *Commun. ACM*, vol. 35, no. 12, pp. 61–70, 1992.
- [20] D. Gross, "Recovering low-rank matrices from few coefficients in any basis," *IEEE Trans. Inf. Theory*, vol. 57, no. 3, pp. 1548–1566, 2011.
- [21] J. P. Haldar, "Low-rank modeling of local  $k$ -space neighborhoods (loraks) for constrained mri," *IEEE Trans. Med. Imag.*, vol. 33, no. 3, pp. 668–681, 2014.
- [22] E. C. Hall and R. M. Willett, "Online convex optimization in dynamic environments," *IEEE J. Sel. Topics Signal Process.*, vol. 9, no. 4, pp. 647–662, 2015.
- [23] R. A. Horn and C. R. Johnson, "Matrix analysis," *Cambridge University Press, New York*, 1985.
- [24] P. Jain, P. Netrapalli, and S. Sanghavi, "Low-rank matrix completion using alternating minimization," in *Proceedings of the 45th Annual ACM Symposium on Theory of Computing*, 2013, pp. 665–674.
- [25] K. H. Jin, D. Lee, and J. C. Ye, "A general framework for compressed sensing and parallel mri using annihilating filter based low-rank hankel matrix," *IEEE Transactions on Computational Imaging*, vol. 2, no. 4, pp. 480–495, 2016.
- [26] K. H. Jin and J. C. Ye, "Sparse and low-rank decomposition of a hankel structured matrix for impulse noise removal," *IEEE Trans. Image Process.*, vol. PP, no. 99, pp. 1–1, 2017.
- [27] Y. Li and Y. Chi, "Off-the-grid line spectrum denoising and estimation with multiple measurement vectors," *IEEE Trans. Signal Process.*, vol. 64, no. 5, pp. 1257–1269, March 2016.
- [28] W. Liao and A. Fannjiang, "Music for single-snapshot spectral estimation: Stability and super-resolution," *Applied and Computational Harmonic Analysis*, vol. 40, no. 1, pp. 33–67, 2016.
- [29] J. Liu, A. Eryilmaz, N. B. Shroff, and E. S. Bentley, "Heavy-ball: A new approach to tame delay and convergence in wireless network optimization," in *The 35th Annual IEEE International Conference on Computer Communications*. IEEE, 2016, pp. 1–9.
- [30] M. Lustig, D. Donoho, and J. M. Pauly, "Sparse MRI: The application of compressed sensing for rapid MR imaging," *Magnetic Resonance in Medicine*, vol. 58, no. 6, pp. 1182–1195, 2007.
- [31] J. Ma, P. Zhang, H. Fu, B. Bo, and Z. Dong, "Application of phasor measurement unit on locating disturbance source for low-frequency oscillation," *IEEE Trans. Smart Grid*, vol. 1, no. 3, pp. 340–346, 2010.
- [32] J. K. Merikoski and R. Kumar, "Inequalities for spreads of matrix sums and products," *Applied Mathematics E-Notes*, vol. 4, pp. 150–159, 2004.
- [33] N. Mohammadiha and A. Leijon, "Nonnegative hmm for babble noise derived from speech hmm: Application to speech enhancement," *IEEE Trans. Audio, Speech, Language Process.*, vol. 21, no. 5, pp. 998–1011, 2013.
- [34] N. Mohammadiha, P. Smaragdakis, G. Panahandeh, and S. Doclo., "A state-space approach to dynamic nonnegative matrix factorization," *IEEE Trans. Signal Process.*, vol. 63, no. 4, pp. 949–959, 2015.
- [35] G. J. Mysore, P. Smaragdakis, and B. Raj, "Non-negative hidden markov modeling of audio with application to source separation," in *International Conference on Latent Variable Analysis and Signal Separation*. Springer, 2010, pp. 140–148.
- [36] P. Netrapalli, U. Niranjan, S. Sanghavi, A. Anandkumar, and P. Jain, "Non-convex robust pca," in *NIPS*, 2014, pp. 1107–1115.
- [37] G. Ongie, S. Biswas, and M. Jacob, "Convex recovery of continuous domain piecewise constant images from nonuniform fourier samples," *IEEE Trans. Signal Process.*, vol. 66, no. 1, pp. 236–250, 2018.
- [38] G. Ongie and M. Jacob, "Off-the-grid recovery of piecewise constant images from few fourier samples," *SIAM Journal on Imaging Sciences*, vol. 9, no. 3, pp. 1004–1041, 2016.

- [39] A. Phadke and J. Thorp, *Synchronized phasor measurements and their applications*. Springer, 2008.
- [40] B. T. Polyak, "Introduction to optimization. 1987," *Optimization Software, Inc, New York*.
- [41] L. C. Potter, E. Ertin, J. T. Parker, and M. Cetin, "Sparsity and compressed sensing in radar imaging," *Proc. IEEE*, vol. 98, no. 6, pp. 1006–1020, 2010.
- [42] B. Recht, "A simpler approach to matrix completion," *The Journal of Machine Learning Research*, pp. 3413–3430, 2011.
- [43] P. J. Shin, P. E. Z. Larson, M. A. Ohliger, M. Elad, J. M. Pauly, D. B. Vigneron, and M. Lustig, "Calibrationless parallel imaging reconstruction based on structured low-rank matrix completion," *Magnetic Resonance in Medicine*, vol. 72, no. 4, pp. 959–970, 2014.
- [44] J. A. Tropp, "User-friendly tail bounds for sums of random matrices," *Foundations of computational mathematics*, vol. 12, no. 4, pp. 389–434, 2012.
- [45] K. Usevich and P. Comon, "Hankel low-rank matrix completion: Performance of the nuclear norm relaxation," *IEEE J. Sel. Topics Signal Process.*, vol. 10, no. 4, pp. 637–646, 2016.
- [46] K. Wei, J.-F. Cai, T. F. Chan, and S. Leung, "Guarantees of riemannian optimization for low rank matrix completion," *arXiv preprint arXiv:1603.06610*, 2016.
- [47] Z. Yang and L. Xie, "Exact joint sparse frequency recovery via optimization methods," *IEEE Trans. Signal Process.*, vol. 64, no. 19, pp. 5145–5157, 2016.
- [48] J. C. Ye, J. M. Kim, K. H. Jin, and K. Lee, "Compressive sampling using annihilating filter-based low-rank interpolation," *IEEE Trans. Inf. Theory*, vol. 63, no. 2, pp. 777–801, 2017.
- [49] S. Zhang, Y. Hao, M. Wang, and J. H. Chow, "Multi-channel missing data recovery by exploiting the low-rank hankel structures," in *Proc. CAMSAP*, 2017.

## APPENDIX

### A. Notations and assumptions

We first introduce notations used in the proof. For matrix  $\mathbf{Z}_1 \in \mathbb{C}^{n_c \times n}$ , we define the **Block Hankel Operator**  $\tilde{\mathcal{H}}$  as

$$\tilde{\mathcal{H}}\mathbf{Z}_1 = [\mathcal{H}\mathbf{Z}_1 \quad \mathcal{H}\mathbf{Z}_1 \quad \cdots \quad \mathcal{H}\mathbf{Z}_1] \in \mathbb{C}^{n_c n_1 \times n_c n_2}.$$

$\tilde{\mathcal{H}}\mathbf{Z}_1$  is an  $n_c$ -block Hankel matrix.  $\tilde{\mathcal{H}}^*$  is the adjoint operator of  $\tilde{\mathcal{H}}$ . For any matrix  $\mathbf{Z}_2 \in \mathbb{C}^{n_c n_1 \times n_c n_2}$ ,  $(\tilde{\mathcal{H}}^*\mathbf{Z}_2) \in \mathbb{C}^{n_c \times n}$  satisfies

$$\langle \tilde{\mathcal{H}}^*\mathbf{Z}_2, \mathbf{e}_k \mathbf{e}_t^* \rangle = \sum_{l=0}^{n_c-1} \sum_{k_1+k_2=t+1} \langle \mathbf{Z}_2, \mathbf{e}_{(k_1-1)n_c+k} \mathbf{e}_{k_2+ln_2}^* \rangle.$$

Define  $\tilde{\mathcal{D}}^2 := \tilde{\mathcal{H}}^* \tilde{\mathcal{H}}$ , an operator from an  $n_c \times n$  matrix  $\mathbf{Z}$  to an  $n_c \times n$  matrix with

$$\tilde{\mathcal{D}}^2 \mathbf{Z} = \sum_{t=1}^n \sum_{k=1}^{n_c} n_c w_t \langle \mathbf{Z}, \mathbf{e}_k \mathbf{e}_t^* \rangle \mathbf{e}_k \mathbf{e}_t^*,$$

where  $w_t$  is defined in (15). Then the Moore-Penrose pseudoinverse of  $\tilde{\mathcal{H}}$ , denoted as  $\tilde{\mathcal{H}}^\dagger$ , equals to  $\tilde{\mathcal{D}}^{-2} \tilde{\mathcal{H}}^*$ . Further, we define  $\tilde{\mathcal{G}} = \tilde{\mathcal{H}} \tilde{\mathcal{D}}^{-1}$ , then the adjoint operator of  $\tilde{\mathcal{G}}$  is defined as  $\tilde{\mathcal{G}}^* = \tilde{\mathcal{D}}^{-1} \tilde{\mathcal{H}}^*$ . Additionally,

$$\mathbf{Y} := \tilde{\mathcal{D}} \mathbf{X} \quad \text{and} \quad \mathbf{Y}_l := \tilde{\mathcal{D}} \mathbf{X}_l. \quad (29)$$

For any matrix  $\mathbf{Z} \in \mathbb{C}^{n_c \times n}$ , one can check that  $\|\tilde{\mathcal{H}}\mathbf{X}\| = \sqrt{n_c} \|\mathcal{H}\mathbf{X}\|$  and  $\|\tilde{\mathcal{H}}\mathbf{X}\|_F = \sqrt{n_c} \|\mathcal{H}\mathbf{X}\|_F$ . Immediately,  $\tilde{\mathcal{H}}\mathbf{X}$  and  $\mathcal{H}\mathbf{X}$  share the same conditional number  $\kappa$ . Moreover, it is clear that  $\tilde{\mathcal{G}}$  is a unit operator as  $\tilde{\mathcal{G}}^* \tilde{\mathcal{G}} = \mathcal{I}$ , and  $\tilde{\mathcal{H}}\mathbf{X} = \tilde{\mathcal{G}}\mathbf{Y}$ .

The following proofs will be established on Block Hankel Operator  $\tilde{\mathcal{H}}$ . Consider AM-FIHT in terms of  $\tilde{\mathcal{H}}$ , the initialization can be written as  $\tilde{\mathbf{L}}_0 = p^{-1} \mathcal{Q}_r(\tilde{\mathcal{H}} \mathcal{P}_\Omega(\mathbf{X}))$ . Further, the major steps can be represented as

$$\tilde{\mathbf{W}}_l = \mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{H}}(\mathbf{X}_l + p^{-1} \mathbf{G}_l + \beta \Delta \tilde{\mathbf{W}}), \quad \tilde{\mathbf{L}}_{l+1} = \mathcal{Q}_r(\tilde{\mathbf{W}}_l),$$

where  $\tilde{\mathcal{S}}_l$  is the tangent subspace at  $\tilde{\mathbf{L}}_l$ . The resulting  $\tilde{\mathbf{X}}_l$  returned by AM-FIHT in terms of  $\tilde{\mathcal{H}}$  is given as  $\tilde{\mathcal{H}}^\dagger \tilde{\mathbf{L}}_l$ .

Moreover, by replacing  $\mathcal{H}$  with  $\tilde{\mathcal{H}}$ , AM-FIHT returns the same result as  $\tilde{\mathbf{X}}_l = \mathbf{X}_l$ . In other words, we can show that

$$\tilde{\mathbf{L}}_l = [\mathbf{L}_l \quad \mathbf{L}_l \quad \cdots \quad \mathbf{L}_l], \quad \forall l \geq 0. \quad (30)$$

To see this, it is clear that (30) holds for  $l = 0$  from the definition of  $\tilde{\mathbf{L}}_0$ . Then suppose (30) holds when  $l = k$ . Immediately, we have  $\tilde{\mathbf{W}}_k = [\mathbf{W}_k \quad \mathbf{W}_k \quad \cdots \quad \mathbf{W}_k]$ . Let  $\mathbf{W}_k = \sum_{i=1}^{\min(n_c n_1, n_2)} \sigma_i \mathbf{u}_i \mathbf{v}_i^*$  be the SVD with  $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_{\min(n_c n_1, n_2)}$ . Then for  $l = k + 1$ ,

$$\begin{aligned} \tilde{\mathbf{L}}_{k+1} &= \mathcal{Q}_r(\tilde{\mathbf{W}}_k) = \sum_{i=1}^r \sqrt{n_c} \sigma_i \mathbf{u}_i \frac{1}{\sqrt{n_c}} [\mathbf{v}_i^* \quad \cdots \quad \mathbf{v}_i^*] \\ &= \sum_{i=1}^r [\sigma_i \mathbf{u}_i \mathbf{v}_i^* \quad \cdots \quad \sigma_i \mathbf{u}_i \mathbf{v}_i^*]. \\ &= [\mathbf{L}_{k+1} \quad \mathbf{L}_{k+1} \quad \cdots \quad \mathbf{L}_{k+1}]. \end{aligned}$$

Hence, the connection between  $\mathbf{X}_l$  and  $\tilde{\mathbf{L}}_l$  can be given as

$$\begin{aligned} \|\mathbf{X}_l - \mathbf{X}\|_F &= \|\tilde{\mathcal{D}}^{-1}(\mathbf{Y}_l - \mathbf{Y})\|_F \leq \frac{1}{\sqrt{n_c}} \|\mathbf{Y}_l - \mathbf{Y}\|_F \\ &= \frac{1}{\sqrt{n_c}} \|\tilde{\mathcal{G}}^*(\tilde{\mathbf{L}}_l - \tilde{\mathcal{G}}\mathbf{Y})\|_F \leq \frac{1}{\sqrt{n_c}} \|\tilde{\mathbf{L}}_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F. \end{aligned} \quad (31)$$

For RAM-FIHT, similarly define an  $n_c$ -block matrix  $\tilde{\mathbf{L}}'_l = [\mathbf{L}'_l \quad \mathbf{L}'_l \quad \cdots \quad \mathbf{L}'_l]$ . From the discussion above, one can verify that (30) also holds in RAM-FIHT. Then the SVD of  $\tilde{\mathbf{L}}'_l$  in (30) is  $\tilde{\mathbf{L}}'_l = \tilde{\mathbf{U}}_l \tilde{\Sigma}_l \tilde{\mathbf{V}}_l^* = \mathbf{U}_l (\sqrt{n_c} \Sigma_l) (\frac{1}{\sqrt{n_c}} [\mathbf{V}_l^* \quad \cdots \quad \mathbf{V}_l^*])$ . Then  $\tilde{\mathbf{A}}_l$  and  $\tilde{\mathbf{B}}_l$  are defined as

$$(\tilde{\mathbf{A}}_l)_{i*} = \frac{(\tilde{\mathbf{U}}_l)_{i*}}{\|(\tilde{\mathbf{U}}_l)_{i*}\|} \min \left\{ \|(\tilde{\mathbf{U}}_l)_{i*}\|, \sqrt{\frac{\mu r}{n_c n_1}} \right\}, \quad (32)$$

$$(\tilde{\mathbf{B}}_l)_{i*} = \frac{(\tilde{\mathbf{V}}_l)_{i*}}{\|(\tilde{\mathbf{V}}_l)_{i*}\|} \min \left\{ \|(\tilde{\mathbf{V}}_l)_{i*}\|, \sqrt{\frac{\mu r}{n_c n_2}} \right\}. \quad (33)$$

**Sampling model with replacement.** As shown in [42], due to the duplications, the number of observed entries in a sampling model with replacement is less than or equal to that in a sampling model without replacement. Thus, it is sufficient to study the sampling model with replacement. To distinguish  $\tilde{\Omega}$  in (10), which represents the sampling set without replacement, let  $\Omega$  be  $m$  unions of indices sampled uniformly from set  $\{1, 2, \dots, n_c\} \times \{1, 2, \dots, n\}$  with replacement, and

$$\mathcal{P}_\Omega(\mathbf{Z}_1) = \sum_{a=1}^m \langle \mathbf{Z}_1, \mathbf{e}_{k_a} \mathbf{e}_{t_a}^* \rangle \mathbf{e}_{k_a} \mathbf{e}_{t_a}^*, \quad (34)$$

for any  $\mathbf{Z}_1 \in \mathbb{C}^{n_c \times n}$ . By changing the sampling model, several critical lemmas can be derived from Bernstein Inequality.

**Lemma 1** ([44], Theorem 1.6). *Consider a finite sequence  $\{\mathbf{Z}_k\}$  of independent, random matrices with dimensions  $d_1 \times d_2$ . Assume that such random matrix satisfies*

$$\mathbb{E}(\mathbf{Z}_k) = 0 \quad \text{and} \quad \|\mathbf{Z}_k\| \leq R \quad \text{almost surely.}$$

*Define*

$$\sigma^2 := \max \left\{ \left\| \sum_k \mathbb{E}(\mathbf{Z}_k \mathbf{Z}_k^*) \right\|, \left\| \sum_k \mathbb{E}(\mathbf{Z}_k^* \mathbf{Z}_k) \right\| \right\}.$$

*Then for all  $t \geq 0$ ,*

$$\mathbb{P} \left\{ \left\| \sum_k \mathbf{Z}_k \right\| \geq t \right\} \leq (d_1 + d_2) \exp \left( \frac{-t^2/2}{\sigma^2 + Rt/3} \right).$$

Suppose that  $t \leq \sigma^2/R$ , then the right hand side can be released as  $(d_1 + d_2) \exp(-\frac{3}{8}t^2/\sigma^2)$ . Such kind of manipulation will be adopted in several proofs.

Note that the set of matrices

$$\{\widetilde{\mathbf{H}}_{k,t} | \widetilde{\mathbf{H}}_{k,t} = \frac{1}{\sqrt{n_c w_t}} \widetilde{\mathcal{H}}(\mathbf{e}_k \mathbf{e}_t^*), 1 \leq k \leq n_c, 1 \leq t \leq n\}$$

forms an orthonormal basis of the  $n_c$ -block Hankel matrix, where  $\widetilde{\mathcal{H}}\mathbf{X} = \sum_{k=1}^{n_c} \sum_{t=1}^n \langle \widetilde{\mathbf{H}}_{k,t}, \widetilde{\mathcal{H}}\mathbf{X} \rangle \widetilde{\mathbf{H}}_{k,t}$ . Then, for all  $(k_a, t_a) \in \Omega$ ,  $\mathcal{P}_\Omega$  is also used as the operator

$$\mathcal{P}_\Omega(\mathbf{Z}_2) = \sum_{a=1}^m \langle \mathbf{Z}_2, \widetilde{\mathbf{H}}_{k_a, t_a} \rangle \widetilde{\mathbf{H}}_{k_a, t_a},$$

for any  $\mathbf{Z}_2 \in \mathbb{C}^{n_c n_1 \times n_c n_2}$ . In spite of abuse of notation, the meaning of  $\mathcal{P}_\Omega$  is clear from context. By such definition,  $\mathcal{P}_\Omega(\widetilde{\mathcal{H}}\mathbf{Z}_1) = \widetilde{\mathcal{H}}\mathcal{P}_\Omega(\mathbf{Z}_1)$  for any  $\mathbf{Z}_1 \in \mathbb{C}^{n \times n_c}$ . Additionally,  $\widetilde{\mathbf{H}}_{k,t}$  only has  $n_c w_t$  nonzero entries of magnitude  $1/\sqrt{n_c w_t}$ , so  $\|\widetilde{\mathbf{H}}_{k,t}\|_F = 1$ . The following lemma can be established directly from the definition of incoherence.

**Lemma 2.** Let  $\widetilde{\mathcal{H}}\mathbf{X} = \widetilde{\mathbf{U}}\widetilde{\Sigma}\widetilde{\mathbf{V}}^*$  be the SVD of  $\widetilde{\mathcal{H}}\mathbf{X}$ . Assume  $\widetilde{\mathcal{H}}\mathbf{X}$  is  $\mu$ -incoherent. Then

$$\|\mathbf{e}_{k_1}^* \widetilde{\mathbf{U}}\|^2 \leq \frac{\mu c_s r}{n_c n}, \quad \|\mathbf{e}_{k_2}^* \widetilde{\mathbf{V}}\|^2 \leq \frac{\mu c_s r}{n_c n}, \quad (35)$$

$$\|\mathcal{P}_{\widetilde{\mathbf{U}}}(\widetilde{\mathbf{H}}_{k,t})\|_F^2 \leq \frac{\mu c_s r}{n_c n}, \quad \|\mathcal{P}_{\widetilde{\mathbf{V}}}(\widetilde{\mathbf{H}}_{k,t})\|_F^2 \leq \frac{\mu c_s r}{n_c n}. \quad (36)$$

where  $\mathbf{e}_{k_1}, \mathbf{e}_{k_2}$  are the coordinate unit vectors.

### B. Supporting Lemmas for Theorem 1

We first present some supporting lemmas to prove Theorem 1. Lemma 3 shows that the maximum number of repetitions is bounded by  $O(\log n)$  with high probability in uniform sampling. Lemma 4 derives the properties of  $p^{-1}\mathcal{P}_{\widetilde{\mathcal{S}}}\widetilde{\mathcal{G}}\mathcal{P}_\Omega\widetilde{\mathcal{G}}^*\mathcal{P}_{\widetilde{\mathcal{S}}}$ , and the random operator can be close enough to its mean  $\mathcal{P}_{\widetilde{\mathcal{S}}}\widetilde{\mathcal{G}}\mathcal{P}_\Omega\widetilde{\mathcal{G}}^*\mathcal{P}_{\widetilde{\mathcal{S}}}$  with a significant large amount of observed entries. Lemma 5 connects the angle of two subspaces, represented as  $\|\mathcal{P}_{\widetilde{\mathcal{S}}_l} - \mathcal{P}_{\widetilde{\mathcal{S}}}\|$ , with  $\|\widetilde{\mathbf{L}}_l - \widetilde{\mathcal{G}}\mathbf{Y}\|_F$ , and shows the angle decreases as  $\mathbf{L}_l$  approaches to the ground truth. Lemma 6 indicates the distance between the initial point and ground truth. Lemmas 4 and 6 are built upon Lemmas 5 and 2 of [7], respectively, by extending from single-channel signals to multi-channel signals. When  $n_c = 1$ , Lemmas 4 and 6 are reduced to corresponding lemmas in [7].

**Lemma 3.** Under sampling with replacement model, the maximum number of repetitions of any entry in  $\Omega$  is less than  $3\log(n)$  with probability at least  $1 - n_c n^{-2}$  for  $n \geq 12$ .

**Lemma 4.** Let  $\widetilde{\mathcal{S}}$  be the tangent subspace of  $\widetilde{\mathcal{H}}\mathbf{X}$ . Assume  $\widetilde{\mathcal{H}}\mathbf{X}$  is  $\mu$ -incoherent. Then with  $m \geq 32\mu c_s r \log(n)$ ,

$$\|\mathcal{P}_{\widetilde{\mathcal{S}}}\widetilde{\mathcal{G}}\mathcal{P}_\Omega\widetilde{\mathcal{G}}^*\mathcal{P}_{\widetilde{\mathcal{S}}} - p^{-1}\mathcal{P}_{\widetilde{\mathcal{S}}}\widetilde{\mathcal{G}}\mathcal{P}_\Omega\widetilde{\mathcal{G}}^*\mathcal{P}_{\widetilde{\mathcal{S}}}\| \leq \sqrt{\frac{32\mu c_s r \log(n)}{m}}$$

holds with probability at least  $1 - n_c n^{-2}$ .

**Lemma 5** ([46], Lemma 4.1). Let  $\mathbf{Z}_l$  be a rank- $r$  matrix and  $\mathcal{S}_l$  be the tangent subspace of  $\mathbf{Z}_l$ . If  $\mathbf{Z}$  is also a rank- $r$  matrix and its tangent subspace is denoted as  $\mathcal{S}$ , then

$$\|(\mathcal{I} - \mathcal{P}_{\mathcal{S}_l})(\mathbf{Z}_l - \mathbf{Z})\|_F \leq \frac{\|\mathbf{Z}_l - \mathbf{Z}\|_F^2}{\sigma_{\min}(\mathbf{Z})}, \quad (37)$$

$$\|\mathcal{P}_{\mathcal{S}_l} - \mathcal{P}_{\mathcal{S}}\| \leq \frac{2\|\mathbf{Z}_l - \mathbf{Z}\|_F}{\sigma_{\min}(\mathbf{Z})}. \quad (38)$$

**Lemma 6.** Assume  $\widetilde{\mathcal{H}}\mathbf{X}$  is  $\mu$ -incoherent. With the initial point  $\widetilde{\mathbf{L}}_0 = \mathcal{Q}_r(p^{-1}\widetilde{\mathcal{H}}\mathcal{P}_\Omega(\mathbf{X}))$ , if  $m \geq 16\mu c_s r \log(n)$ , we have

$$\|\widetilde{\mathbf{L}}_0 - \widetilde{\mathcal{H}}\mathbf{X}\| \leq \sqrt{\frac{64\mu c_s r \log(n)}{m}} \|\widetilde{\mathcal{H}}\mathbf{X}\|$$

holds with probability at least  $1 - n_c n^{-2}$ .

### C. Proof of Theorem 1

The proof of Theorem 1 is extended from that of Theorem 3 in [7] with some modifications. The majority of the efforts are devoted to the ‘‘Inductive Step’’ to build the connections between  $\mathbf{W}_{l-1}$  and  $\mathbf{W}_l$  through (43). In (43), the major issue is to bound  $I_1, I_2, I_3$  and  $I_4$ , and (44) and (47) provide critical steps in bounding these items. This inductive step is built upon a similar analysis for  $\mathbf{L}_l$ 's in [7]. Here we study  $\mathbf{W}_l$ 's instead of  $\mathbf{L}_l$ 's since the analysis of Theorem 2 is based on analyzing  $\mathbf{W}_l$ 's. Although Lemma 6 provides the theoretical bound for  $\mathbf{L}_0$  directly, a similar result for  $\mathbf{W}_0$  is lacking. Thus, some efforts to analyze  $\mathbf{W}_0$  is needed in the ‘‘Base Case’’ part. (24) in Theorem 1 is obtained from (51), which provides the theoretical bound for the required number of observations to ensure successful recovery. We include detailed steps as follows for the completeness of the proof.

*Proof of Theorem 1.* As  $\widetilde{\mathbf{L}}_{l+1} = \mathcal{Q}_r(\widetilde{\mathbf{W}}_l)$ ,  $\widetilde{\mathbf{L}}_{l+1}$  is the best rank- $r$  approximation to  $\widetilde{\mathbf{W}}_l$ . Then we have

$$\begin{aligned} \|\widetilde{\mathbf{L}}_{l+1} - \widetilde{\mathcal{G}}\mathbf{Y}\|_F &\leq \|\widetilde{\mathbf{W}}_l - \widetilde{\mathbf{L}}_{l+1}\|_F + \|\widetilde{\mathbf{W}}_l - \widetilde{\mathcal{G}}\mathbf{Y}\|_F \\ &\leq 2\|\widetilde{\mathbf{W}}_l - \widetilde{\mathcal{G}}\mathbf{Y}\|_F. \end{aligned} \quad (39)$$

Therefore, it is sufficient to bound  $\|\widetilde{\mathbf{W}}_l - \widetilde{\mathcal{G}}\mathbf{Y}\|_F$ . Lemma 3 suggests that with probability at least  $1 - n_c n^{-2}$ ,

$$\|\mathcal{P}_\Omega\| \leq 3\log(n) \quad (40)$$

holds. Lemma 4 suggests as long as  $m \geq 32\varepsilon_0^{-2}\mu c_s r \log(n)$ ,

$$\|\mathcal{P}_{\widetilde{\mathcal{S}}}\widetilde{\mathcal{G}}\mathcal{P}_\Omega\widetilde{\mathcal{G}}^*\mathcal{P}_{\widetilde{\mathcal{S}}} - p^{-1}\mathcal{P}_{\widetilde{\mathcal{S}}}\widetilde{\mathcal{G}}\mathcal{P}_\Omega\widetilde{\mathcal{G}}^*\mathcal{P}_{\widetilde{\mathcal{S}}}\| \leq \varepsilon_0 \quad (41)$$

holds with probability at least  $1 - n_c n^{-2}$ .

Now we will show that the following inequality holds by mathematical induction,

$$\frac{\|\widetilde{\mathbf{W}}_k - \widetilde{\mathcal{G}}\mathbf{Y}\|_F}{\sigma_{\min}(\widetilde{\mathcal{G}}\mathbf{Y})} \leq \frac{p^{1/2}\varepsilon_0}{12\log(n)(1 + \varepsilon_0)}. \quad (42)$$

**Inductive Step:** Suppose (42) holds for  $k = l - 1$ . Recall that

$$\begin{aligned} \widetilde{\mathbf{W}}_l &= \mathcal{P}_{\widetilde{\mathcal{S}}_l} \widetilde{\mathcal{H}}(\mathbf{X}_l + p^{-1}\mathcal{P}_\Omega(\mathbf{X} - \mathbf{X}_l)) \\ &= \mathcal{P}_{\widetilde{\mathcal{S}}_l} \widetilde{\mathcal{G}}(\mathbf{Y}_l + p^{-1}\mathcal{P}_\Omega(\mathbf{Y} - \mathbf{Y}_l)). \end{aligned}$$

Then, for all  $l \geq 1$ , we have

$$\begin{aligned}
& \left\| \widetilde{\mathbf{W}}_l - \widetilde{\mathcal{G}}\mathbf{Y} \right\|_F = \left\| \mathcal{P}_{\widetilde{\mathcal{S}}_l} \widetilde{\mathcal{G}}(\mathbf{Y}_l + p^{-1}\mathcal{P}_\Omega(\mathbf{Y} - \mathbf{Y}_l) - \widetilde{\mathcal{G}}\mathbf{Y}) \right\|_F \\
& = \left\| \mathcal{P}_{\widetilde{\mathcal{S}}_l} \widetilde{\mathcal{G}}\mathbf{Y} - \widetilde{\mathcal{G}}\mathbf{Y} + (\mathcal{P}_{\widetilde{\mathcal{S}}_l} \widetilde{\mathcal{G}} - p^{-1}\mathcal{P}_{\widetilde{\mathcal{S}}_l} \widetilde{\mathcal{G}}\mathcal{P}_\Omega)(\mathbf{Y}_l - \mathbf{Y}) \right\|_F \\
& \stackrel{(a)}{\leq} \left\| (\mathcal{I} - \mathcal{P}_{\widetilde{\mathcal{S}}_l})(\widetilde{\mathbf{L}}_l - \widetilde{\mathcal{G}}\mathbf{Y}) \right\|_F \\
& \quad + \left\| (\mathcal{P}_{\widetilde{\mathcal{S}}_l} \widetilde{\mathcal{G}}\widetilde{\mathcal{G}}^* - p^{-1}\mathcal{P}_{\widetilde{\mathcal{S}}_l} \widetilde{\mathcal{G}}\mathcal{P}_\Omega \widetilde{\mathcal{G}}^*)(\mathcal{P}_{\widetilde{\mathcal{S}}_l} \widetilde{\mathbf{W}}_{l-1} - \widetilde{\mathcal{G}}\mathbf{Y}) \right\|_F \\
& \leq \left\| (\mathcal{I} - \mathcal{P}_{\widetilde{\mathcal{S}}_l})(\widetilde{\mathbf{L}}_l - \widetilde{\mathcal{G}}\mathbf{Y}) \right\|_F + \left\| \mathcal{P}_{\widetilde{\mathcal{S}}_l} \widetilde{\mathcal{G}}\widetilde{\mathcal{G}}^*(\mathcal{I} - \mathcal{P}_{\widetilde{\mathcal{S}}_l})(\widetilde{\mathbf{L}}_l - \widetilde{\mathcal{G}}\mathbf{Y}) \right\|_F \\
& \quad + \left\| (\mathcal{P}_{\widetilde{\mathcal{S}}_l} \widetilde{\mathcal{G}}\widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}_l} - p^{-1}\mathcal{P}_{\widetilde{\mathcal{S}}_l} \widetilde{\mathcal{G}}\mathcal{P}_\Omega \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}_l})(\widetilde{\mathbf{W}}_{l-1} - \widetilde{\mathcal{G}}\mathbf{Y}) \right\|_F \\
& \quad + p^{-1} \left\| \mathcal{P}_{\widetilde{\mathcal{S}}_l} \widetilde{\mathcal{G}}\mathcal{P}_\Omega \widetilde{\mathcal{G}}^*(\mathcal{I} - \mathcal{P}_{\widetilde{\mathcal{S}}_l})(\widetilde{\mathbf{L}}_l - \widetilde{\mathcal{G}}\mathbf{Y}) \right\|_F \\
& := I_1 + I_2 + I_3 + I_4, \tag{43}
\end{aligned}$$

where (a) holds since  $\widetilde{\mathbf{L}}_l = \mathcal{P}_{\widetilde{\mathcal{S}}_l} \widetilde{\mathbf{W}}_{l-1}$ . By applying (37),

$$\begin{aligned}
& I_1 + I_2 + I_4 \\
& \leq \frac{2\|\widetilde{\mathbf{L}}_l - \widetilde{\mathcal{G}}\mathbf{Y}\|_F^2}{\sigma_{\min}(\widetilde{\mathcal{G}}\mathbf{Y})} + p^{-1} \left\| \mathcal{P}_\Omega \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}_l} \right\| \frac{\|\widetilde{\mathbf{L}}_l - \widetilde{\mathcal{G}}\mathbf{Y}\|_F^2}{\sigma_{\min}(\widetilde{\mathcal{G}}\mathbf{Y})} \\
& \leq \frac{8\|\widetilde{\mathbf{W}}_{l-1} - \widetilde{\mathcal{G}}\mathbf{Y}\|_F^2}{\sigma_{\min}(\widetilde{\mathcal{G}}\mathbf{Y})} + 4p^{-1} \left\| \mathcal{P}_\Omega \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}_l} \right\| \frac{\|\widetilde{\mathbf{W}}_{l-1} - \widetilde{\mathcal{G}}\mathbf{Y}\|_F^2}{\sigma_{\min}(\widetilde{\mathcal{G}}\mathbf{Y})}. \tag{44}
\end{aligned}$$

Next, we will bound  $\|\mathcal{P}_\Omega \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}_l}\|$ . For any  $\mathbf{Z} \in \mathbb{C}^{n_c n_1 \times n_c n_2}$ ,

$$\begin{aligned}
& \|\mathcal{P}_\Omega \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}_l}(\mathbf{Z})\|^2 = \langle \mathcal{P}_\Omega \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}_l}(\mathbf{Z}), \mathcal{P}_\Omega \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}_l}(\mathbf{Z}) \rangle \\
& \stackrel{(b)}{\leq} 3 \log(n) \langle \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}_l}(\mathbf{Z}), \mathcal{P}_\Omega \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}_l}(\mathbf{Z}) \rangle \\
& = 3 \log(n) \langle \mathbf{Z}, \mathcal{P}_{\widetilde{\mathcal{S}}_l} \widetilde{\mathcal{G}} \mathcal{P}_\Omega \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}_l}(\mathbf{Z}) \rangle \\
& \stackrel{(c)}{\leq} 3 \log(n) (1 + \varepsilon_0) p \|\mathbf{Z}\|_F^2.
\end{aligned} \tag{45}$$

where (b) holds due to (40), and (c) holds due to (41). Hence,

$$\left\| \mathcal{P}_{\widetilde{\mathcal{S}}_l} \widetilde{\mathcal{G}} \mathcal{P}_\Omega \right\| = \left\| \mathcal{P}_\Omega \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}_l} \right\| \leq \sqrt{3 \log(n) (1 + \varepsilon_0) p},$$

and

$$\begin{aligned}
& \left\| \mathcal{P}_\Omega \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}_l} \right\| \leq \left\| \mathcal{P}_\Omega \widetilde{\mathcal{G}}^* (\mathcal{P}_{\widetilde{\mathcal{S}}_l} - \mathcal{P}_{\widetilde{\mathcal{S}}}) \right\| + \left\| \mathcal{P}_\Omega \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}} \right\| \\
& \stackrel{(a)}{\leq} 3 \log(n) \frac{2\|\widetilde{\mathbf{L}}_l - \widetilde{\mathcal{G}}\mathbf{Y}\|_F}{\sigma_{\min}(\widetilde{\mathcal{G}}\mathbf{Y})} + \left\| \mathcal{P}_\Omega \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}} \right\| \\
& \stackrel{(b)}{\leq} 3 \log(n) \frac{p^{1/2} \varepsilon_0}{3 \log(n) (1 + \varepsilon_0)} + \sqrt{3 \log(n) (1 + \varepsilon_0) p} \\
& \leq 3 \log(n) (1 + \varepsilon_0) p^{1/2},
\end{aligned} \tag{46}$$

where (a) holds due to (38) and (40), and (b) holds due to (39) and the inductive hypothesis.

Hence,  $I_1 + I_2 + I_4 \leq 2\varepsilon_0 \|\widetilde{\mathbf{W}}_{l-1} - \widetilde{\mathcal{G}}\mathbf{Y}\|_F$ . Moreover,

$$\begin{aligned}
& \left\| \mathcal{P}_{\widetilde{\mathcal{S}}_l} \widetilde{\mathcal{G}} \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}_l} - p^{-1} \mathcal{P}_{\widetilde{\mathcal{S}}_l} \widetilde{\mathcal{G}} \mathcal{P}_\Omega \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}_l} \right\| \\
& \leq \left\| \mathcal{P}_{\widetilde{\mathcal{S}}} \widetilde{\mathcal{G}} \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}} - p^{-1} \mathcal{P}_{\widetilde{\mathcal{S}}} \widetilde{\mathcal{G}} \mathcal{P}_\Omega \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}} \right\| + \left\| (\mathcal{P}_{\widetilde{\mathcal{S}}} - \mathcal{P}_{\widetilde{\mathcal{S}}_l}) \widetilde{\mathcal{G}} \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}_l} \right\| \\
& \quad + \left\| \mathcal{P}_{\widetilde{\mathcal{S}}} \widetilde{\mathcal{G}} \widetilde{\mathcal{G}}^* (\mathcal{P}_{\widetilde{\mathcal{S}}} - \mathcal{P}_{\widetilde{\mathcal{S}}_l}) \right\| + \left\| p^{-1} (\mathcal{P}_{\widetilde{\mathcal{S}}} - \mathcal{P}_{\widetilde{\mathcal{S}}_l}) \widetilde{\mathcal{G}} \mathcal{P}_\Omega \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}_l} \right\| \\
& \quad + \left\| p^{-1} \mathcal{P}_{\widetilde{\mathcal{S}}} \widetilde{\mathcal{G}} \mathcal{P}_\Omega \widetilde{\mathcal{G}}^* (\mathcal{P}_{\widetilde{\mathcal{S}}} - \mathcal{P}_{\widetilde{\mathcal{S}}_l}) \right\| \\
& \stackrel{(c)}{\leq} \varepsilon_0 + \frac{2\|\widetilde{\mathbf{L}}_l - \widetilde{\mathcal{G}}\mathbf{Y}\|_F}{\sigma_{\min}(\widetilde{\mathcal{G}}\mathbf{Y})} \left( 2 + p^{-1} \|\mathcal{P}_\Omega \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}_l}\| + p^{-1} \|\mathcal{P}_{\widetilde{\mathcal{S}}} \widetilde{\mathcal{G}} \mathcal{P}_\Omega\| \right) \\
& \stackrel{(d)}{\leq} 4\varepsilon_0. \tag{47}
\end{aligned}$$

where (c) comes from Lemma 5, and (d) comes from (42) and (46). Then,  $I_3$  can be bounded as

$$I_3 \leq 4\varepsilon_0 \|\widetilde{\mathbf{W}}_{l-1} - \widetilde{\mathcal{G}}\mathbf{Y}\|_F. \tag{48}$$

By putting pieces together, we have

$$\|\widetilde{\mathbf{W}}_l - \widetilde{\mathcal{G}}\mathbf{Y}\|_F \leq \nu \|\widetilde{\mathbf{W}}_{l-1} - \widetilde{\mathcal{G}}\mathbf{Y}\|_F. \tag{49}$$

Hence, (42) still holds for  $k = l$ .

**Base Case:** Let us assume

$$\|\widetilde{\mathbf{L}}_0 - \widetilde{\mathcal{G}}\mathbf{Y}\|_F \leq \frac{p^{1/2} \varepsilon_0}{6 \log(n) (1 + \varepsilon_0)}. \tag{50}$$

Then, similar to  $I_1 + I_2 + I_3 + I_4$  in (43), we have

$$\begin{aligned}
& \left\| \widetilde{\mathbf{W}}_0 - \widetilde{\mathcal{G}}\mathbf{Y} \right\|_F = \left\| \mathcal{P}_{\widetilde{\mathcal{S}}_0} \widetilde{\mathcal{G}}(\mathbf{Y}_0 + p^{-1}\mathcal{P}_\Omega(\mathbf{Y} - \mathbf{Y}_0)) - \widetilde{\mathcal{G}}\mathbf{Y} \right\|_F \\
& \leq \left\| (\mathcal{I} - \mathcal{P}_{\widetilde{\mathcal{S}}_0})(\widetilde{\mathbf{L}}_0 - \widetilde{\mathcal{G}}\mathbf{Y}) \right\|_F + \left\| \mathcal{P}_{\widetilde{\mathcal{S}}_0} \widetilde{\mathcal{G}} \widetilde{\mathcal{G}}^*(\mathcal{I} - \mathcal{P}_{\widetilde{\mathcal{S}}_0})(\widetilde{\mathbf{L}}_0 - \widetilde{\mathcal{G}}\mathbf{Y}) \right\|_F \\
& \quad + \left\| (\mathcal{P}_{\widetilde{\mathcal{S}}_0} \widetilde{\mathcal{G}} \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}_0} - p^{-1} \mathcal{P}_{\widetilde{\mathcal{S}}_0} \widetilde{\mathcal{G}} \mathcal{P}_\Omega \widetilde{\mathcal{G}}^* \mathcal{P}_{\widetilde{\mathcal{S}}_0})(\widetilde{\mathbf{L}}_0 - \widetilde{\mathcal{G}}\mathbf{Y}) \right\|_F \\
& \quad + p^{-1} \left\| \mathcal{P}_{\widetilde{\mathcal{S}}_0} \widetilde{\mathcal{G}} \mathcal{P}_\Omega \widetilde{\mathcal{G}}^*(\mathcal{I} - \mathcal{P}_{\widetilde{\mathcal{S}}_0})(\widetilde{\mathbf{L}}_0 - \widetilde{\mathcal{G}}\mathbf{Y}) \right\|_F \\
& \leq 5\varepsilon_0 \|\widetilde{\mathbf{L}}_0 - \widetilde{\mathcal{G}}\mathbf{Y}\|_F.
\end{aligned}$$

Since  $\varepsilon_0 \in (0, 1/10)$ , we have

$$\|\widetilde{\mathbf{W}}_0 - \widetilde{\mathcal{G}}\mathbf{Y}\|_F \leq \frac{1}{2} \|\widetilde{\mathbf{L}}_0 - \widetilde{\mathcal{G}}\mathbf{Y}\|_F \leq \frac{p^{1/2} \varepsilon_0}{12 \log(n) (1 + \varepsilon_0)},$$

which completes the induction part.

Then the only thing is to check the assumption (50). Using Lemma 6 and  $\|\widetilde{\mathbf{L}}_0 - \widetilde{\mathcal{G}}\mathbf{Y}\|_F \leq \sqrt{2r} \|\widetilde{\mathbf{L}}_0 - \widetilde{\mathcal{G}}\mathbf{Y}\|$ , with probability at least  $1 - n_c n^{-2}$ ,

$$\frac{\|\widetilde{\mathbf{L}}_0 - \widetilde{\mathcal{G}}\mathbf{Y}\|_F}{\sigma_{\min}(\widetilde{\mathcal{G}}\mathbf{Y})} \leq \frac{\sqrt{2r} \|\widetilde{\mathbf{L}}_0 - \widetilde{\mathcal{G}}\mathbf{Y}\|}{\sigma_{\min}(\widetilde{\mathcal{G}}\mathbf{Y})} = \kappa \sqrt{\frac{128\mu_c s r^2 \log(n)}{m}}.$$

Therefore, to guarantee (50), we need

$$\kappa \sqrt{\frac{128\mu_c s r^2 \log(n)}{m}} \leq \frac{p^{1/2} \varepsilon_0}{6 \log(n) (1 + \varepsilon_0)}, \tag{51}$$

That is  $m \geq C_1 (1 + \varepsilon_0) \varepsilon_0^{-1} n_c^{1/2} \mu^{1/2} c_s^{1/2} \kappa r n^{1/2} \log^{3/2}(n)$  with  $C_1 = 48\sqrt{2}$ .

Hence, with probability at least  $1 - (2l+1)n_c n^{-2}$ , from (49),

$$\begin{aligned}
\|\mathbf{Y}_l - \mathbf{Y}\|_F & = \|\widetilde{\mathcal{G}}^*(\widetilde{\mathbf{L}}_l - \widetilde{\mathcal{G}}\mathbf{Y})\|_F \leq \|\widetilde{\mathbf{L}}_l - \widetilde{\mathcal{G}}\mathbf{Y}\|_F \\
& \leq 2\|\widetilde{\mathbf{W}}_{l-1} - \widetilde{\mathcal{G}}\mathbf{Y}\|_F \leq 2\nu^{l-1} \|\widetilde{\mathbf{W}}_0 - \widetilde{\mathcal{G}}\mathbf{Y}\|_F \\
& \leq \nu^{l-1} \|\widetilde{\mathbf{L}}_0 - \widetilde{\mathcal{G}}\mathbf{Y}\|_F, \tag{52}
\end{aligned}$$

where  $\mathbf{Y}_l = \widetilde{\mathcal{G}}^* \widetilde{\mathbf{L}}_l$ ,  $\widetilde{\mathcal{G}}^* \widetilde{\mathcal{G}} = \mathcal{I}$  and  $\|\widetilde{\mathcal{G}}^*\| \leq 1$ .  $\square$

## D. Proof of Theorem 2

First, we extend the eigenvalues and eigenvectors of linear operators on vector spaces to the eigenvalues and eigenmatrices of linear operators on matrix spaces, defined as follows.

**Definition 2.** Let  $\mathcal{A}$  denote a linear operator from  $\mathbb{C}^{l_1 \times l_2}$  to  $\mathbb{C}^{l_1 \times l_2}$ , for any matrix  $\mathbf{M}$  in the space and  $\mathbf{M} \neq \mathbf{0}$ , if  $\mathcal{A}\mathbf{M} = \lambda\mathbf{M}$  holds, then  $\mathbf{M}$  is one eigenmatrix of operator  $\mathcal{A}$ , and  $\lambda$  is the corresponding eigenvalue.

Let  $\mathcal{L}$  denote the following linear operator on the matrix space  $\mathbb{C}^{n_c n_1 \times n_c n_2}$ ,

$$\mathcal{L} = \mathcal{P}_{\tilde{\mathcal{S}}}\tilde{\mathcal{G}}\tilde{\mathcal{G}}^*\mathcal{P}_{\tilde{\mathcal{S}}} - p^{-1}\mathcal{P}_{\tilde{\mathcal{S}}}\tilde{\mathcal{G}}\mathcal{P}_{\Omega}\tilde{\mathcal{G}}^*\mathcal{P}_{\tilde{\mathcal{S}}}.$$

We first introduce Lemmas 7, 8 and 9 that are useful in the proof of Theorem 2.

**Lemma 7.** *Suppose that for any  $\epsilon > 0$ , there always exists an integer  $s_\epsilon$  such that for any integer  $k \geq 0$ , the iterates  $\tilde{\mathbf{W}}_{s_\epsilon+k}$  generated by AM-FIHT satisfy  $\|\tilde{\mathbf{W}}_{s_\epsilon+k} - \tilde{\mathcal{G}}\mathbf{Y}\|_F \leq \epsilon$ . Then with any  $l > s_\epsilon + 1$ , the updated rule can be denoted as*

$$\begin{bmatrix} \tilde{\mathbf{W}}_l - \tilde{\mathcal{G}}\mathbf{Y} \\ \tilde{\mathbf{W}}_{l-1} - \tilde{\mathcal{G}}\mathbf{Y} \end{bmatrix} = \begin{bmatrix} \mathcal{L}(\tilde{\mathbf{W}}_{l-1} - \tilde{\mathcal{G}}\mathbf{Y}) + \beta\mathcal{P}_{\tilde{\mathcal{S}}}(\tilde{\mathbf{W}}_{l-1} - \tilde{\mathbf{W}}_{l-2}) \\ \tilde{\mathbf{W}}_{l-1} - \tilde{\mathcal{G}}\mathbf{Y} \end{bmatrix} + \tilde{\mathbf{Z}}_{l-1},$$

where

$$\|\tilde{\mathbf{Z}}_{l-1}\|_F = o\left(\left\|\begin{bmatrix} \tilde{\mathbf{W}}_{l-1} - \tilde{\mathcal{G}}\mathbf{Y} \\ \tilde{\mathbf{W}}_{l-2} - \tilde{\mathcal{G}}\mathbf{Y} \end{bmatrix}\right\|_F\right).$$

**Lemma 8.** *All the eigenvalues of operator  $\mathcal{L}$  are real numbers.*

**Lemma 9** ([40], Lemma 2.1). *Let  $\mathcal{A}$  be a linear operator from  $\mathbb{C}^{l_1 \times l_2}$  to  $\mathbb{C}^{l_1 \times l_2}$ , and let  $\lambda_1, \dots, \lambda_n$  be its eigenvalues, let  $\rho(\mathcal{A}) = \max_{1 \leq i \leq n} |\lambda_i|$ , if  $\rho(\mathcal{A}) < 1$ , then there exists some constant  $c(\delta)$  such that  $\|\mathcal{A}^k\| \leq c(\delta)(\rho(\mathcal{A}) + \delta)^k$  holds for all integers  $k$ , where  $0 < \delta < 1 - \rho(\mathcal{A})$ .*

*Proof of Theorem 2.* First, we claim that AM-FIHT is still convergent with a small  $\beta \in (0, 1)$ . Based on the proof of Theorem 1,

$$\|\tilde{\mathbf{W}}_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F \leq \nu \|\tilde{\mathbf{W}}_{l-1} - \tilde{\mathcal{G}}\mathbf{Y}\|_F,$$

where  $\nu = 6\epsilon_0, 0 < \epsilon_0 < 1/10$ . A loose bound from a direct derivation with  $\beta \neq 0$  is

$$\|\tilde{\mathbf{W}}_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F \leq (\nu + \beta)\|\tilde{\mathbf{W}}_{l-1} - \tilde{\mathcal{G}}\mathbf{Y}\|_F + \beta\|\tilde{\mathbf{W}}_{l-2} - \tilde{\mathcal{G}}\mathbf{Y}\|_F.$$

Thus if  $\nu + 2\beta < 1$ , i.e.,  $\beta \in (0, \frac{1}{5})$ , the iteration is still convergent. Thus for any  $\epsilon > 0$ , we can always find such an  $l$  that  $\|\tilde{\mathbf{W}}_{l-2+k} - \tilde{\mathcal{G}}\mathbf{Y}\|_F \leq \epsilon, \forall k \geq 0$ . Then following Lemma 7, if we ignore  $\tilde{\mathbf{Z}}_{l-1}$ , then

$$\tilde{\mathbf{W}}_l - \tilde{\mathcal{G}}\mathbf{Y} = \mathcal{L}(\tilde{\mathbf{W}}_{l-1} - \tilde{\mathcal{G}}\mathbf{Y}) + \beta\mathcal{P}_{\tilde{\mathcal{S}}}(\tilde{\mathbf{W}}_{l-1} - \tilde{\mathbf{W}}_{l-2}).$$

Thus  $\mathcal{P}_{\tilde{\mathcal{S}}}(\tilde{\mathbf{W}}_l - \tilde{\mathcal{G}}\mathbf{Y}) = \tilde{\mathbf{W}}_l - \tilde{\mathcal{G}}\mathbf{Y}$ . With  $\mathcal{P}_{\tilde{\mathcal{S}}}(\tilde{\mathcal{G}}\mathbf{Y}) = \tilde{\mathcal{G}}\mathbf{Y}$ , we have  $\mathcal{P}_{\tilde{\mathcal{S}}}(\tilde{\mathbf{W}}_l) = \tilde{\mathbf{W}}_l$ . The update rule of AM-FIHT can be further simplified as

$$\begin{aligned} \begin{bmatrix} \tilde{\mathbf{W}}_l - \tilde{\mathcal{G}}\mathbf{Y} \\ \tilde{\mathbf{W}}_{l-1} - \tilde{\mathcal{G}}\mathbf{Y} \end{bmatrix} &= \begin{bmatrix} \mathcal{L}(\tilde{\mathbf{W}}_{l-1} - \tilde{\mathcal{G}}\mathbf{Y}) + \beta(\tilde{\mathbf{W}}_{l-1} - \tilde{\mathbf{W}}_{l-2}) \\ \tilde{\mathbf{W}}_{l-1} - \tilde{\mathcal{G}}\mathbf{Y} \end{bmatrix} \\ &:= \tilde{\mathcal{L}} \begin{bmatrix} \tilde{\mathbf{W}}_{l-1} - \tilde{\mathcal{G}}\mathbf{Y} \\ \tilde{\mathbf{W}}_{l-2} - \tilde{\mathcal{G}}\mathbf{Y} \end{bmatrix}. \end{aligned}$$

Following Lemma 4, we have  $\|\mathcal{L}\| < 1$ , if  $m > 32\mu_c s r \log(n)$ . Based on the definitions of  $\rho(\mathcal{L})$  and  $\|\mathcal{L}\|$ , we have  $\rho(\mathcal{L}) \leq \|\mathcal{L}\| < 1$ . Our aim is to prove  $\rho(\tilde{\mathcal{L}}) < \rho(\mathcal{L})$ . Let  $\lambda$  denote one nonzero eigenvalue of  $\tilde{\mathcal{L}}$ , the corresponding eigenmatrix is  $\begin{bmatrix} \mathbf{M}_1 \\ \mathbf{M}_2 \end{bmatrix}$ , then

$$\tilde{\mathcal{L}} \begin{bmatrix} \mathbf{M}_1 \\ \mathbf{M}_2 \end{bmatrix} = \begin{bmatrix} \mathcal{L}\mathbf{M}_1 + \beta(\mathbf{M}_1 - \mathbf{M}_2) \\ \mathbf{M}_1 \end{bmatrix} = \lambda \begin{bmatrix} \mathbf{M}_1 \\ \mathbf{M}_2 \end{bmatrix},$$

we have  $\mathbf{M}_1 = \lambda\mathbf{M}_2, \mathcal{L}(\mathbf{M}_1) + \beta\mathbf{M}_1 - \beta\mathbf{M}_2 = \lambda\mathbf{M}_1$ . Therefore,  $\lambda\mathcal{L}(\mathbf{M}_2) + \lambda\beta\mathbf{M}_2 - \beta\mathbf{M}_2 = \lambda^2\mathbf{M}_2$ . With  $\lambda \neq 0$ ,

$$\mathcal{L}(\mathbf{M}_2) = (\lambda - \beta + \beta/\lambda)\mathbf{M}_2 := \eta_i\mathbf{M}_2,$$

thus  $\mathbf{M}_2$  is an eigenmatrix of operator  $\mathcal{L}$ , with the corresponding eigenvalue as  $\eta_i$ . Lemma 8 shows  $\eta_i \in \mathbb{R}$ , then we have

$$\begin{aligned} \lambda^2 - \eta_i\lambda - \beta\lambda + \beta &= 0, \\ \lambda_{i1} &= \frac{\eta_i + \beta + \sqrt{(\eta_i + \beta)^2 - 4\beta}}{2}, \\ \lambda_{i2} &= \frac{\eta_i + \beta - \sqrt{(\eta_i + \beta)^2 - 4\beta}}{2}. \end{aligned}$$

Here we analyze in two cases:

- 1) for any  $\eta_i$  that satisfies  $(\eta_i + \beta)^2 - 4\beta \leq 0$ , the modulus  $|\lambda_{i1}| = |\lambda_{i2}| = \sqrt{\beta}$ .
- 2) for any  $\eta_i$  that satisfies  $(\eta_i + \beta)^2 - 4\beta > 0$ ,  $\eta_i$  cannot be zero for any  $\beta \in (0, 1)$ . With  $\rho(\mathcal{L}) = \max_i |\eta_i| < 1$ , it holds that  $\eta_i < 1$ . In this case, with  $\beta \in (0, 1)$ , we have

$$\begin{aligned} (\eta_i + \beta)^2 - 4\beta &= (\eta_i - \beta)^2 - 4(1 - \eta_i)\beta < (\eta_i - \beta)^2, \\ \max\{|\lambda_{i1}|, |\lambda_{i2}|\} &= \frac{|\eta_i + \beta| + \sqrt{(\eta_i + \beta)^2 - 4\beta}}{2} \\ &< \frac{|\eta_i + \beta| + |\eta_i - \beta|}{2} = \max\{|\eta_i|, \beta\} \leq \max\{|\eta_i|, \sqrt{\beta}\}. \end{aligned}$$

Combining the two cases, if we choose a positive  $\beta$  that satisfies  $\beta < (\max_i \{|\eta_i|\})^2 = \rho^2(\mathcal{L})$ , let

$$q(0) = \rho(\mathcal{L}), q(\beta) = \rho(\tilde{\mathcal{L}}), \tau = \min\{1/5, \rho^2(\mathcal{L})\}, \quad (53)$$

then we have  $q(0) > q(\beta), \forall \beta \in (0, \tau)$ .  $\square$

### E. Proof of Theorem 3

Lemma 10 derives the properties of  $\hat{p}^{-1}\mathcal{P}_{\tilde{\mathcal{S}}_i}\tilde{\mathcal{G}}\mathcal{P}_{\Omega_{i+1}}\tilde{\mathcal{G}}^*(\mathcal{P}_{\tilde{\mathcal{U}}_i} - \mathcal{P}_{\tilde{\mathcal{U}}_i})$ , and the random operator can be close enough to its mean  $\mathcal{P}_{\tilde{\mathcal{S}}_i}\tilde{\mathcal{G}}\tilde{\mathcal{G}}^*(\mathcal{P}_{\tilde{\mathcal{U}}_i} - \mathcal{P}_{\tilde{\mathcal{U}}_i})$  with a significant large amount of observed entries. Lemma 11 illustrates the relation between  $\tilde{\mathbf{L}}_l$  and  $\tilde{\mathbf{L}}'_l$  and gives the bound on the incoherence of  $\tilde{\mathbf{L}}'_l$ , which is obtained after the trimming part (line 5 to 10). Lemmas 10 and 11 are built upon Lemmas 9 and 10 in [7] by extending from single-channel signals to multi-channel signals. Similar to the proof of Theorem 1, the proof of Theorem 3 is built upon that of Lemma 3 in [7], which is originally proposed as an initialization strategy. The major steps are devoted to bounding  $I_5, I_6$  and  $I_7$  in (58), and the corresponding results are presented in (59), (60), and (61). We include some details for the completeness of this proof.

**Lemma 10.** *Let  $\tilde{\mathbf{L}}'_l = \hat{\mathbf{U}}_l \hat{\Sigma}_l \hat{\mathbf{V}}_l^*$  and  $\tilde{\mathcal{G}}\mathbf{Y} = \tilde{\mathbf{U}} \tilde{\Sigma} \tilde{\mathbf{V}}^*$  be the SVD of  $\tilde{\mathbf{L}}'_l$  and  $\tilde{\mathcal{G}}\mathbf{Y}$ . Further let  $\hat{\mathcal{S}}_l$  be the tangent subspace of  $\tilde{\mathbf{L}}'_l$ . Assume there exists a constant  $\mu$  such that*

$$\left\|\mathcal{P}_{\hat{\mathcal{U}}_l} \tilde{\mathbf{H}}_{k,t}\right\|_F^2 \leq \frac{\mu c_s r}{n_c n}, \quad \left\|\mathcal{P}_{\hat{\mathcal{V}}_l} \tilde{\mathbf{H}}_{k,t}\right\|_F^2 \leq \frac{\mu c_s r}{n_c n},$$

and

$$\left\|\mathcal{P}_{\tilde{\mathcal{U}}} \tilde{\mathbf{H}}_{k,t}\right\|_F^2 \leq \frac{\mu c_s r}{n_c n}, \quad \left\|\mathcal{P}_{\tilde{\mathcal{V}}} \tilde{\mathbf{H}}_{k,t}\right\|_F^2 \leq \frac{\mu c_s r}{n_c n}.$$

for all  $1 \leq t \leq n, 1 \leq k \leq n_c$ . Let  $\Omega_{l+1} = \{(k_a, t_a) | a = 1, \dots, \hat{m}\}$  be a set of indices sampled with replacement. If  $\mathcal{P}_{\Omega_{l+1}}$  is independent of  $\tilde{\mathbf{U}}, \tilde{\mathbf{V}}, \tilde{\mathbf{U}}_l$  and  $\tilde{\mathbf{V}}_l$ , then

$$\left\| \mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}} (\mathcal{I} - \hat{p}^{-1} \mathcal{P}_{\Omega_{l+1}}) \tilde{\mathcal{G}}^* (\mathcal{P}_{\tilde{\mathcal{U}}} - \mathcal{P}_{\tilde{\mathcal{U}}_l}) \right\| \leq \sqrt{\frac{160 \mu_{c_s} r \log(n)}{\hat{m}}}$$

with probability at least  $1 - n_c n^{-2}$ , if  $\hat{m} \geq \frac{125}{18} \mu_{c_s} r \log(n)$ .

**Lemma 11.** Let  $\tilde{\mathbf{L}}_l = \tilde{\mathbf{U}}_l \tilde{\Sigma}_l \tilde{\mathbf{V}}_l^*$  and  $\tilde{\mathcal{G}}\mathbf{Y} = \tilde{\mathbf{U}} \tilde{\Sigma} \tilde{\mathbf{V}}^*$  be the SVD of  $\tilde{\mathbf{L}}_l$  and  $\tilde{\mathcal{G}}\mathbf{Y}$ . Assume

$$\max_{k_1} \|\tilde{\mathbf{U}}_{k_1^*}\|^2 \leq \frac{\mu_{c_s} r}{n_c n} \text{ and } \max_{k_2} \|\tilde{\mathbf{V}}_{k_2^*}\|^2 \leq \frac{\mu_{c_s} r}{n_c n}.$$

Suppose  $\tilde{\mathbf{L}}_l$  and  $\tilde{\mathcal{G}}\mathbf{Y}$  are both rank- $r$  matrices satisfying

$$\|\tilde{\mathbf{L}}_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F \leq \frac{\sigma_{\min}(\tilde{\mathcal{G}}\mathbf{Y})}{10\sqrt{2}}.$$

Then the matrix  $\tilde{\mathbf{L}}'_l = \hat{\mathbf{U}}_l \hat{\Sigma}_l \hat{\mathbf{V}}_l^*$ , denoting the SVD of  $\tilde{\mathbf{L}}'_l$ , that is obtained after trimming in RAM-FIHT satisfies

$$\|\tilde{\mathbf{L}}'_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F \leq 8\kappa \|\tilde{\mathbf{L}}_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F, \quad (54)$$

$$\max_{k_1, k_2} \left\{ \|\hat{\mathbf{U}}_{k_1^*}\|^2, \|\hat{\mathbf{V}}_{k_2^*}\|^2 \right\} \leq \frac{100 \mu_{c_s} r}{81 n_c n}, \quad (55)$$

where  $\kappa$  denotes the condition number of  $\tilde{\mathcal{G}}\mathbf{Y}$ .

*Proof of Theorem 3.* First, we show the following inequality holds with high probability by mathematical induction.

$$\|\tilde{\mathbf{L}}_k - \tilde{\mathcal{G}}\mathbf{Y}\|_F \leq \frac{\varepsilon_0 \sigma_{\min}(\tilde{\mathcal{G}}\mathbf{Y})}{128\kappa^2}. \quad (56)$$

**Inductive Step:** Suppose (56) holds when  $k = l$  and  $l \geq 0$ . Then (55) in Lemma 11 holds. Further, we can conclude

$$\left\| \mathcal{P}_{\tilde{\mathcal{U}}_l} \tilde{\mathbf{H}}_{k,t} \right\|_F^2 \leq \frac{100 \mu_{c_s} r}{81 n_c n} \text{ and } \left\| \mathcal{P}_{\tilde{\mathcal{V}}_l} \tilde{\mathbf{H}}_{k,t} \right\|_F^2 \leq \frac{100 \mu_{c_s} r}{81 n_c n}. \quad (57)$$

Recall that  $\mathbf{Y} = \tilde{\mathcal{D}}\mathbf{X}$  and  $\tilde{\mathcal{G}}\mathbf{Y} = \tilde{\mathcal{H}}\mathbf{X}$ . Define  $\hat{\mathbf{Y}}_l = \tilde{\mathcal{D}}\hat{\mathbf{X}}_l$ . Since the measurements are noiseless, then  $\mathbf{M} = \mathbf{X}$  and

$$\tilde{\mathcal{H}}(\hat{\mathbf{X}}_l + \hat{p}^{-1} \tilde{\mathcal{P}}_{\Omega_{l+1}}(\mathbf{X} - \hat{\mathbf{X}}_l)) = \tilde{\mathcal{G}}(\hat{\mathbf{Y}}_l + \hat{p}^{-1} \tilde{\mathcal{P}}_{\Omega_{l+1}}(\mathbf{Y} - \hat{\mathbf{Y}}_l)).$$

Then,

$$\begin{aligned} & \left\| \tilde{\mathbf{L}}_{l+1} - \tilde{\mathcal{G}}\mathbf{Y} \right\|_F \leq 2 \left\| \mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}}(\hat{\mathbf{Y}}_l + \hat{p}^{-1} \mathcal{P}_{\Omega_{l+1}}(\mathbf{Y} - \hat{\mathbf{Y}}_l)) - \tilde{\mathcal{G}}\mathbf{Y} \right\|_F \\ & \leq 2 \left\| \mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}}\mathbf{Y} - \tilde{\mathcal{G}}\mathbf{Y} \right\|_F + 2 \left\| (\mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}} - \hat{p}^{-1} \mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}} \mathcal{P}_{\Omega_{l+1}})(\mathbf{Y} - \hat{\mathbf{Y}}_l) \right\|_F \\ & = 2 \left\| (\mathcal{I} - \mathcal{P}_{\tilde{\mathcal{S}}_l}) \tilde{\mathcal{G}}\mathbf{Y} \right\|_F \\ & \quad + 2 \left\| (\mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}} \tilde{\mathcal{G}}^* - \hat{p}^{-1} \mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}} \mathcal{P}_{\Omega_{l+1}} \tilde{\mathcal{G}}^*)(\tilde{\mathbf{L}}'_l - \tilde{\mathcal{G}}\mathbf{Y}) \right\|_F \\ & \leq 2 \left\| (\mathcal{I} - \mathcal{P}_{\tilde{\mathcal{S}}_l})(\tilde{\mathbf{L}}'_l - \tilde{\mathcal{G}}\mathbf{Y}) \right\|_F \\ & \quad + 2 \left\| (\mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}} \tilde{\mathcal{G}}^* \mathcal{P}_{\tilde{\mathcal{S}}_l} - \hat{p}^{-1} \mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}} \mathcal{P}_{\Omega_{l+1}} \tilde{\mathcal{G}}^* \mathcal{P}_{\tilde{\mathcal{S}}_l})(\tilde{\mathbf{L}}'_l - \tilde{\mathcal{G}}\mathbf{Y}) \right\|_F \\ & \quad + 2 \left\| \mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}}(\mathcal{I} - \hat{p}^{-1} \mathcal{P}_{\Omega_{l+1}}) \tilde{\mathcal{G}}^*(\mathcal{I} - \mathcal{P}_{\tilde{\mathcal{S}}_l})(\tilde{\mathbf{L}}'_l - \tilde{\mathcal{G}}\mathbf{Y}) \right\|_F \\ & := I_5 + I_6 + I_7. \end{aligned} \quad (58)$$

where the first inequality comes from (39).

With (54) and (56),  $I_5$  can be bounded as

$$I_5 \leq \frac{2 \|\tilde{\mathbf{L}}'_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F^2}{\sigma_{\min}(\tilde{\mathcal{G}}\mathbf{Y})} \leq \varepsilon_0 \|\tilde{\mathbf{L}}_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F, \quad (59)$$

As for the item  $I_6$ , Lemma 4 along with (54) suggests

$$\begin{aligned} I_6 & \leq 2 \sqrt{\frac{3200 \mu_{c_s} r \log(n)}{81 \hat{m}}} \|\tilde{\mathbf{L}}'_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F \\ & \leq 16\kappa \sqrt{\frac{3200 \mu_{c_s} r \log(n)}{81 \hat{m}}} \|\tilde{\mathbf{L}}_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F \end{aligned} \quad (60)$$

with probability at least  $1 - n_c n^{-2}$ . To bound  $I_7$ ,

$$\begin{aligned} & (\mathcal{I} - \mathcal{P}_{\tilde{\mathcal{S}}_l})(\tilde{\mathbf{L}}'_l - \tilde{\mathcal{G}}\mathbf{Y}) = (\mathbf{I} - \hat{\mathbf{U}}_l \hat{\mathbf{U}}_l^*)(-\tilde{\mathcal{G}}\mathbf{Y})(\mathbf{I} - \hat{\mathbf{V}}_l \hat{\mathbf{V}}_l^*) \\ & = (\tilde{\mathbf{U}} \tilde{\mathbf{U}}^* - \hat{\mathbf{U}}_l \hat{\mathbf{U}}_l^*)(\tilde{\mathbf{L}}'_l - \tilde{\mathcal{G}}\mathbf{Y})(\mathbf{I} - \hat{\mathbf{V}}_l \hat{\mathbf{V}}_l^*) \\ & = (\mathcal{P}_{\tilde{\mathcal{U}}} - \mathcal{P}_{\tilde{\mathcal{U}}_l})(\mathcal{I} - \mathcal{P}_{\tilde{\mathcal{V}}})(\tilde{\mathbf{L}}'_l - \tilde{\mathcal{G}}\mathbf{Y}). \end{aligned}$$

Hence, by Lemma 10, with probability at least  $1 - n_c n^{-2}$ ,

$$\begin{aligned} I_7 & = 2 \left\| \mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}}(\mathcal{I} - \hat{p}^{-1} \mathcal{P}_{\Omega_{l+1}}) \tilde{\mathcal{G}}^* (\mathcal{P}_{\tilde{\mathcal{U}}} - \mathcal{P}_{\tilde{\mathcal{U}}_l})(\mathcal{I} - \mathcal{P}_{\tilde{\mathcal{V}}})(\tilde{\mathbf{L}}'_l - \tilde{\mathcal{G}}\mathbf{Y}) \right\|_F \\ & \leq 2 \left\| \mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}}(\mathcal{I} - \hat{p}^{-1} \mathcal{P}_{\Omega_{l+1}}) \tilde{\mathcal{G}}^* (\mathcal{P}_{\tilde{\mathcal{U}}} - \mathcal{P}_{\tilde{\mathcal{U}}_l}) \right\| \|\tilde{\mathbf{L}}'_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F \\ & \leq 16\kappa \sqrt{\frac{16000 \mu_{c_s} r \log(n)}{81 \hat{m}}} \|\tilde{\mathbf{L}}'_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F. \end{aligned} \quad (61)$$

Therefore, if  $\hat{m} \geq C_4 \varepsilon_0^{-2} \mu_{c_s} \kappa^2 r \log(n)$  for some constant  $C_4$ ,

$$I_6 + I_7 \leq 326\kappa \sqrt{\frac{\mu_{c_s} r \log(n)}{\hat{m}}} \|\tilde{\mathbf{L}}_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F \leq \varepsilon_0 \|\tilde{\mathbf{L}}_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F.$$

Putting pieces together gives

$$\|\tilde{\mathbf{L}}_{l+1} - \tilde{\mathcal{G}}\mathbf{Y}\|_F \leq 2\varepsilon_0 \|\tilde{\mathbf{L}}_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F \quad (62)$$

with probability at least  $1 - 2n_c n^{-2}$ . Hence, (56) also holds when  $k = l + 1$ .

**Base Case:** Since  $\tilde{\mathbf{L}}_0 = \mathcal{Q}_r(\hat{p}^{-1} \tilde{\mathcal{H}} \mathcal{P}_{\Omega_0}(\mathbf{X}))$ , we can follow the same idea in the proof of base case in Theorem 1. Thus, when  $k = 0$ , (56) is valid with probability at least  $1 - n_c n^{-2}$  provided  $\hat{m} \geq C_5 \varepsilon_0^{-2} \mu_{c_s} \kappa^6 r^2 \log(n)$  for some constant  $C_5$ .

Let  $C_2 = \max\{C_4, C_5\}$ . If  $\hat{m} \geq C_2 \varepsilon_0^{-2} \mu_{c_s} \kappa^6 r^2 \log(n)$ , then for each  $l \geq 0$ , we have

$$\|\tilde{\mathbf{L}}_{l+1} - \tilde{\mathcal{G}}\mathbf{Y}\|_F \leq 2\varepsilon_0 \|\tilde{\mathbf{L}}_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F.$$

with probability at least  $1 - 2n_c n^{-2}$ . Directly the following inequality is obtained with probability  $1 - (2L + 1)n_c n^{-2}$ ,

$$\|\tilde{\mathbf{L}}_L - \tilde{\mathcal{G}}\mathbf{Y}\|_F \leq \nu^L \|\tilde{\mathbf{L}}_0 - \tilde{\mathcal{G}}\mathbf{Y}\|_F \leq \nu^L \frac{\varepsilon_0 \sigma_{\min}(\tilde{\mathcal{G}}\mathbf{Y})}{128\kappa^2}. \quad (63)$$

If we take  $L = \left\lceil \varepsilon_0^{-1} \log \left( \frac{\sigma_{\max}(\mathcal{H}\mathbf{X})}{128\kappa^3 \varepsilon} \right) \right\rceil$  with an arbitrarily small positive constant  $\varepsilon$ , since  $\sigma_{\max}(\tilde{\mathcal{G}}\mathbf{Y}) = \sqrt{n_c} \sigma_{\max}(\mathcal{H}\mathbf{X})$ ,

$$\|\tilde{\mathbf{L}}_L - \tilde{\mathcal{G}}\mathbf{Y}\|_F \leq n_c^{1/2} \varepsilon, \quad (64)$$

which completes the proof of Theorem 3.  $\square$

#### F. Proof of Theorem 4

The proof of Theorem 4 is similar to Theorem 3, except some modification to handle the noise matrix  $\mathbf{N}$ . We first present the following lemma that will be useful in the proof.

**Lemma 12.** Suppose  $m \geq 16 \log(n)$ , then

$$\left\| \hat{p}^{-1} \tilde{\mathcal{H}} \mathcal{P}_{\Omega}(\mathbf{N}) - \tilde{\mathcal{H}}\mathbf{N} \right\| \leq \sqrt{\frac{16 \log(n)}{m}} n_c n \|\tilde{\mathcal{H}}\mathbf{N}\|_{\infty}$$

with probability at least  $1 - n_c n^{-2}$ .

*Proof of Theorem 4.* For the noisy case where  $\mathbf{M} = \mathbf{X} + \mathbf{N}$ , we have assumed  $\|\mathbf{N}\|_\infty \leq \frac{\varepsilon_0 \|\mathcal{H}\mathbf{X}\|}{2048\kappa^3 r^{1/2} n_c^{1/2} n}$ . Recall (29), define  $\mathbf{S} = \tilde{\mathbf{D}}\mathbf{N}$ , then

$$\|\tilde{\mathcal{G}}\mathbf{S}\|_\infty = \|\tilde{\mathcal{H}}\mathbf{N}\|_\infty = \|\mathbf{N}\|_\infty \leq \frac{\varepsilon_0 \sigma_{\min}(\tilde{\mathcal{G}}\mathbf{Y})}{2048\kappa^2 r^{1/2} n_c n}.$$

Similar to the derivation of (39), we have

$$\begin{aligned} & \|\tilde{\mathbf{L}}_{l+1} - \tilde{\mathcal{G}}\mathbf{Y}\| \leq 2 \left\| \mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}}(\hat{\mathbf{Y}}_l + \hat{p}^{-1} \mathcal{P}_{\Omega_{l+1}}(\mathbf{Y} + \mathbf{S} - \hat{\mathbf{Y}}_l)) - \tilde{\mathcal{G}}\mathbf{Y} \right\| \\ & \leq 2 \left\| \mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}}(\hat{\mathbf{Y}}_l + \hat{p}^{-1} \mathcal{P}_{\Omega_{l+1}}(\mathbf{Y} - \hat{\mathbf{Y}}_l)) - \tilde{\mathcal{G}}\mathbf{Y} \right\| \\ & \quad + 2 \left\| \hat{p}^{-1} \mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}} \mathcal{P}_{\Omega_{l+1}}(\mathbf{S}) \right\|. \\ & \|\tilde{\mathbf{L}}_{l+1} - \tilde{\mathcal{G}}\mathbf{Y}\|_F \leq \sqrt{2r} \|\tilde{\mathbf{L}}_{l+1} - \tilde{\mathcal{G}}\mathbf{Y}\| \\ & \leq 2\sqrt{2r} \left\| \mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}}(\hat{\mathbf{Y}}_l + \hat{p}^{-1} \mathcal{P}_{\Omega_{l+1}}(\mathbf{Y} + \mathbf{S} - \hat{\mathbf{Y}}_l)) - \tilde{\mathcal{G}}\mathbf{Y} \right\| \\ & \leq 2\sqrt{2r} \left\| \mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}}(\hat{\mathbf{Y}}_l + \hat{p}^{-1} \mathcal{P}_{\Omega_{l+1}}(\mathbf{Y} - \hat{\mathbf{Y}}_l)) - \tilde{\mathcal{G}}\mathbf{Y} \right\| \\ & \quad + 2\sqrt{2r} \left\| \hat{p}^{-1} \mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}} \mathcal{P}_{\Omega_{l+1}}(\mathbf{S}) \right\| \\ & \leq 2\sqrt{2r} \left\| (\mathcal{I} - \mathcal{P}_{\tilde{\mathcal{S}}_l})(\tilde{\mathbf{L}}'_l - \tilde{\mathcal{G}}\mathbf{Y}) \right\|_F \\ & \quad + 2\sqrt{2r} \left\| (\mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}} \tilde{\mathcal{G}}^* \mathcal{P}_{\tilde{\mathcal{S}}_l} - \hat{p}^{-1} \mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}} \mathcal{P}_{\Omega_{l+1}} \tilde{\mathcal{G}}^* \mathcal{P}_{\tilde{\mathcal{S}}_l})(\tilde{\mathbf{L}}'_l - \tilde{\mathcal{G}}\mathbf{Y}) \right\|_F \\ & \quad + 2\sqrt{2r} \left\| \mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}}(\mathcal{I} - \hat{p}^{-1} \mathcal{P}_{\Omega_{l+1}}) \tilde{\mathcal{G}}^*(\mathcal{I} - \mathcal{P}_{\tilde{\mathcal{S}}_l})(\tilde{\mathbf{L}}'_l - \tilde{\mathcal{G}}\mathbf{Y}) \right\|_F \\ & \quad + 2\sqrt{2r} \left\| \hat{p}^{-1} \mathcal{P}_{\tilde{\mathcal{S}}_l} \tilde{\mathcal{G}} \mathcal{P}_{\Omega_{l+1}}(\mathbf{S}) \right\| \\ & := \sqrt{2r}(I_5 + I_6 + I_7) + I_9, \end{aligned} \quad (65)$$

where  $I_5 + I_6 + I_7$  has been defined in (58).

Similar to the proof of Theorem 3, we show that the following inequality holds with high probability by induction.

$$\|\tilde{\mathbf{L}}_k - \tilde{\mathcal{G}}\mathbf{Y}\|_F \leq \frac{\varepsilon_0 \sigma_{\min}(\tilde{\mathcal{G}}\mathbf{Y})}{128\sqrt{2}\kappa^2 r^{1/2}}. \quad (66)$$

**Inductive Step:** Suppose (66) holds when  $k = l$  and  $l \geq 0$ . By Lemma 5 and (56), we have

$$\begin{aligned} \sqrt{2r}I_5 & \leq \frac{2\sqrt{2r} \|\tilde{\mathbf{L}}'_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F^2}{\sigma_{\min}(\tilde{\mathcal{G}}\mathbf{Y})} \leq \frac{128\sqrt{2}\kappa^2 \sqrt{r} \|\tilde{\mathbf{L}}_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F^2}{\sigma_{\min}(\tilde{\mathcal{G}}\mathbf{Y})} \\ & \leq \varepsilon_0 \|\tilde{\mathbf{L}}_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F. \end{aligned}$$

Similar to  $I_6$  and  $I_7$ ,

$$\sqrt{2r}(I_6 + I_7) \leq 326\kappa \sqrt{\frac{2\mu c_s r^2 \log(n)}{\hat{m}}} \|\tilde{\mathbf{L}}_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F.$$

Hence, if  $\hat{m} \geq C_6 \varepsilon_0^{-2} \mu c_s \kappa^2 r^2 \log(n)$  for some constant  $C_6$ ,

$$\sqrt{2r}(I_5 + I_6 + I_7) \leq 2\varepsilon_0 \|\tilde{\mathbf{L}}_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F$$

with probability at least  $1 - 2n_c n^{-2}$ . On the other hand, if  $m \geq 256r \log(n)$ , then with probability at least  $1 - n_c n^{-2}$

$$\begin{aligned} I_9 & \leq 2\sqrt{2r} \left\| \hat{p}^{-1} \tilde{\mathcal{G}} \mathcal{P}_{\Omega_{l+1}}(\mathbf{S}) \right\| \\ & \leq 2\sqrt{2r} \left\| \hat{p}^{-1} \tilde{\mathcal{G}} \mathcal{P}_{\Omega_{l+1}}(\mathbf{S}) - \tilde{\mathcal{G}}\mathbf{S} \right\| + 2\sqrt{2r} \|\tilde{\mathcal{G}}\mathbf{S}\| \\ & \leq 8\sqrt{2} \sqrt{\frac{r \log(n)}{m}} n_c n \|\tilde{\mathcal{G}}\mathbf{S}\|_\infty + 2\sqrt{2r} n_c n \|\tilde{\mathcal{G}}\mathbf{S}\|_\infty \\ & \leq \frac{1}{16\sqrt{2}} \sqrt{\frac{r \log(n)}{m}} \frac{\varepsilon_0 \sigma_{\min}(\tilde{\mathcal{G}}\mathbf{Y})}{\kappa^2 r^{1/2}} + \frac{\varepsilon_0 \sigma_{\min}(\tilde{\mathcal{G}}\mathbf{Y})}{512\sqrt{2}\kappa^2 r^{1/2}} \\ & \leq \frac{\varepsilon_0 \sigma_{\min}(\tilde{\mathcal{G}}\mathbf{Y})}{256\sqrt{2}\kappa^2 r^{1/2}}, \end{aligned} \quad (67)$$

where the second last inequality comes from Lemma 12. Following  $\nu = 2\varepsilon_0 \leq 1/2$  and (66), with probability at least  $1 - 3n_c n^{-2}$ , we can bound  $\|\tilde{\mathbf{L}}_{l+1} - \tilde{\mathcal{G}}\mathbf{Y}\|_F$  by

$$\frac{1}{2} \|\tilde{\mathbf{L}}_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F + \frac{\varepsilon_0 \sigma_{\min}(\tilde{\mathcal{G}}\mathbf{Y})}{256\sqrt{2}\kappa^2 r^{1/2}} \leq \frac{\varepsilon_0 \sigma_{\min}(\tilde{\mathcal{G}}\mathbf{Y})}{128\sqrt{2}\kappa^2 r^{1/2}}.$$

Hence, (66) also holds when  $k = l + 1$ .

**Base Case:** Since  $\tilde{\mathbf{L}}_0 = \mathcal{Q}_r(\hat{p}^{-1} \mathcal{H} \mathcal{P}_{\Omega_0}(\mathbf{X} + \mathbf{N}))$ , then with probability at least  $1 - n_c n^{-2}$ ,

$$\begin{aligned} & \|\tilde{\mathbf{L}}_0 - \tilde{\mathcal{G}}\mathbf{Y}\|_F \leq \sqrt{2r} \|\tilde{\mathbf{L}}_0 - \tilde{\mathcal{G}}\mathbf{Y}\| \\ & \leq \sqrt{2r} \left\| \hat{p}^{-1} \tilde{\mathcal{G}} \mathcal{P}_{\Omega_0}(\mathbf{Y} + \mathbf{S}) - \tilde{\mathbf{L}}_0 \right\| + \sqrt{2r} \left\| \hat{p}^{-1} \tilde{\mathcal{G}} \mathcal{P}_{\Omega_0}(\mathbf{Y} + \mathbf{S}) - \tilde{\mathcal{G}}\mathbf{Y} \right\| \\ & \leq 2\sqrt{2r} \left\| \hat{p}^{-1} \tilde{\mathcal{G}} \mathcal{P}_{\Omega_0}(\mathbf{Y}) - \tilde{\mathcal{G}}\mathbf{Y} \right\| + 2\sqrt{2r} \left\| \hat{p}^{-1} \tilde{\mathcal{G}} \mathcal{P}_{\Omega_0}(\mathbf{S}) \right\| \\ & \leq \sqrt{\frac{512\mu c_s r^2 \log(n)}{m}} \|\tilde{\mathcal{G}}\mathbf{Y}\| + \frac{\varepsilon_0 \sigma_{\min}(\tilde{\mathcal{G}}\mathbf{Y})}{256\sqrt{2}\kappa^2 \sqrt{r}}. \end{aligned}$$

where the last inequality comes from (67) and Lemma 6. To guarantee that (66) holds with  $k = 0$ , we need

$$\sqrt{\frac{512\mu c_s r^2 \log(n)}{m}} \|\tilde{\mathcal{G}}\mathbf{Y}\| \leq \frac{\varepsilon_0 \sigma_{\min}(\tilde{\mathcal{G}}\mathbf{Y})}{256\sqrt{2}\kappa^2 \sqrt{r}}. \quad (68)$$

That is  $\hat{m} \geq C_7 \varepsilon_0^{-2} \mu c_s \kappa^6 r^3 \log(n)$  for some constant  $C_7$ .

Let  $C_3 = \max\{C_6, C_7\}$ , if  $\hat{m} \geq C_3 \varepsilon_0^{-2} \mu c_s \kappa^6 r^3 \log(n)$ , for each  $l \geq 0$ , with probability  $1 - 2n_c n^{-2}$ , we have

$$\|\tilde{\mathbf{L}}_{l+1} - \tilde{\mathcal{G}}\mathbf{Y}\|_F \leq 2\varepsilon_0 \|\tilde{\mathbf{L}}_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F + \Delta. \quad (69)$$

where  $\Delta = 32\sqrt{2}n_c n \|\tilde{\mathcal{G}}\mathbf{S}\|_\infty + 2\sqrt{2}r^{1/2} \|\tilde{\mathcal{G}}\mathbf{S}\|$ . Then

$$\|\tilde{\mathbf{L}}_{l+1} - \tilde{\mathcal{G}}\mathbf{Y}\|_F - \frac{\Delta}{1-\nu} \leq \nu \left( \|\tilde{\mathbf{L}}_l - \tilde{\mathcal{G}}\mathbf{Y}\|_F - \frac{\Delta}{1-\nu} \right).$$

Therefore, with probability  $1 - (3L + 1)n_c n^{-2}$ ,

$$\|\tilde{\mathbf{L}}_L - \tilde{\mathcal{G}}\mathbf{Y}\|_F \leq \nu^L \|\tilde{\mathbf{L}}_0 - \tilde{\mathcal{G}}\mathbf{Y}\|_F + \frac{\Delta}{1-\nu}. \quad (70)$$

Similar to (64), take  $L = \left\lceil \varepsilon_0^{-1} \log\left(\frac{\sigma_{\max}(\mathcal{H}\mathbf{X})}{128\kappa^3 \varepsilon}\right) \right\rceil$  with an arbitrarily small positive constant  $\varepsilon$ , since  $\nu \leq 1/2$ ,

$$\begin{aligned} \|\tilde{\mathbf{L}}_L - \tilde{\mathcal{G}}\mathbf{Y}\|_F & \leq n_c^{1/2} \varepsilon + 64\sqrt{2}n_c n \|\tilde{\mathcal{G}}\mathbf{S}\|_\infty + 4\sqrt{2}r^{1/2} \|\tilde{\mathcal{G}}\mathbf{S}\| \\ & \leq n_c^{1/2} \varepsilon + 128n_c n \|\tilde{\mathcal{G}}\mathbf{S}\|_\infty + 8r^{1/2} \|\tilde{\mathcal{G}}\mathbf{S}\|. \end{aligned}$$

which completes the proof of Theorem 4.  $\square$

### G. Proof of Theorem 5

We first introduce some useful lemmas.

**Lemma 13** ([23], Corollary 7.7.4(a)). *If  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^n$  are positive-definite, then  $\mathbf{A} \succeq \mathbf{B}$  if and only if  $\mathbf{B}^{-1} \succeq \mathbf{A}^{-1}$*

**Lemma 14** ([32], Theorem 7). *Let  $\lambda_1 \geq \dots \geq \lambda_n$  be eigenvalues of  $\mathbf{A}$ , denoted by  $\lambda_i(\mathbf{A}) = \lambda_i$ . Let  $\mathbf{A}$  and  $\mathbf{B}$  be Hermitian positive semi-definite  $n \times n$  matrices. If  $1 \leq k \leq i \leq n$  and  $1 \leq l \leq n - i + 1$ , then*

$$\lambda_{i+l-1}(\mathbf{A}) \lambda_{n-l+1}(\mathbf{B}) \leq \lambda_i(\mathbf{AB}) \leq \lambda_{i-k+1}(\mathbf{A}) \lambda_k(\mathbf{B}).$$

In particular,

$$\lambda_n(\mathbf{A}) \lambda_n(\mathbf{B}) \leq \lambda_n(\mathbf{AB}), \quad \lambda_1(\mathbf{AB}) \leq \lambda_1(\mathbf{A}) \lambda_1(\mathbf{B}).$$

*Proof of Theorem 5.* Consider  $n_c = 1$ , the definition of  $\mathcal{H}$  can be extended to a row vector, which corresponds with one

channel data. Let  $\mathcal{H}X_{k*} = U_k \Sigma_k V_k$  and  $\mathcal{H}X = U \Sigma V$  be the SVD of  $\mathcal{H}X_{k*}$  and  $\mathcal{H}X$ . Then,  $\mu_0$  is defined as

$$\max_{k_1} \|e_{k_1}^* U\|^2 \leq \frac{\mu_0 r}{n_1}, \quad \max_{k_2} \|e_{k_2}^* V\|^2 \leq \frac{\mu_0 r}{n_2}.$$

Notice that all  $\mathcal{H}X_{k*}$  share the same column space and row space, they have the same incoherence  $\mu_0$ . It is trivial  $\mu = \mu_0$  if we consider the incoherence of row spaces. Hence, we only focus on the incoherence of column space.

By (6),  $\mathcal{H}X = P_L \Gamma P_R^T$ . Define a series of diagonal matrices  $D_k$  as  $D_k = \text{diag}(d_k)$  with  $1 \leq k \leq n_c$ , and  $d_k = [d_{k,1}, \dots, d_{k,r}]$ , where  $d_{k,i} = r_i^* s_1 C_{k*} l_i$ . We need one mild assumption that  $d_{k,i} \neq 0$ . It guarantees that each  $D_k$  is full rank. Thus,  $\mathcal{H}X_{k*} = E_L D_k P_R$ , where  $E_L = P_L$  with  $n_c = 1$ . There exists a row-switching matrix  $Q_1$  satisfying

$$Q_1(\mathcal{H}X) = \begin{bmatrix} \mathcal{H}X_{1*} \\ \mathcal{H}X_{2*} \\ \vdots \\ \mathcal{H}X_{n_c*} \end{bmatrix} = \begin{bmatrix} E_L D_1 \\ E_L D_2 \\ \vdots \\ E_L D_{n_c} \end{bmatrix} P_R := \tilde{E}_L P_R.$$

Define a mapping  $f: \{1, 2, \dots, n_c n_1\} \mapsto \{1, 2, \dots, n_c n_1\}$ ,  $f(z) = w$  with  $e_z = Q_1 e_w$ , then  $f$  is a bijective mapping. Hence, we have

$$\begin{aligned} \max_{i_1} \|e_{i_1}^* U\|^2 &= \max_{k_1} (Q_1 e_{k_1})^* \tilde{E}_L (\tilde{E}_L^* \tilde{E}_L)^{-1} \tilde{E}_L^* (Q_1 e_{k_1}) \\ &= \max_{k_1} e_{k_1}^* \tilde{E}_L \left( \sum_{k=1}^{n_c} D_k^* E_L^* E_L D_k \right)^{-1} \tilde{E}_L^* e_{k_1}. \end{aligned} \quad (71)$$

Consider  $1 \leq k_1 \leq n_1$ , we know that  $e_{k_1}^* \tilde{E}_L = \hat{e}_{k_1}^* E_L D_1$ , where  $e_{k_1} \in \mathbb{C}^{n_c n_1}$  and  $\hat{e}_{k_1} \in \mathbb{C}^{n_1}$  are both coordinate vectors. Additionally, it is easy to show that symmetric matrices  $\{D_k^* E_L^* E_L D_k\}_{k=1}^{n_c}$  are positive definite since  $\{E_L D_k\}_{k=1}^{n_c}$  are full rank. Also, following Lemma 13, we have

$$\begin{aligned} \sum_{k=1}^{n_c} D_k^* E_L^* E_L D_k &\succ D_1^* E_L^* E_L D_1 \succ 0, \\ \left( \sum_{k=1}^{n_c} D_k^* E_L^* E_L D_k \right)^{-1} &\prec (D_1^* E_L^* E_L D_1)^{-1}. \end{aligned} \quad (72)$$

Then,

$$\begin{aligned} \|e_{k_1}^* U\|^2 &= \hat{e}_{k_1}^* E_L D_1 \left( \sum_{k=1}^{n_c} D_k^* E_L^* E_L D_k \right)^{-1} D_1^* E_L^* \hat{e}_{k_1} \\ &< \hat{e}_{k_1}^* E_L D_1 (D_1^* E_L^* E_L D_1)^{-1} D_1^* E_L^* \hat{e}_{k_1} \\ &= \hat{e}_{k_1}^* E_L (E_L^* E_L)^{-1} E_L^* \hat{e}_{k_1} \leq \frac{u_0 r}{n_1} = \frac{(n_c u_0) r}{n_c n_1}. \end{aligned}$$

Similarly, we can prove  $\|e_{k_1}^* U\|^2 < \frac{(n_c u_0) r}{n_c n_1}$  for all  $i$  satisfying  $1 \leq i \leq n_c n_1$ , which leads to (25).

Moreover, we can provide a tighter bound on  $\mu$  with a stronger assumption. Suppose there exists a  $\hat{d} \in \mathbb{C}$  and a real number  $\delta \in (0, 1)$  satisfying  $(1 - \delta)|\hat{d}| \leq |d_{k,i}| \leq (1 + \delta)|\hat{d}|$ . By Lemma 14, define  $\kappa_L = \frac{\sigma_{\max}(E_L)}{\sigma_{\min}(E_L)}$ , then:

$$\begin{aligned} \lambda_{\max}(D_1^* E_L^* E_L D_1) &= \lambda_{\max}(E_L^* E_L D_1 D_1^*) \\ &\leq \lambda_{\max}(E_L^* E_L) \lambda_{\max}(D_1 D_1^*) \\ &\leq \frac{\kappa_L^2 (1 + \delta)^2}{(1 - \delta)^2} \lambda_{\min}(E_L^* E_L) \lambda_{\min}(D_k D_k^*) \\ &\leq \frac{\kappa_L^2 (1 + \delta)^2}{(1 - \delta)^2} \lambda_{\min}(D_k^* E_L^* E_L D_k). \end{aligned}$$

since even the minimum eigenvalue of  $D_k^* E_L^* E_L D_k$  is larger than the maximum one of  $\frac{(1 - \delta)^2}{\kappa_L^2 (1 + \delta)^2} D_1^* E_L^* E_L D_1$ , we have

$$\sum_{k=1}^{n_c} D_k^* E_L^* E_L D_k \succeq [1 + (n_c - 1) \frac{(1 - \delta)^2}{\kappa_L^2 (1 + \delta)^2}] D_1^* E_L^* E_L D_1.$$

Similarly, we can establish the following relation between  $\mu$  and  $\mu_0$ ,

$$\mu \leq \frac{n_c \mu_0}{1 + (n_c - 1) \frac{(1 - \delta)^2}{\kappa_L^2 (1 + \delta)^2}}. \quad \square$$



**Shuai Zhang** received the B.E. degree from University of Science and Technology of China, Hefei, China, in 2016.

He is pursuing the Ph.D. degree in electrical engineering at Rensselaer Polytechnic Institute, Troy, NY. His research interests include signal processing and high dimensional data analysis.



**Yingshuai Hao** (S'14) received the B.E. degree from Shandong University, Jinan, China, in 2011 and the M.S. degree in electrical engineering from Shanghai Jiao Tong University, Shanghai, China, in 2014.

He is pursuing the Ph.D. degree in electrical engineering at Rensselaer Polytechnic Institute, Troy, NY. His research interests include cyber security of power systems, and PMU data quality improvement.



**Meng Wang** (M'12) received the Ph.D. degree from Cornell University, Ithaca, NY, USA, in 2012.

She is an Assistant Professor in the department of Electrical, Computer, and Systems Engineering at Rensselaer Polytechnic Institute. Her research interests include high dimensional data analysis and their applications in power systems monitoring and network inference.



**Joe H. Chow** (F'92) received the M.S. and Ph.D. degrees from the University of Illinois, Urbana-Champaign, Urbana, IL, USA.

After working in the General Electric power system business in Schenectady, NY, USA, he joined Rensselaer Polytechnic Institute, Troy, NY, USA, in 1987, where he is a Professor of Electrical, Computer, and Systems Engineering. His research interests include multivariable control, power system dynamics and control, FACTS controllers, and synchronized phasor data.