

Fast and Efficient Compressive Sensing using Structurally Random Matrices

Thong T. Do, *Student Member, IEEE*, Lu Gan, *Member, IEEE*, Nam Nguyen and Trac D.
Tran, *Member, IEEE*

This work has been supported in part by the National Science Foundation under Grant CCF-0728893.
Thong T. Do, Nam Nguyen and Trac D. Tran are with the Johns Hopkins University, Baltimore, MD, 21218 USA.
Lu Gan is with the Brunel University, London, UK.

Abstract

This paper introduces a fast and efficient framework for practical compressive sensing. Our framework is mainly based on a novel design called *Structurally Random Matrix* (SRM). It is highly promising for large-scale, real-time compressive sensing applications because it can be realized as a product of simple and fast operators and thus, there is no need for storing the sensing matrix explicitly. The introduced framework is flexible and provides relevant features such as universality, block-based processing and hardware friendliness to analog and optical domain implementation. Despite all of these practical advantages, the framework can be shown to approach optimal performance, i.e. the number of measurements for exact signal reconstruction is at the minimum bound. Simulation results with several interesting SRM under various practical settings are also presented to verify the validity of the theory as well as to illustrate the promising potentials of the proposed framework.

Index Terms

compressed sensing, compressive sensing, random projection, sparse reconstruction, fast and efficient algorithm

I. INTRODUCTION

Compressed sensing (CS) [1], [2] has attracted a lot of interests over the past few years as a revolutionary signal sampling paradigm. Suppose that \mathbf{x} is a length- N signal. It is said to be K -sparse (or compressible) if \mathbf{x} can be well approximated using only $K \ll N$ coefficients under some linear transform:

$$\mathbf{x} = \Psi\boldsymbol{\alpha},$$

where Ψ is the sparsifying basis and $\boldsymbol{\alpha}$ is the transform coefficient vector that has K (significant) nonzero entries.

According to the CS theory, such a signal can be acquired through the following random linear projection:

$$\mathbf{y} = \Phi\mathbf{x} + \mathbf{e},$$

where \mathbf{y} is the sampled vector with $M \ll N$ data points, Φ represents a $M \times N$ random matrix and \mathbf{e} is the acquisition noise. The CS framework is attractive as it implies that \mathbf{x} can be faithfully recovered from only $M = \mathcal{O}(K \log N)$ measurements, suggesting the potential of significant cost reduction in digital data acquisition.

While the sampling process is simply a random linear projection, the reconstruction to find the sparsest signal from the received measurements is highly non-linear process. More precisely, the reconstruction algorithm is to solve the l_1 -minimization of a transform coefficient vector:

$$\min \|\boldsymbol{\alpha}\|_1 \quad \text{s.t.} \quad \mathbf{y} = \Phi\Psi\boldsymbol{\alpha}.$$

Linear programming [1], [2] and other convex optimization algorithms [3], [4], [5] have been proposed to solve the l_1 minimization. Furthermore, there also exists a family of greedy pursuit algorithms [6], [7], [8], [9], [10] offering another promising option for sparse reconstruction. These algorithms all need to compute $\Phi\Psi$ and $(\Phi\Psi)^T$ multiple times. Thus, computational complexity of the system depends on the structure of sensing matrix Φ and its transpose Φ^T .

Preferably, the sensing matrix Φ should be highly incoherent with sparsifying basis Ψ , i.e. rows of Φ do not have any sparse representation in the basis Ψ . Incoherence between two matrices is mathematically quantified by the mutual coherence coefficient [11].

Definition I.1. The mutual coherence of an orthonormal matrix $N \times N$ Φ and another orthonormal matrix $N \times N$ Ψ is defined as:

$$\mu(\Phi, \Psi) = \max_{1 \leq i, j \leq N} |\langle \Phi_i, \Psi_j \rangle|$$

where Φ_i are rows of Φ and Ψ_j are columns of Ψ , respectively.

If Φ and Ψ are two orthonormal matrices, $\|\Phi\Psi_j\|_2 = \|\Psi_j\|_2 = 1$. Thus, it is easy to see that for two orthonormal matrices Φ and Ψ , $1/\sqrt{N} \leq \mu \leq 1$. Incoherence implies that the mutual coherence or the maximum magnitude of entries of the product matrix $\Phi\Psi$ is relatively small. Two matrices are completely incoherent if their mutual coherence coefficient approaches the lower bound value of $1/\sqrt{N}$.

A popular family of sensing matrices is a random projection or a random matrix of i.i.d random variables from a sub-Gaussian distribution such as Gaussian or Bernoulli [12], [13]. This family of sensing matrix is well-known as it is universally incoherent with all other sparsifying basis. For example, if Φ is a random matrix of Gaussian i.i.d entries and Ψ is an arbitrary orthonormal sparsifying basis, the sensing matrix in the transform domain $\Phi\Psi$ is also Gaussian i.i.d matrix. The universal property of a sensing matrix is important because it enables us to sense a signal directly in its original domain without significant loss of sensing efficiency and without any other prior knowledge. In addition, it can be shown that random projection approaches the optimal sensing performance of $M = \mathcal{O}(K \log N)$.

However, it is quite costly to realize random matrices in practical sensing applications as they require very high computational complexity and huge memory buffering due to their completely unstructured nature [14]. For example, to process a 512×512 image with $64K$ measurements (i.e., 25% of the original sampling rate), a Bernoulli random matrix requires nearly gigabytes storage and giga-flop operations, which makes both the sampling and recovery processes very expensive and in many cases, unrealistic.

Another class of sensing matrices is a uniformly random subset of rows of an orthonormal matrix in which the partial Fourier matrix (or the partial FFT) is a special case [13], [14]. While the partial FFT is well known for having fast and efficient implementation, it only works well in the transform domain or in the case that the sparsifying basis is the identity matrix. More specifically, it is shown in [[14], *Theorem 1.1*] that the minimal number of measurements required for exact recovery depends on the incoherence of Φ and Ψ :

$$M = \mathcal{O}(\mu_n^2 K \log N) \quad (1)$$

where μ_n is the normalized mutual coherence: $\mu_n = \sqrt{N}\mu$ and $1 \leq \mu_n \leq \sqrt{N}$. With many well-known sparsifying basis such as wavelets, this mutual coherence coefficient might be large and thus, resulting in performance loss. Another approach is to design a sensing matrix to be incoherent with a given sparsifying basis. For example, *Noiselets* is designed to be incoherent with the Haar wavelet basis in [15], i.e. $\mu_n = 1$ when Φ is Noiselets transform and Ψ

is the Haar wavelet basis. Noiselets also has low-complexity implementation $\mathcal{O}(N \log N)$ although it is unknown if noiselets is also incoherent with other bases.

II. COMPRESSIVE SENSING WITH STRUCTURALLY RANDOM MATRICES

A. Overview

One of remaining challenges for CS in practice is to design a CS framework that has the following features:

- *Optimal or near optimal sensing performance*: the number of measurements for exact recovery is almost minimal, i.e. on the order of $\mathcal{O}(K \log N)$;
- *Universality*: sensing performance is equally good with all sparsifying bases;
- *Low complexity, fast implementation that can support block-based processing*: this is necessary for large-scale, realtime sensing applications;
- *Easy and cheap to implement in hardware and optics domain*: Preferably, entries of the sensing matrix should only take values in the set $\{0, 1, -1\}$.

In this paper, we propose a framework that aims to satisfy the above wish-list. Lying at the heart of our framework is the concept of *Structurally Random Matrix*(SRM) that is defined as a product of three matrices:

$$\Phi = \sqrt{\frac{N}{M}} \mathbf{D} \mathbf{F} \mathbf{R} \quad (2)$$

where:

- $\mathbf{R} \in N \times N$ is either a uniform random permutation matrix or a diagonal random matrix whose diagonal entries R_{ii} are i.i.d Bernoulli random variables with identical distribution $P(R_{ii} = \pm 1) = 1/2$. A uniformly random permutation matrix scrambles signal's sample locations globally while a diagonal matrix of Bernoulli random variables flips signal's sample signs locally. Hence, we often refer the former as the *global randomizer* and the latter as the *local randomizer*.
- $\mathbf{F} \in N \times N$ is an orthonormal matrix that, in practice, is selected to be fast computable such as popular fast transforms: FFT, DCT, WHT or their block diagonal versions. The purpose of the matrix \mathbf{F} is to spread *information* (or energy) of the signal's samples over all measurements
- $\mathbf{D} \in M \times N$ is a subsampling matrix/operator. The operator \mathbf{D} selects a random subset of rows of the matrix $\mathbf{F} \mathbf{R}$. If the probability of selecting a row P (a row is selected) is M/N , the number of rows selected would be M in average. In matrix representation, \mathbf{D} is simply a random subset of M rows of the identity matrix of size $N \times N$. The scale coefficient $\sqrt{\frac{N}{M}}$ is to normalize the transform so that energy of the measurement vector is almost similar to that of the input signal vector.

The proposed sensing algorithm can be described step by step as follows:

- Step 1 (Signal pre-randomization): Randomize a target signal by either flipping its sample signs or uniformly permuting its sample locations. This step corresponds to multiplying the signal with the matrix \mathbf{D}
- Step 2 (Signal transform): Apply a fast transform \mathbf{F} to the randomized signal

- Step 3 (Signal subsampling): randomly pick up M measurements out of N transform coefficients. This step corresponds to multiplying the transform coefficients with the matrix \mathbf{D}

Conventional CS reconstruction algorithm is employed to recover the transform coefficient vector $\boldsymbol{\alpha}$ by solving the l_1 minimization:

$$\hat{\boldsymbol{\alpha}} = \operatorname{argmin} \|\boldsymbol{\alpha}\|_1 \quad \text{s.t.} \quad \mathbf{y} = \Phi \Psi \boldsymbol{\alpha}. \quad (3)$$

Finally, the signal is recovered as $\hat{\mathbf{x}} = \Psi \hat{\boldsymbol{\alpha}}$. The framework can achieve perfect reconstruction if $\hat{\mathbf{x}} = \mathbf{x}$.

From the best of our knowledge, the proposed sensing algorithm is distinct from currently existing methods such as random projection [16], random filters [17], structured Toeplitz [18], random convolution [19] via the step of pre-randomization. The main idea of this step is to deliberately scramble the structure of the signal, converting the signal to be sampled into a white noise-like one. Detail analysis in the following section will show that pre-randomization is necessary for obtaining universally incoherent sensing. The intuition behind this pre-randomization strategy is that scrambling a signal into a white noise-like form enables the sensing process to be independent of the signal's sparsifying basis.

The remaining of the paper is organized as follows. We first discuss about incoherence between SRMs and sparsifying transforms in Section III. More specifically, Section III-A will give us a rough intuition of why SRM could work as well as a random Gaussian matrix. Detail quantitative analysis of the incoherence for SRM with local randomizer and global randomizer is presented in Section III-B and Section III-C, respectively. Based on these incoherence results, theoretical performance of the proposed framework is analyzed in Section IV and then followed by experiment validation in Section V. Finally, Section VI concludes the paper with detail discussion of practical advantages of the proposed framework and relationship between the proposed framework and other related works.

B. Notations

We reserve a bold letter for a vector, a capital and bold letter for a matrix, a capital and bold letter with one sub-index for a row or a column of a matrix and a capital letter with two sub-indices for an entry of a matrix. We often employ $\mathbf{x} \in \mathbb{R}^N$ for the input signal, $\mathbf{y} \in \mathbb{R}^M$ for the measurement vector, $\Phi \in \mathbb{R}^{M \times N}$ for the sensing matrix, $\Psi \in \mathbb{R}^{N \times N}$ for the sparsifying matrix and $\boldsymbol{\alpha} \in \mathbb{R}^N$ for the transform coefficient vector ($\mathbf{x} = \Psi \boldsymbol{\alpha}$). We use the notation $\operatorname{supp}(\mathbf{z})$ to indicate the index set (or coordinate set) of nonzero entries of the vector \mathbf{z} . Occasionally, we also use \mathcal{T} to alternatively refer to this index set of nonzero entries (i.e., $\mathcal{T} = \operatorname{supp}(\mathbf{z})$). In this case, $\mathbf{z}_{\mathcal{T}}$ denotes the portion of vector \mathbf{z} indexed by the set \mathcal{T} and $\mathbf{A}_{\mathcal{T}}$ denotes the submatrix of \mathbf{A} whose columns are indexed by the set \mathcal{T} .

Let S_{ij} , F_{ij} be the entry at the i^{th} row and j^{th} column of $\mathbf{A}\Psi$ and \mathbf{F} , R_{kk} be the k^{th} entry on the diagonal of the diagonal matrix \mathbf{R} , \mathbf{A}_i and Ψ_j be the i^{th} row of \mathbf{A} and j^{th} column of Ψ , respectively.

In addition, we also employ the following notations:

- x_n is on the order of $o(z_n)$, denoted as $x_n = o(z_n)$, if

$$\lim_{n \rightarrow \infty} \frac{x_n}{z_n} = 0.$$

- x_n is on the order of $\mathcal{O}(z_n)$, denoted as $x_n = \mathcal{O}(z_n)$, if

$$\lim_{n \rightarrow \infty} \frac{x_n}{z_n} = c.$$

where c is some positive constant.

- A random variable X_n is called asymptotically normally distributed $\mathcal{N}(0, \sigma^2)$, if

$$\lim_{n \rightarrow \infty} P\left(\frac{X_n}{\sigma} \leq x\right) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{y^2}{2}} dy.$$

III. INCOHERENCE ANALYSIS

A. Asymptotical Distribution Analysis

If Φ is an i.i.d Gaussian matrix $\mathcal{N}(0, \frac{1}{N})$ and Ψ is an arbitrarily orthonormal matrix, $\Phi\Psi$ is also i.i.d Gaussian matrix $\mathcal{N}(0, \frac{1}{N})$, implying that with overwhelming probability, a Gaussian matrix is highly incoherent with all orthonormal Ψ . In other words, i.i.d. Gaussian matrices are universally incoherent with fixed transforms (with overwhelming probability). In this section, we will argue that under some mild conditions, with $\Phi = \mathbf{D}\mathbf{F}\mathbf{R}$, where $\mathbf{D}, \mathbf{F}, \mathbf{R}$ are defined as in the previous section, entries of $\Phi\Psi$ are asymptotically normally distributed $\mathcal{N}(0, \sigma^2)$, where $\sigma^2 \leq \mathcal{O}(\frac{1}{N})$. This claim is illustrated in Fig. 1, which depicts the quantile-quantile (QQ) plots of entries of $\Phi\Psi$, where $N = 256$, \mathbf{F} is the 256×256 DCT matrix and Ψ is the Daubechies-8 orthogonal wavelet basis. Fig. 1(a) and Fig. 1(b) correspond to the case \mathbf{R} is the local and global randomizer, respectively. In both cases, the QQ-plots appear straight, as the Gaussian model demands.

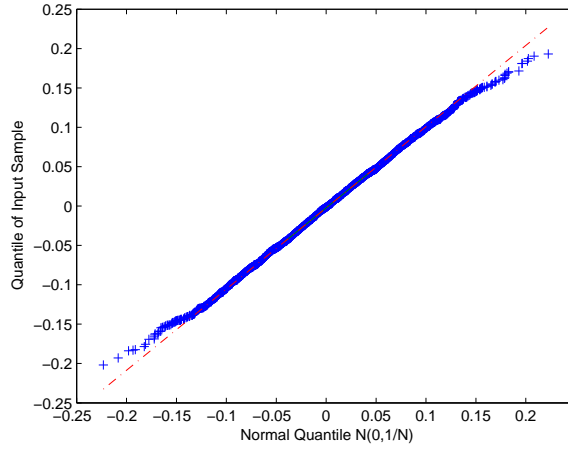
Note that Φ is a submatrix of $\mathbf{A} = \mathbf{F}\mathbf{R}$. Thus, asymptotical distribution of the entries of $\mathbf{A}\Psi$ is similar to that of entries of $\Phi\Psi$.

Theorem III.1. *Let $\mathbf{A} = \mathbf{F}\mathbf{R}$, where \mathbf{R} is an $N \times N$ random diagonal matrix of i.i.d Bernoulli random variables along its diagonal $P(R_{ii} = \pm 1) = 1/2$. Let \mathbf{F} be an $N \times N$ unit-norm row matrix with absolute magnitude of all entries on the order of $\mathcal{O}(\frac{1}{\sqrt{N}})$. Let Ψ be an $N \times N$ unit-norm column matrix with the maximal absolute magnitude of entries on the order of $o(1)$. Then, entries of $\mathbf{A}\Psi$ are asymptotically normally distributed $\mathcal{N}(0, \sigma^2)$ with $\sigma^2 \leq \mathcal{O}(\frac{1}{N})$.*

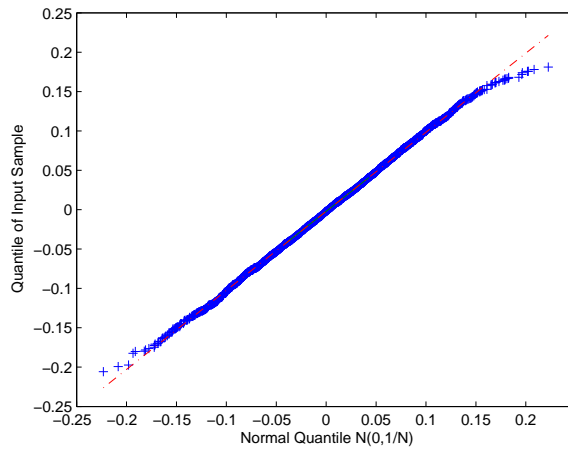
Proof. With notations being defined in Section II-B, we have:

$$S_{ij} = \langle \mathbf{A}_i, \Psi_j \rangle = \sum_{k=1}^N F_{ik} \Psi_{kj} R_{kk} \quad (4)$$

Denote $Z_k = F_{ik} \Psi_{kj} R_{kk}$. Because R_{kk} are i.i.d Bernoulli random variables, Z_k are i.i.d zero-mean random variables with $E(Z_k) = 0$. The assumption that $|F_{ik}|$ are on the order of $\mathcal{O}(\frac{1}{\sqrt{N}})$ implies that there exist two positive constants c_1 and c_2 such that:



(a)



(b)

Fig. 1. QQ plots comparing distribution of entries of $\Phi\Psi$ and Gaussian distribution. (a) \mathbf{R} is the local randomizer. (b) \mathbf{R} is the global randomizer. The plots all appear nearly linear, indicating that entries of $\Phi\Psi$ are nearly Normal distributed

$$\frac{c_1}{N} \Psi_{kj}^2 \leq \text{Var}(Z_k) = F_{ik}^2 \Psi_{kj}^2 \leq \frac{c_2}{N} \Psi_{kj}^2. \quad (5)$$

The variance of S_{ij} , σ^2 , can be bounded as the follows:

$$\frac{c_1}{N} = \frac{c_1}{N} \sum_{k=1}^N \Psi_{kj}^2 \leq \sigma^2 = \sum_{k=1}^N \text{Var}(Z_k) \leq \frac{c_2}{N} \sum_{k=1}^N \Psi_{kj}^2 = \frac{c_2}{N}. \quad (6)$$

Because S_{ij} is a sum of i.i.d zero-mean random variables $\{Z_k\}_{k=1}^N$, according to the Central Limit Theorem (CLT)(see Appendix I), $S_{ij} \rightarrow \mathcal{N}(0, \mathcal{O}(\frac{1}{N}))$. To apply CLT, we need to verify its convergence condition: for a given $\epsilon > 0$ and there exists N that is sufficiently large such that the $\text{Var}(Z_k)$ satisfy:

$$\text{Var}(Z_k) < \epsilon \sigma^2, k = 1, 2, \dots, N. \quad (7)$$

To show that this convergence condition is met, we use the counterproof method. Assume there exists ϵ_0 such that $\forall N$, there exists at least $k_0 \in \{1, 2, \dots, N\}$:

$$\text{Var}(Z_{k_0}) > \epsilon_0 \sigma^2. \quad (8)$$

From (5), (6) and (8), we achieve:

$$\epsilon \frac{c_1}{N} \leq \text{Var}(Z_{k_0}) \leq \frac{c_2}{N} \Psi_{k_0 j}^2. \quad (9)$$

This inequality can be true if all $\Psi_{k_0 j}$ are in the order of $o(1)$. The underlying intuition of the convergence condition is to guarantee that there is no random variable with dominant variance in the sum S_{ij} . In this case, it simply requires that there is no dominant entry on each column of Ψ . \square

Given a few more restrictions, we can establish a similar result when \mathbf{R} is a uniformly random permutation matrix,

Theorem III.2. *Let $\mathbf{A} = \mathbf{FR}$, where \mathbf{R} is an $N \times N$ uniformly random permutation matrix. Let \mathbf{F} be an $N \times N$ unit-norm row matrix with the maximal absolute magnitude of entries on the order of $\mathcal{O}(\frac{1}{\sqrt{N}})$. Let Ψ be an $N \times N$ unit-norm column matrix with the maximal absolute magnitude of entries on the order of $o(1)$ and the average sum of entries on each column is on the order of $o(\frac{1}{\sqrt{N}})$. Assume that sum of entries on each row of \mathbf{F} is zero. Also, assume that entries on each row of \mathbf{F} and on each column of Ψ are not all equal. Then, entries of $\mathbf{A}\Psi$ are asymptotically normally distributed $\mathcal{N}(0, \sigma^2)$, where $\sigma^2 \leq \mathcal{O}(\frac{1}{N})$.*

Proof. Let $[\omega_1, \omega_2, \dots, \omega_N]$ be a uniform random permutation of $[1, 2, \dots, N]$. Note that $\{\omega_k\}_{k=1}^N$ can be viewed as a sequence of random variables with identical distribution. In particular, for a fixed k :

$$P(\omega_k = i) = \frac{1}{N}, i = 1, 2, \dots, N.$$

Denote $Z_k = F_{i\omega_k} \Psi_{kj}$ (we omit the dependence of Z_k on i and j to simplify the notation), we have:

$$S_{ij} = \langle \mathbf{A}_i, \Psi_j \rangle = \sum_{k=1}^N F_{i\omega_k} \Psi_{kj} = \sum_{k=1}^N Z_k.$$

Using the assumption that the vector \mathbf{F}_i has zero average sum and unit norm, we derive:

$$E(Z_k) = \Psi_{kj} E(F_{i\omega_k}) = \frac{\Psi_{kj}}{N} \sum_{j=1}^N F_{ij} = 0.$$

and also,

$$E(Z_k^2) = \Psi_{kj}^2 E(F_{i\omega_k}^2) = \frac{\Psi_{kj}^2}{N} \sum_{j=1}^N F_{ij}^2 = \frac{\Psi_{kj}^2}{N}.$$

In addition, note that although $\{\omega_k\}_{k=1}^N$ have the identical distribution, they are correlated random variables because of the uniformly random permutation *without* replacement. Thus, with a pair of k and l such that $1 \leq k \neq$

$l \leq N$, we have:

$$\begin{aligned}
E(Z_k Z_l) &= \Psi_{kj} \Psi_{lj} E(F_{i\omega_k} F_{i\omega_l}) \\
&= \frac{\Psi_{kj} \Psi_{lj}}{N(N-1)} \sum_{1 \leq p \neq q \leq N} F_{ip} F_{iq} \\
&= \frac{\Psi_{kj} \Psi_{lj}}{N(N-1)} \left(\left(\sum_{p=1}^N F_{ip} \right)^2 - \sum_{p=1}^N F_{ip}^2 \right) \\
&= -\frac{\Psi_{kj} \Psi_{lj}}{N(N-1)}.
\end{aligned}$$

The last equation holds because the vector \mathbf{F}_i has zero average sum and unit-norm. Then, we derive the expectation and variance of S_{ij} as follows:

$$\begin{aligned}
E(S_{ij}) &= 0; \\
\text{Var}(S_{ij}) &= \sum_{k=1}^N E(Z_k^2) + \sum_{1 \leq k \neq l \leq N} E(Z_k Z_l) \\
&= \frac{1}{N} \sum_{k=1}^N \Psi_{kj}^2 - \frac{1}{N(N-1)} \sum_{1 \leq k \neq l \leq N} \Psi_{kj} \Psi_{lj} \\
&= \frac{1}{N} - \frac{1}{N(N-1)} \left(\left(\sum_{k=1}^N \Psi_{kj} \right)^2 - \sum_{k=1}^N \Psi_{kj}^2 \right) \\
&= \frac{1}{N} - \frac{1}{N(N-1)} \left(\sum_{k=1}^N \Psi_{kj} \right)^2 - 1 \\
&\leq \frac{1}{N} + \frac{1}{N(N-1)} = \mathcal{O}\left(\frac{1}{N}\right).
\end{aligned}$$

The forth equations holds because the column Ψ_j has unit-norm. The theorem is then a simple corollary of the Combinatorial Central Limit Theorem [20] (see Appendix 1), provided that its convergence condition can be verified that is:

$$\lim_{N \rightarrow \infty} N \frac{\max_{1 \leq k \leq N} (F_{ik} - \bar{F}_i)^2}{\sum_{k=1}^N (F_{ik} - \bar{F}_i)^2} \frac{\max_{1 \leq k \leq N} (\Psi_{kj} - \bar{\Psi}_j)^2}{\sum_{k=1}^N (\Psi_{kj} - \bar{\Psi}_j)^2} = 0, \quad (10)$$

where

$$\bar{F}_i = \frac{1}{N} \sum_{k=1}^N F_{ik}; \quad \bar{\Psi}_j = \frac{1}{N} \sum_{k=1}^N \Psi_{kj}.$$

Because $\bar{F}_i = 0$, $\|F_i\|_2^2 = 1$ and $\max_{1 \leq k \leq N} F_{ik}^2 = \mathcal{O}\left(\frac{1}{N}\right)$, the equation (10) holds if the following equation holds:

$$\lim_{N \rightarrow \infty} \frac{\max_{1 \leq k \leq N} (\Psi_{jk} - \bar{\Psi}_j)^2}{\sum_{k=1}^N (\Psi_{jk} - \bar{\Psi}_j)^2} = 0. \quad (11)$$

Because $\{|\bar{\Psi}_j|\}_{j=1}^N$ are on the order of $o\left(\frac{1}{\sqrt{N}}\right)$:

$$\sum_{k=1}^N (\Psi_{jk} - \bar{\Psi}_j)^2 = \|b\|_2^2 - N\bar{\Psi}_j^2 = 1 - N\bar{\Psi}_j^2 = \mathcal{O}(1). \quad (12)$$

Also, due to $|\bar{\Psi}_j| \leq \max_{1 \leq k \leq N} |\Psi_{jk}|$ and $|\Psi_{jk}|$ are on the order of $o(1)$:

$$\max_{1 \leq k \leq N} (\Psi_{jk} - \bar{\Psi}_j)^2 \leq 4 \max_{1 \leq k \leq N} \Psi_{jk}^2 = o(1). \quad (13)$$

Combination of (12) and (13) implies (11) and thus the convergence condition of the Combinatorial Central Limit Theorem is verified. \square

The condition that each row of \mathbf{F} has zero average sum is to guarantee that entries of $\mathbf{F}\Psi$ have zero mean while the condition that entries on each row of \mathbf{F} and on each column of Ψ are not all equal is to prevent the degenerate case that entries of $\mathbf{F}\Psi$ might become a deterministic quantity. For example, when entries of a row \mathbf{F}_i are all equal $\frac{1}{\sqrt{N}}$, $S_{ij} = \frac{1}{\sqrt{N}} \sum_{k=1}^N \Psi_{kj}$, which is a deterministic quantity, not a random variable. Note that these conditions are not necessary when \mathbf{R} is a diagonal matrix of Bernoulli random entries.

If \mathbf{F} is a DCT matrix, a (normalized) WHT matrix or a (normalized) DFT matrix, all the rows (except for the first one) have zero average sum due to the symmetry in these matrices. The first row, whose entries are all equal $\frac{1}{\sqrt{N}}$, can be considered as the averaging row, or a lowpass filtering operation. When the input signal is zero-mean, this row might be chosen or not without affecting quality of the reconstructed signal. Otherwise, it should be included in the chosen row set to encode the signal's mean. Lastly, the condition that absolute average sum of every column of the sparsifying basis Ψ are on the order of $o(\frac{1}{\sqrt{N}})$ is also close to the reality because the majority of columns of the sparsifying basis Ψ can be roughly viewed as bandpass and highpass filters whose average sum of the coefficients are always zero. For example, if Ψ is a wavelet basis (with at least one vanishing moment), then all columns of Ψ (except one at DC) has column sum of zero.

The aforementioned theorems show that under certain conditions, the majority of entries of $\mathbf{A}\Psi$ (also $\Phi\Psi$) behave like Gaussian random variables $\mathcal{N}(0, \sigma^2)$, where $\sigma^2 \leq \mathcal{O}(\frac{1}{N})$. Roughly speaking, this behavior constitutes to a good sensing performance for the proposed framework. However, these asymptotic results are not sufficient for an explicit measurement of sensing performance because in general, entries of $\mathbf{A}\Psi$ are not stochastically independent, violating a condition of a Gaussian i.i.d matrix. In fact, the sensing performance might be quantitatively analyzed by employing a powerful analysis framework of a random subset of rows of an orthonormal matrix [14]. Note that \mathbf{A} is also an orthonormal matrix when \mathbf{R} is either a random permutation matrix or a diagonal matrix of Bernoulli random entries.

Based on the Gaussian tail probability and a union bound for a supreme (i.e., maximum absolute value) of a random sequence, the maximum absolute magnitude of $\mathbf{A}\Psi$ can be asymptotically bounded as follows:

$$P(\max_{1 \leq i, j \leq N} |S_{ij}| \geq t) \leq 2N^2 \exp(-\frac{t^2}{2\sigma^2})$$

where $\sigma^2 \leq \frac{c}{N}$ and c is some positive constant and \preceq stands for "asymptotically smaller or equal", i.e., when N goes to infinity, \preceq becomes \leq .

If we choose $t = \sqrt{\frac{2c \log 2(N/\delta)^2}{N}}$, the above inequality is equivalent to:

$$P(\max_{1 \leq i, j \leq N} |S_{ij}| \leq \sqrt{\frac{c \log 2(N/\delta)^2}{N}}) \succeq 1 - \delta$$

which implies that with probability at least $1 - \delta$, the mutual coherence of \mathbf{A} and Ψ is upper bounded by $\mathcal{O}(\sqrt{\frac{\log(N/\delta)}{N}})$, which is close to the optimal value, except the $\log N$ factor.

In the following section, we will employ a more powerful tool from the theory of concentration inequalities to analyze the coherence between \mathbf{A} and Ψ when N is finite.

B. Incoherence Analysis With The Local Randomizer

The first theorem is about the mutual coherence of $\mathbf{A} = \mathbf{FR}$ and Ψ when \mathbf{R} is a diagonal matrix of i.i.d Bernoulli random variables.

Theorem III.3. *Let $\mathbf{A} = \mathbf{FR}$, where \mathbf{R} is an $N \times N$ random diagonal matrix of i.i.d Bernoulli random variables along its diagonal $P(R_{ii} = \pm 1) = 1/2$. Let \mathbf{F} be an $N \times N$ unit-norm row matrix with the maximal absolute magnitude of entries on the order of $\mathcal{O}(\frac{1}{\sqrt{B}})$, where $1 \leq B \leq N$. Let Ψ be an $N \times N$ unit-norm column matrix.*

- *With probability at least $1 - \delta$, the mutual coherence of \mathbf{A} and Ψ is upper bounded by $\mathcal{O}(\sqrt{\frac{\log(N/\delta)}{B}})$.*
- *In addition, if the maximal absolute magnitude of entries of Ψ is on the order of $\mathcal{O}(\frac{1}{\sqrt{N}})$, the mutual coherence is upper bounded by $\mathcal{O}(\sqrt{\frac{\log(N/\delta)}{N}})$, which is independent of B .*

Proof. A common proof strategy for this theorem as well as for other theorems in this paper is to establish a large deviation inequality that implies the quantity of our interest is concentrated around its expected value with high probability. Proof steps include:

- Showing that the quantity of our interest is a sum of independent random variables;
- Bounding the expectation and variance of the quantity;
- Applying a relevant concentration inequality of a sum of random variables;
- Applying a union bound for the supreme of a random sequence.

In this case, the quantity of interest is:

$$S_{ij} = \langle \mathbf{A}_i, \Psi_j \rangle = \sum_{k \in \text{supp}(\mathbf{F}_i)} F_{ik} \Psi_{kj} R_{kk}$$

Denote $Z_k = F_{ik} \Psi_{kj} R_{kk}$, for $k \in \text{supp}(\mathbf{F}_i)$ (in the support set of the row \mathbf{F}_i). Because R_{kk} are i.i.d Bernoulli random variables, Z_k are also i.i.d random variables with $E(Z_k) = 0$. Z_{kk} are also bounded because $Z_k = \pm F_{ik} \Psi_{kj}$

S_{ij} is a sum of independent, bounded random variables. Applying the Hoeffding's inequality (see Appendix 2) yields:

$$\Pr(|S_{ij}| \geq t) \leq 2 \exp\left(-\frac{t^2}{\sum_{k \in \text{supp}(\mathbf{F}_i)} F_{ik}^2 \Psi_{jk}^2}\right).$$

The next step is to evaluate $\sigma^2 = \sum_{k \in \text{supp}(\mathbf{F}_i)} F_{ik}^2 \Psi_{jk}^2$. Here, σ^2 can be roughly viewed as an approximation of the variance of S_{ij} .

$$\sigma^2 \leq \max_{1 \leq i, j \leq N} |F_{ij}|^2 \sum_{k \in \text{supp}(\mathbf{F}_i)} \Psi_{kj}^2 \leq \max_{1 \leq i, j \leq N} |F_{ij}|^2 = \frac{c}{B} \quad (14)$$

If the maximal absolute magnitude of entries of Ψ is on the order of magnitude of $\mathcal{O}(\frac{1}{\sqrt{N}})$:

$$\max_{1 \leq i, j \leq N} |\Psi_{ij}| = \frac{c}{\sqrt{N}},$$

where c is some positive constant, then

$$\sigma^2 \leq \max_{1 \leq i, j \leq N} |\Psi_{ij}|^2 \sum_{1 \leq k \leq N} F_{ik}^2 \leq \max_{1 \leq i, j \leq N} |\Psi_{ij}|^2 = \frac{c}{N}. \quad (15)$$

Finally, we derive an upper bound of the mutual coherence $\mu = \max_{1 \leq i, j \leq N} |S_{ij}|$ by taking a union bound for the supreme of a random sequence:

$$P\left(\max_{1 \leq i, j \leq N} |S_{ij}| \geq t\right) \leq 2N^2 \exp\left(\frac{-t^2}{\sigma^2}\right).$$

Choose $t = \sqrt{\sigma^2 \log(2N^2/\delta)}$, after simplifying the inequality, we get:

$$P\left(\max_{1 \leq i, j \leq N} |S_{ij}| \leq \sqrt{\sigma^2 \log(2N^2/\delta)}\right) \geq 1 - \delta.$$

Thus, with an arbitrary Ψ , (14) holds and we achieve the first claim of the Theorem:

$$P\left(\max_{1 \leq i, j \leq N} |S_{ij}| \leq \sqrt{\frac{c \log(2N^2/\delta)}{B}}\right) \geq 1 - \delta.$$

In the case that (15) holds, we achieve the second claim of the Theorem:

$$P\left(\max_{1 \leq i, j \leq N} |S_{ij}| \leq \sqrt{\frac{c \log(2N^2/\delta)}{N}}\right) \geq 1 - \delta.$$

□

Remark III.1. When \mathbf{A} is a popular transform such as the DCT or the normalized WHT, the maximal absolute magnitude of entries is on the order of $\mathcal{O}(\frac{1}{\sqrt{N}})$. As a result, the mutual coherence of the \mathbf{A} and an *arbitrary* Ψ is upper bounded by $\mathcal{O}(\sqrt{\frac{\log(N/\delta)}{N}})$, which is also consistent with our asymptotic analysis above. In other words, when at least Φ or Ψ is a *dense and uniform* matrix, i.e. the maximal absolute magnitude of their entries is on the order of $\mathcal{O}(\frac{1}{\sqrt{N}})$, their mutual coherence is nearly minimal, except the $\log N$ factor. Otherwise, mutual coherence between any arbitrary Ψ and a sparse matrix \mathbf{A} (e.g. block diagonal matrix of block size B) might be $\sqrt{\frac{N}{B}}$ times larger.

Cumulative coherence is more subtle way to quantify incoherence between two matrices [21].

Definition III.1. The cumulative coherence of an $N \times N$ \mathbf{A} and an $N \times K$ \mathbf{B} is defined as:

$$\mu_c(\mathbf{A}, \mathbf{B}) = \max_{1 \leq i \leq N} \sqrt{\sum_{1 \leq j \leq K} \langle \mathbf{A}_i, \mathbf{B}_j \rangle^2}$$

where \mathbf{A}_i and \mathbf{B}_j are rows of \mathbf{A} and columns of \mathbf{B} , respectively.

The cumulative coherence $\mu_c(\mathbf{A}, \mathbf{B})$ measures the *average* incoherence between two matrices \mathbf{A} and \mathbf{B} while mutual coherence $\mu(\mathbf{A}, \mathbf{B})$ measures the entry-wise incoherence. As a result, the cumulative coherence seems to be a better indicator of average sensing performance. In many cases, we are only interested in cumulative coherence between \mathbf{A} and Ψ_T , where T is the support of the transform coefficient vector. As will be shown in the following section, the cumulative coherence provides a more powerful tool to obtain a tighter bound of the number of measurements required for exact recovery.

From the definition of cumulative coherence, it is easy to verify that $\mu_c \leq \sqrt{K}\mu$. If we directly apply the result of the Theorem III.3, we obtain a trivial bound of the cumulative coherence: $\mu_c = \mathcal{O}(\sqrt{\frac{K \log N}{B}})$ for any arbitrary basis Ψ and $\mu_c = \mathcal{O}(\sqrt{\frac{K \log N}{N}})$ for any dense and uniform Ψ . In fact, we can get rid of the factor $\log N$ by directly measuring the cumulative coherence from its definition.

Theorem III.4. *Let $\mathbf{A} = \mathbf{FR}$, where \mathbf{R} is an $N \times N$ random diagonal matrix of i.i.d Bernoulli random variables along its diagonal $P(R_{ii} = \pm 1) = 1/2$. Let \mathbf{F} be an $N \times N$ unit-norm row matrix with the maximal absolute magnitude of entries is on the order of $\mathcal{O}(\frac{1}{\sqrt{B}})$, i.e. $\max_{1 \leq i, j \leq N} |F_{ij}| = \frac{c}{\sqrt{B}}$, where $1 \leq B \leq N$ and c is some positive constant. Let Ψ be an $N \times N$ unit-norm column matrix. With probability at least $1 - \delta$, the cumulative coherence of \mathbf{A} and $\Psi_{\mathcal{T}}$, where $|\mathcal{T}| = K$, is upper bounded by $\mathcal{O}(\sqrt{\frac{K}{B}})$ if $K > 16c^2 \log(2N/\delta)$.*

Proof. Denote $\mathbf{U} = \Psi_{\mathcal{T}}^*$ and \mathbf{U}_k are columns of \mathbf{U} . Let \mathbf{A}_i and Ψ_j ($j \in \mathcal{T}$) be rows of \mathbf{A} and columns of $\Psi_{\mathcal{T}}$, respectively.

$$S_i = \sqrt{\sum_{j \in \mathcal{T}} \langle \mathbf{A}_i, \Psi_j \rangle^2} = \|\mathbf{A}_i \Psi_{\mathcal{T}}\|_2 = \left\| \sum_{k \in \text{supp}(\mathbf{F}_i)} R_{kk} F_{ik} \mathbf{U}_k \right\|_2.$$

Denote $\mathbf{V}_k = F_{ik} \mathbf{U}_k$ and \mathbf{V} is the matrix of columns \mathbf{V}_k , $k \in \text{supp}(\mathbf{F}_i)$. First, we derive upper bound for the Frobenius of \mathbf{V} :

$$\|\mathbf{V}\|_F^2 \leq \max_{1 \leq i, j \leq N} F_{ij}^2 \|\mathbf{U}\|_F^2 = \frac{c^2 K}{B}.$$

The last equation holds because $\|\mathbf{U}\|_F^2 = K$. Also, the bound for the spectral norm is:

$$\begin{aligned} \|\mathbf{V}\|_2^2 &= \sup_{\|\boldsymbol{\beta}\|_2=1} \sum_{k \in \text{supp}(\mathbf{F}_i)} |\langle \boldsymbol{\beta}, \mathbf{V}_k \rangle|^2 \\ &= \sup_{\|\boldsymbol{\beta}\|_2=1} \sum_{k \in \text{supp}(\mathbf{F}_i)} F_{ik}^2 \left(\sum_{j=1}^K \boldsymbol{\beta}_j U_{kj} \right)^2 \\ &\leq \max_{1 \leq i, j \leq N} F_{ij}^2 \sup_{\|\boldsymbol{\beta}\|_2=1} \sum_{1 \leq k \leq N} |\langle \boldsymbol{\beta}, \mathbf{U}_k \rangle|^2 \\ &\leq \frac{c^2}{B} \|\mathbf{U}\|_2^2 = \frac{c^2}{B}. \end{aligned}$$

The last equation holds because $\|\mathbf{U}\|_2^2 = 1$. Now, we have:

$$S_i = \left\| \sum_{k \in \text{supp}(\mathbf{F}_i)} R_{kk} F_{ik} \mathbf{U}_k \right\|_2 = \left\| \sum_{k \in \text{supp}(\mathbf{F}_i)} R_{kk} \mathbf{V}_k \right\|_2.$$

Let us denote $\mathbf{Z} = \sum_{k \in \text{supp}(\mathbf{F}_i)} R_{kk} \mathbf{V}_k$.

\mathbf{Z} is a Rademacher sum of vectors and $S_i = \|\mathbf{Z}\|_2$ is a random variable. To show that S_i is concentrated around its expectation, we first derive bound of $E(\|\mathbf{Z}\|_2)$. It is easy to verify that for a random variable X , $E(X) \leq \sqrt{E(X^2)}$.

Thus, we will derive the upper bound for the simpler quantity $E(\|\mathbf{Z}\|_2^2)$

$$\begin{aligned} E(\|\mathbf{Z}\|_2^2) &= E(\mathbf{Z}^* \mathbf{Z}) = \sum_{k, l \in \text{supp}(\mathbf{F}_i)} E(R_{kk} R_{ll}) \langle \mathbf{V}_k, \mathbf{V}_l \rangle \\ &= \sum_{k \in \text{supp}(\mathbf{F}_i)} \langle \mathbf{V}_k, \mathbf{V}_k \rangle = \|\mathbf{V}\|_F^2 = \frac{c^2 K}{B}. \end{aligned}$$

The third equality holds because R_{kk} are i.i.d Bernoulli random variables and thus, $E(R_{kk}R_{ll}) = 0 \forall k \neq l$. As a result,

$$E(S_i) = E(\|\mathbf{Z}\|_2) \leq c\sqrt{\frac{K}{B}}.$$

Applying Ledoux's concentration inequality of the norm of a Rademacher sum of vectors [22] (see Appendix 2). Noting that $\|\mathbf{V}\|_2^2$ can be viewed as the variance of S_i , yields:

$$\Pr(S_i \geq c\sqrt{\frac{K}{B}} + t) \leq 2 \exp(-t^2 \frac{B}{16c^2})$$

Finally, apply a union bound for the supreme of a random process,we obtain:

$$\Pr(\max_{1 \leq i \leq N} S_i \geq c\sqrt{\frac{K}{B}} + t) \leq 2N \exp(-t^2 \frac{B}{16c^2}).$$

Choose $t = \sqrt{\frac{K}{B}}$. If $K > 16c^2 \log(2N/\delta)$, we get:

$$\Pr(\max_{1 \leq i \leq N} S_i \geq \mathcal{O}(\sqrt{\frac{K}{B}})) \leq \delta.$$

□

Remark III.2. When \mathbf{F} is some popular transform such as the DCT or the normalized WHT, the maximum absolute magnitude of entries is on the order of $\mathcal{O}(\frac{1}{\sqrt{N}})$. As a result, the cumulative coherence of \mathbf{A} and any arbitrary $\Psi_{\mathcal{T}}$, where $|\mathcal{T}| = K$, is upper bounded by $\mathcal{O}(\sqrt{\frac{K}{N}})$ if $K > 16c^2 \log(\frac{2N}{\delta})$, where c is some positive constant.

Remark III.3. The above theorem represents the worst-case analysis because Ψ can be an arbitrary matrix (the worst case corresponds to the case when Ψ is the identity matrix). When Ψ is known to be dense and uniform, the upper bound of cumulative coherence, according to the Theorem III.3 and the fact that $\mu_c \leq \mu\sqrt{K}$, is $\mathcal{O}(\sqrt{\frac{K \log N}{N}})$, which is, in general, better than $\mathcal{O}(\sqrt{\frac{K}{B}})$.

C. Incoherence Analysis With The Global Randomizer

The asymptotical analysis above reveals a significant technical difference for two cases: when \mathbf{R} is the local randomizer and when \mathbf{R} is the global randomizer. With the local randomizer, entries of $\mathbf{A}\Psi$ are sums of *independent* random variables while with global randomizer they are sums of *dependent* random variables. Stochastic dependence among random variables makes it much harder to set up similar arguments of their sum's concentration. In this case, we will show that the incoherence of \mathbf{A} and Ψ might depend on an extra quantity, the *heterogeneity coefficient* of the matrix Ψ .

Definition III.2. Assume Ψ is an $N \times N$ matrix. Let \mathcal{T}_k be the support of the column Ψ_k . Define:

$$\rho_k = \frac{\max_{1 \leq i \leq N} |\Psi_{ki}|}{\sqrt{\frac{1}{|\mathcal{T}_k|} \sum_{i \in \mathcal{T}_k} \Psi_{ki}^2}}. \quad (16)$$

The column-wise heterogeneity coefficient of the matrix Ψ is defined as:

$$\rho_{\Psi} = \max_{1 \leq k \leq N} \rho_k. \quad (17)$$

Obviously, $1 \leq \rho_k \leq \sqrt{|\mathcal{T}_k|}$. ρ_k illustrates the difference between the largest entry's magnitude and the average energy of *nonzero* entries. Roughly speaking, it indicates heterogeneity of nonzero entries of the vector Ψ_k . If nonzero entries of a column Ψ_k are homogeneous, i.e. they are on the same order of magnitude, ρ_k is on the order of a constant. If all nonzero entries of a matrix are homogeneous, the heterogeneity coefficient is also on the order of a constant, $C_{\Psi} = \mathcal{O}(1)$ and Ψ is referred as a uniform matrix. Note that a uniform matrix is not necessarily dense, for example, a block-diagonal matrix of DCT or WHT blocks

The following theorem indicates that when the global randomizer is employed, the mutual coherence between \mathbf{A} and Ψ is upper-bounded by $\mathcal{O}(\rho_{\Psi} \sqrt{\frac{\log(N/\delta)}{B}})$, where B is the block size of Φ and Ψ is an arbitrarily matrix with the heterogeneity coefficient ρ_{Ψ} .

Theorem III.5. *Let $\mathbf{A} = \mathbf{F}\mathbf{R}$, where \mathbf{R} is an $N \times N$ uniformly random permutation matrix. Let \mathbf{F} be an $N \times N$ unit-norm row matrix with the maximal absolute magnitude of entries is on the order of $\mathcal{O}(\frac{1}{\sqrt{B}})$. Assume that all rows of \mathbf{F} have zero average sum. Let Ψ be an $N \times N$ unit-norm column matrix with \mathcal{T}_k and ρ_{Ψ} defined as in (16) and (17). Assume that $\rho_k \geq 4 \log(2N^2/\delta) \forall k \in \{1, 2, \dots, N\}$.*

- *With probability at least $1 - \delta$, the mutual coherence of \mathbf{A} and Ψ is upper-bounded by $\mathcal{O}(\rho_{\Psi} \sqrt{\frac{\log(N/\delta)}{B}})$.*
- *In addition, if Ψ is dense and uniform, i.e. the maximum absolute magnitude of its entries is on the order of $\mathcal{O}(\frac{1}{\sqrt{N}})$ and $B \geq 4 \log(2N^2/\delta)$, the mutual coherence is upper-bounded by $\mathcal{O}(\sqrt{\frac{\log(N/\delta)}{N}})$, which is independent of B .*

Proof. Let $[\omega_1, \omega_2, \dots, \omega_N]$ be a uniformly random permutation of $[1, 2, \dots, N]$. With the same notation of $\mathbf{A}_i, \Psi_j, \mathbf{F}_i, S_{ij}, F_{ik}, \Psi_j$ as denoted previously:

$$S_{ij} = \langle \mathbf{A}_i, \Psi_j \rangle = \sum_{k=1}^N F_{i\omega_k} \Psi_{jk}.$$

As in the proof of the Theorem III.2, $\{\omega_k\}_{k=1}^N$ can be viewed as a sequence of dependent random variables with identical distribution, i.e. for a fixed $k \in \{1, 2, \dots, N\}$:

$$P(\omega_k = i) = \frac{1}{N}, \quad i \in \{1, 2, \dots, N\}.$$

The condition of \mathbf{F} is equivalent to $\max_{1 \leq i, j \leq N} |F_{ij}| = \frac{c}{\sqrt{B}}$, where c is some positive constant. Define $\{w_{k\omega_k}\}_{k=1}^N$ as the follows:

$$w_{k\omega_k} = \begin{cases} \frac{\sqrt{B|\mathcal{T}_k|}}{2c\rho_{\Psi}} F_{i\omega_k} \Psi_{jk} + \frac{1}{2} & \text{if } \Psi_{jk} \neq 0 \\ \frac{\sqrt{B|\mathcal{T}_k|}}{2c\rho_{\Psi}} F_{i\omega_k} \Psi_{jk} & \text{if } \Psi_{jk} = 0. \end{cases}$$

It is easy to verify that $0 \leq w_{k\omega_k} \leq 1$. Define W_k as the sum of dependent random variables $w_{k\omega_k}$

$$\begin{aligned} W_k &= \sum_{k=1}^N w_{k\omega_k} = \frac{\sqrt{B|\mathcal{T}_k|}}{2c\rho_\Psi} \sum_{k=1}^N F_{i\omega_k} \Psi_{jk} + \frac{|\mathcal{T}_k|}{2} \\ &= \frac{\sqrt{B|\mathcal{T}_k|}}{2c\rho_\Psi} S_{ij} + \frac{|\mathcal{T}_k|}{2}. \end{aligned}$$

Note that $\{F_{i\omega_k}\}_{k=1}^N$ are zero-mean random variables because \mathbf{F}_i has zero average sum. Thus, $E(S_{ij}) = 0$ and $E(W_k) = \frac{|\mathcal{T}_k|}{2}$. Then, applying the Sourav's theorem of concentration inequality for a sum of *dependent* random variables [23] (see Appendix 2) results in:

$$P\left\{\frac{\sqrt{B|\mathcal{T}_k|}}{2c\rho_\Psi} |S_{ij}| \geq \epsilon\right\} \leq 2 \exp\left(-\frac{\epsilon^2}{2|\mathcal{T}_k| + 2\epsilon}\right).$$

Denote $t = \frac{2c\rho_\Psi}{\sqrt{B|\mathcal{T}_k|}} \epsilon$. The above inequality is equivalent to:

$$P\{|S_{ij}| \geq t\} \leq 2 \exp\left(-\frac{B|\mathcal{T}_k|}{4c^2\rho_\Psi^2} \frac{t^2}{2|\mathcal{T}_k| + \frac{t}{c\rho_\Psi} \sqrt{B|\mathcal{T}_k|}}\right).$$

By choosing $t = 4c\rho_\Psi \sqrt{\frac{1}{B} \log(\frac{2N^2}{\delta})}$, we achieve:

$$P\{|S_{ij}| \geq t\} \leq 2 \exp\left(\frac{-4|\mathcal{T}_k| \log(\frac{2N^2}{\delta})}{2|\mathcal{T}_k| + 4\sqrt{|\mathcal{T}_k| \log(\frac{2N^2}{\delta})}}\right).$$

If $|\mathcal{T}_k| \geq 4 \log(\frac{2N^2}{\delta})$, the denominator inside the exponent is smaller than $4|\mathcal{T}_k|$. Thus,

$$P\{|S_{ij}| \geq 2c\rho_\Psi \sqrt{\frac{1}{B} \log(\frac{2N^2}{\delta})}\} \leq 2 \exp(-\log(\frac{2N^2}{\delta})) = \frac{\delta}{N^2}.$$

Finally, after taking the union bound for the supreme of a random sequence and simplifying the inequality, we obtain the first claim of the Theorem:

$$P\left\{\max_{1 \leq i, j \leq N} |S_{ij}| \leq \mathcal{O}(\rho_\Psi \sqrt{\frac{\log(N/\delta)}{B}})\right\} \geq 1 - \delta.$$

If Ψ is known to be dense and uniform, i.e. $\max_{1 \leq i, j \leq N} |\Psi_{ij}| = \frac{c_1}{\sqrt{N}}$, where c_1 is some positive constant. We then define $\{w_{k\omega_k}\}_{k=1}^N$ as the following:

$$w_{k\omega_k} = \begin{cases} \frac{\sqrt{BN}}{2cc_1} F_{ik} \Psi_{j\omega_k} + \frac{1}{2} & \text{if } F_{ik} \neq 0 \\ \frac{\sqrt{BN}}{2cc_1} F_{ik} \Psi_{j\omega_k} & \text{if } F_{ik} = 0. \end{cases}$$

Note that $0 \leq w_{k\omega_k} \leq 1$ and $E(w_{k\omega_k}) = \frac{B}{2}$. Repeat the same arguments above, we have:

$$P\{|S_{ij}| \geq t\} \leq 2 \exp\left(-\frac{NB}{4c^2c_1^2} \frac{t^2}{2B + \frac{t}{cc_1} \sqrt{NB}}\right).$$

Similarly, choose $t = 4cc_1 \sqrt{\frac{1}{N} \log(\frac{2N^2}{\delta})}$, we can derive:

$$P\{|S_{ij}| \geq t\} \leq 2 \exp\left(\frac{-4B \log(\frac{2N^2}{\delta})}{2B + 4\sqrt{B \log(\frac{2N^2}{\delta})}}\right).$$

If $B \geq 4 \log(\frac{2N^2}{\delta})$, the denominator inside the exponent is smaller than $4B$. Thus,

$$P\{|S_{ij}| \geq 2cc_1 \sqrt{\frac{1}{N} \log(\frac{2N^2}{\delta})}\} \leq \frac{\delta}{N^2}.$$

After taking the union bound of the supreme of a random sequence, we achieve the second claim of the Theorem. \square

Remark III.4. The first part of theorem implies that when \mathbf{F} is a dense and uniform matrix (e.g. DCT or normalized WHT) and $\mathbf{\Psi}$ is a uniform matrix (not necessarily dense), the mutual coherence closely approaches the minimum $\mathcal{O}(\sqrt{\frac{\log(N/\delta)}{N}})$. Although in this theorem, the mutual coherence depends on the heterogeneity coefficient, one will see in the experimental Section ?? that this dependence is almost negligible in practice.

As a consequence of this theorem, when at least \mathbf{A} or $\mathbf{\Psi}$ is dense and uniform, the mutual coherence of \mathbf{A} and $\mathbf{\Psi}$ is roughly on the order of $\mathcal{O}(\sqrt{\frac{\log N}{N}})$, which is quite close to the lower bound $\frac{1}{\sqrt{N}}$, except for the $\log N$ factor. Otherwise, the coherence linearly depends on the block size B of \mathbf{F} and is on the order of $\mathcal{O}(\sqrt{\frac{\log N}{B}})$. As a matter of fact, this bound is almost optimal because when $\mathbf{\Psi}$ is the identity matrix, the mutual coherence is actually equal the maximum absolute magnitude of entries of \mathbf{A} , which is on the order of $\mathcal{O}(\frac{1}{\sqrt{B}})$.

IV. COMPRESSIVE SAMPLING PERFORMANCE ANALYSIS

Section III demonstrates that under some mild conditions, the matrix \mathbf{A} and $\mathbf{\Psi}$ are highly incoherent, implying that the matrix $\mathbf{A}\mathbf{\Psi}$ is almost dense. When $\mathbf{A}\mathbf{\Psi}$ is dense, energy of nonzero transform coefficients α_T is distributed over all measurements. Commonly speaking, this is good for signal recovery from a small subset of measurements because if energy of some transform coefficients were concentrated in few measurements that happens to be bypassed in the sampling process, there is no hope for exact signal recovery even when employing the most sophisticated reconstruction method. This section shows that a random subset of rows of the matrix $\mathbf{A} = \mathbf{F}\mathbf{R}$ yields almost optimal measurement matrix $\mathbf{\Phi}$ for compressive sensing.

A. Main Assumptions for Theoretical Analysis

We first discuss main assumptions for theoretical results in the next section to hold. A signal \mathbf{x} is assumed to be sparse in some sparsifying basis $\mathbf{\Psi}$: $\mathbf{x} = \mathbf{\Psi}\alpha$, where the vector of transform coefficients α has no more than K nonzero entries. The sign sequence of nonzero transform coefficients α_T which is denoted as \mathbf{z} , is assumed to be a random vector of i.i.d Bernoulli random variables (i.e. $P(z_i = \pm 1) = \frac{1}{2}$). Let $\mathbf{y} = \mathbf{\Phi}\mathbf{x}$ be the measurement vector, where $\mathbf{\Phi} = \sqrt{\frac{N}{M}}\mathbf{D}\mathbf{F}\mathbf{R}$ is a structurally random matrix. For a general analysis, \mathbf{F} is assumed to be a block diagonal (and uniform) matrix with block size B ($1 \leq B \leq N$). If \mathbf{R} is the global randomizer, we also need the additional assumption that $\mathbf{\Psi}$ is uniform so that the theorem III.5 holds with the heterogeneity coefficient of $\mathbf{\Psi}$, $\rho_{\mathbf{\Psi}}$, is on the order of a constant. Note that the sensing operation can be equivalently accomplished via applying the Algorithm 1.

B. Theoretical Results

Theorem IV.1. *With probability at least $1 - \delta$, the proposed sensing framework can recover K -sparse signals exactly if the number of measurements $M \geq \mathcal{O}(\frac{N}{B} K \log^2(\frac{N}{\delta}))$. If \mathbf{F} is a dense and uniform rather than block-diagonal (e.g. DCT or normalized WHT matrix), the number of measurement needed is on the order of $\mathcal{O}(K \log^2(\frac{N}{\delta}))$.*

Proof. This is a simple corollary of the theorem of Candès et. al. [[14] Theorem 1.1] (1) because (i) $\mathbf{A} = \mathbf{FR}$ is an orthonormal matrix, and (ii) our incoherence results between \mathbf{A} and $\mathbf{\Psi}$ in the Theorem III.3 and Theorem III.5. \square

Remark IV.1. If $\mathbf{\Psi}$ is dense and uniform, the number of measurements for exact recovery is always $\mathcal{O}(K \log^2(\frac{N}{\delta}))$ regardless of the block size B . This implies that we can use the identity matrix for the transform \mathbf{F} ($B = 1$). For example, when input signal is known in advance to be spectrally sparse, compressively sampling it in the time domain is as efficient as in any other transform domain.

Compared with the framework that uses random projection, there is an upscale factor of $\log N$ for the number of measurements for exact recovery. In fact, by employing the above result of cumulative coherence, we can eliminate this upscale factor and thus, successfully showing optimal performance guarantee.

Theorem IV.2. *Assume that the sparsity $K > 16c^2 \log(\frac{2N}{\delta})$. With probability at least $1 - \delta$, the proposed framework employing the local randomizer can reconstruct K -sparse signals exactly if the number of measurements $M \geq \mathcal{O}(\frac{N}{B} K \log(\frac{N}{\delta}))$. If \mathbf{F} is a dense and uniform matrix (e.g. DCT or normalized WHT), the minimal number of required measurements is $M = \mathcal{O}(K \log(\frac{N}{\delta}))$.*

Proof. The proof is based on the result of cumulative coherence in the Theorem III.4 and a modification of the proof framework of the compressed sensing [14].

Denote $\mathbf{U} = \sqrt{\frac{N}{M}} \mathbf{FR}\mathbf{\Psi}$, $\mathbf{U}_{\mathcal{T}} = \sqrt{\frac{N}{M}} \mathbf{FR}\mathbf{\Psi}_{\mathcal{T}}$, $\mathbf{U}_{\Omega} = \sqrt{\frac{N}{M}} \mathbf{DFR}\mathbf{\Psi}$ and $\mathbf{U}_{\Omega\mathcal{T}} = \sqrt{\frac{N}{M}} \mathbf{DFR}\mathbf{\Psi}_{\mathcal{T}}$, where the support $\Omega = \{k | \mathbf{D}_{kk} = 1, k = 1, 2, \dots, N\}$. Let \mathbf{v}_k , $k \in \{1, 2, \dots, N\}$, be columns of $\mathbf{U}_{\mathcal{T}}^*$. Denote $\mu_c = \max_{1 \leq k \leq N} \|\mathbf{v}_k\|_2$, where $\mu_c = \mu_c(\mathbf{A}, \mathbf{\Psi}_{\mathcal{T}})$ is the cumulative coherence of $\mathbf{A} = \sqrt{\frac{N}{M}} \mathbf{FR}$ and $\mathbf{\Psi}_{\mathcal{T}}$. According to the above incoherence analysis, $\mu_c \leq \mathcal{O}(\sqrt{\frac{KN}{BM}})$. Also, denote μ as the mutual coherence of \mathbf{A} and $\mathbf{\Psi}_{\mathcal{T}}$, $\mu \leq \mathcal{O}(\sqrt{\frac{N \log N}{BM}})$.

As indicated in [12], [14], to show l_1 minimization exact recovery, it is sufficient to verify the *Exact Recovery Principle*.

Exact Recovery Principle. *With high probability, $|\pi_k| < 1$ for all $k \in \mathcal{T}^c$, where \mathcal{T}^c is the complementary set of the set \mathcal{T} and $\pi = \mathbf{U}_{\Omega}^* \mathbf{U}_{\Omega\mathcal{T}} (\mathbf{U}_{\Omega\mathcal{T}}^* \mathbf{U}_{\Omega\mathcal{T}})^{-1} \mathbf{z}$, where \mathbf{z} is the sign vector of nonzero transform coefficients $\boldsymbol{\alpha}_{\mathcal{T}}$.*

Also note that $\pi_k = \langle \nu_k (\mathbf{U}_{\Omega\mathcal{T}}^* \mathbf{U}_{\Omega\mathcal{T}})^{-1}, \mathbf{z} \rangle$, where ν_k is the k^{th} row of $\mathbf{U}_{\Omega}^* \mathbf{U}_{\Omega\mathcal{T}}$, for some $k \in \mathcal{T}^c$. The proof contains three major steps:

- *Claim 1* (Bound the norm of ν_k): With high probability, $\|\nu_k\|$ on the order of $\mathcal{O}(\mu_c)$
- *Claim 2* (Bound the spectral norm of $\mathbf{U}_{\Omega\mathcal{T}}^* \mathbf{U}_{\Omega\mathcal{T}}$): With probability $1 - \delta$, $\|\mathbf{U}_{\Omega\mathcal{T}}^* \mathbf{U}_{\Omega\mathcal{T}}\| \geq \frac{1}{2}$.

- *Claim 3* (Bound the norm of $w_k = \boldsymbol{\nu}_k(\mathbf{U}_{\Omega T}^* \mathbf{U}_{\Omega T})^{-1}$): With high probability, $\|w_k\|$ is on the order of $\mathcal{O}(\mu_c)$. Finally, exploiting the assumption that \mathbf{z} is a random vector of i.i.d Bernoulli random variables to show that with probability $1 - \mathcal{O}(\delta)$, $|\pi_k| = |\langle w_k, \mathbf{z} \rangle| < 1$

□

We first present proof for the Claim 1.

Proof. Let \mathbf{U}_k be columns of \mathbf{U} . For $k \in \mathcal{T}^c$:

$$\boldsymbol{\nu}_k = \frac{1}{M} \sum_{i=1}^N D_{ii} U_{ik} \mathbf{v}_i = \sum_{i=1}^N (D_{ii} - \frac{M}{N}) U_{ik} \mathbf{v}_i$$

where the second equality holds because $\sum_{i=1}^N U_{ik} \mathbf{v}_i = \mathbf{U}_T^* \mathbf{U}_k = 0$ that results from the orthogonality of columns of \mathbf{U} . Let $Z_i = (D_{ii} - \frac{M}{N})$. Because D_{ii} are i.i.d binary random variables with $P(D_{ii} = 1) = \frac{M}{N}$, Z_i are zero mean i.i.d random variables and $E(Z_i^2) = \frac{M}{N}(1 - \frac{M}{N})$. Let \mathbf{W} be the matrix of columns $\mathbf{W}_i = U_{ik} \mathbf{v}_i$, $i \in \{1, 2, \dots, N\}$. Then, $\boldsymbol{\nu}_k$ can be viewed as a random weighted sum of column vectors \mathbf{w}_i :

$$\boldsymbol{\nu}_k = \frac{1}{M} \sum_{i=1}^N Z_i \mathbf{W}_i$$

and $\|\boldsymbol{\nu}_k\|$ is a random variable. We have:

$$E(\|\boldsymbol{\nu}_k\|^2) = \sum_{1 \leq i, j \leq N} E(Z_i Z_j) \langle \mathbf{W}_i, \mathbf{W}_j \rangle = \sum_{1 \leq i \leq N} E(Z_i^2) \|\mathbf{W}_i\|^2,$$

where the last equality holds due to $E(Z_i Z_j) = 0$ if $i \neq j$. Thus,

$$\begin{aligned} E(\|\boldsymbol{\nu}_k\|^2) &= \frac{M}{N} (1 - \frac{M}{N}) \sum_{1 \leq i \leq N} V_{ki}^2 \|U_i\|^2 \\ &\leq \frac{M}{N} (1 - \frac{M}{N}) \mu_c^2 \sum_{1 \leq i \leq N} U_{ik}^2 \leq \mu_c^2. \end{aligned}$$

where the last inequality holds due to $\|\mathbf{U}_k\|^2 = \frac{N}{M}$. This implies that $E(\|\boldsymbol{\nu}_k\|) \leq \mu_c$. To show that $\|\boldsymbol{\nu}_k\|$ is concentrated around its mean, we use the Talagrand's theorem of concentration inequality [24]. First, we have:

$$\begin{aligned} \|\mathbf{W}\|_2^2 &= \sup_{\|\beta\|=1} \sum_{i=1}^N |\langle \beta, \mathbf{W}_i \rangle|^2 = \sup_{\|\beta\|=1} \sum_{i=1}^N U_{ik}^2 |\langle \beta, \mathbf{v}_i \rangle|^2 \\ &\leq \mu^2 \sup_{\|\beta\|=1} \sum_{i=1}^N |\langle \beta, \mathbf{v}_i \rangle|^2 = \mu^2 \|\mathbf{U}_T\|_2^2 = \frac{N}{M} \mu^2. \end{aligned}$$

where the last equation holds because $\|\mathbf{U}_T\|_2^2 = \frac{N}{M}$. Thus, we derive the upper bound of the variance σ^2 :

$$\sigma^2 = E(Z_k^2) \|\mathbf{W}\|_2^2 \leq \frac{M}{N} (1 - \frac{M}{N}) \frac{N}{M} \mu^2 \leq \mu^2.$$

In addition, it is obvious that $|Z_k| \leq 1$ and thus

$$B = \max_{1 \leq i \leq N} \|\mathbf{W}_i\|_2 \leq \mu \mu_c.$$

The Talagrand's theorem [24] (see Appendix 2) shows that:

$$P(\|\boldsymbol{\nu}_k\| - E(\|\boldsymbol{\nu}_k\|) \geq t) \leq 3 \exp\left(\frac{-t}{cB} \log\left(1 + \frac{Bt}{\sigma^2 + BE(\|\boldsymbol{\nu}_k\|)}\right)\right),$$

where c is some positive constant. Replacing $E(\|\boldsymbol{\nu}_k\|)$, σ^2 and B by their upper bounds in the right-hand side, we obtain:

$$P(\|\boldsymbol{\nu}_k\| - E(\|\boldsymbol{\nu}_k\|) \geq t) \leq 3 \exp\left(\frac{-t}{c\mu\mu_c} \log\left(1 + \frac{\mu\mu_c t}{\mu^2 + \mu\mu_c^2}\right)\right).$$

The next step is to simplify the right-hand side of the above inequality by replacing the denominator inside the log by two times the dominant term and note that $\log(1+x) \geq \frac{x}{2}$ when $x \leq 1$. In particular, there are two cases:

- Case 1: $\mu\mu_c^2 \geq \mu^2$ or equivalently, $\mu_c^2 \geq \mu$, denote $\bar{\sigma}^2 = \mu\mu_c^2$ and $t = a\bar{\sigma}$. If $\mu\mu_c t \leq 2\mu\mu_c^2$ or equivalently, $a \leq 2(1/\mu)^{\frac{1}{2}}$,

$$P(\|\boldsymbol{\nu}_k\| - E(\|\boldsymbol{\nu}_k\|) \geq t) \leq 3 \exp(-\gamma a^2).$$

- Case 2: $\mu^2 \geq \mu\mu_c^2$, denote $\bar{\sigma}^2 = \mu^2$ and $t = a\bar{\sigma}$. If $\mu\mu_c t \leq 2\mu^2$ or equivalently, $a \leq 2/\mu_c$

$$P(\|\boldsymbol{\nu}_k\| - E(\|\boldsymbol{\nu}_k\|) \geq t) \leq 3 \exp(-\gamma a^2).$$

where γ is some positive constant.

In conclusion, we just derive that

$$P(\|\boldsymbol{\nu}_k\| \geq \mu_c + a\bar{\sigma}) \leq 3 \exp(-\gamma a^2), \quad (18)$$

where γ is a positive constant and a is an arbitrary number that satisfies the above conditions. \square

The Theorem 1.2 in [14] shows that the *Claim 2* holds when $M \geq \mu_c^2 \max(c_1 \log K, c_2 \log(3/\delta))$, where c_1 and c_2 are some known positive constants.

Finally, we present proof for the *Claim 3*.

Proof. First, we show that:

$$P\left(\sup_{k \in \mathcal{T}^c} \|\mathbf{W}_k\| \geq 2\mu_c + 2a\bar{\sigma}\right) \leq 3N \exp(-\gamma a^2) + P(\|U_{\Omega\mathcal{T}}^* U_{\Omega\mathcal{T}}\| \leq \frac{1}{2}). \quad (19)$$

where $\mathbf{W}_k = \boldsymbol{\nu}_k (U_{\Omega\mathcal{T}}^* U_{\Omega\mathcal{T}})^{-1}$.

Let \mathcal{A} be the event that $\{\|U_{\Omega\mathcal{T}}^* U_{\Omega\mathcal{T}}\| \geq \frac{1}{2}\}$ or equivalently, $\{\|(U_{\Omega\mathcal{T}}^* U_{\Omega\mathcal{T}})^{-1}\| \leq 2\}$ and \mathcal{B} be the event that $\{\sup_{k \in \mathcal{T}^c} \|\boldsymbol{\nu}_k\| \leq \mu_c + a\bar{\sigma}\}$. Note that

$$\sup_{k \in \mathcal{T}^c} \|\mathbf{W}_k\| \leq \|(U_{\Omega\mathcal{T}}^* U_{\Omega\mathcal{T}})^{-1}\| \sup_{k \in \mathcal{T}^c} \|\boldsymbol{\nu}_k\|.$$

Thus,

$$P\left(\sup_{k \in \mathcal{T}^c} \|\mathbf{W}_k\| \geq 2\mu_c + 2a\bar{\sigma}\right) \leq P(\overline{\mathcal{A} \cap \mathcal{B}}) \leq P(\overline{\mathcal{A}}) + P(\overline{\mathcal{B}}).$$

Note that $P(\overline{\mathcal{B}}) \leq 3N \exp(-\gamma a^2)$ implies (19) holds.

The last step is to show that $\sup_{k \in \mathcal{T}^c} |\langle \mathbf{W}_k, \mathbf{z} \rangle| \leq 1$ with high probability. Note that because \mathbf{z} is assumed to be a vector of i.i.d Bernoulli random variables, $|\langle \mathbf{W}_k, \mathbf{z} \rangle|$ is concentrated around its zero mean. In particular, according to the Hoeffding's inequality:

$$P(|\langle \mathbf{W}_k, \mathbf{z} \rangle| \geq 1) \leq 2 \exp\left(-\frac{1}{2\|\mathbf{W}_k\|^2}\right).$$

$$\Rightarrow P(|\langle \mathbf{W}_k, \mathbf{z} \rangle| \geq 1 \mid \sup_{k \in \mathcal{T}^c} \|\mathbf{W}_k\| \leq \lambda) \leq 2N \exp\left(-\frac{1}{2\lambda^2}\right).$$

Note that with two arbitrary probabilistic events \mathcal{A} and \mathcal{B} :

$$P(\mathcal{A}) = P(\mathcal{A}|\mathcal{B})P(\mathcal{B}) + P(\mathcal{A}|\bar{\mathcal{B}})P(\bar{\mathcal{B}}) \leq P(\mathcal{A}|\mathcal{B}) + P(\bar{\mathcal{B}}).$$

Now, let \mathcal{A} be the event $\{\sup_{k \in \mathcal{T}^c} |\langle \mathbf{W}_k, \mathbf{z} \rangle| \geq 1\}$ and \mathcal{B} be the event $\{\sup_{k \in \mathcal{T}^c} \|\mathbf{W}_k\| \leq \lambda\}$, we can show that

$$P\left(\sup_{k \in \mathcal{T}^c} |\langle \mathbf{W}_k, \mathbf{z} \rangle| \geq 1\right) \leq 2N \exp\left(-\frac{1}{2\lambda^2}\right) + P\left(\sup_{k \in \mathcal{T}^c} \|\mathbf{W}_k\| \geq \lambda\right). \quad (20)$$

Choose $\lambda = 2\mu_c + 2a\bar{\sigma}$, according to (19) and (20), the probability of our interest $P(\sup_{k \in \mathcal{T}^c} |\langle \mathbf{W}_k, \mathbf{z} \rangle| \geq 1)$ is upper bounded by:

$$3N \exp(-\gamma a^2) + 2N \exp\left(-\frac{1}{2\lambda^2}\right) + \delta.$$

To show that $\{\sup_{k \in \mathcal{T}^c} |\langle \mathbf{W}_k, \mathbf{z} \rangle| \leq 1\}$ with probability $1 - \mathcal{O}(\delta)$, it is sufficient to show that the above upper bound is not greater than 3δ . In particular, choose $a^2 = \gamma^{-1} \log(3N/\delta)$ that makes the first term to be equal δ .

To make the second term less than δ , it is required that

$$\frac{1}{2\lambda^2} \geq \log\left(\frac{2N}{\delta}\right). \quad (21)$$

- Case 1: $\mu_c^2 \geq \mu$. The condition that (18) holds is $a \leq 2(1/\mu)^{\frac{1}{2}}$ that is equivalent to:

$$1 \geq \frac{1}{4} \gamma^{-2} \mu^2 \log^2(3N/\delta).$$

It is easy to see $\mu_c \geq a\bar{\sigma}$, where $\bar{\sigma} = (\mu\mu_c^2)^{1/2}$. In this case, $\lambda \leq 4\mu_c$. Thus, (21) holds if

$$1 \geq 32\mu_c^2 \log\left(\frac{2N}{\delta}\right). \quad (22)$$

- Case 2: $\mu \geq \mu_c^2$. The condition that (18) holds is $a \leq 2/\mu_c$ or equivalently,

$$1 \geq \frac{1}{4} \gamma^{-2} \mu_c^2 \log(3N/\delta).$$

If $\mu_c \geq a\bar{\sigma}$, where $\bar{\sigma} = \mu$, $\lambda \leq 4\mu_c$ and the condition is again (22). Otherwise, $\lambda \leq 4a\bar{\sigma}$. In this case, (21) holds if

$$1 \geq 32\gamma^{-1} \mu^2 \log\left(\frac{2N}{\delta}\right).$$

In conclusion, the Exact Recovery Principle is verified if $1 \geq \max(c_1 \mu^2 \log^2(3N/\delta), c_2 \mu_c^2 \log(3N/\delta))$, where c_1 and c_2 are known positive constants.

Finally, note that $\mu^2 \leq \mathcal{O}\left(\frac{N \log N}{BM}\right)$ and $\mu_c^2 \leq \mathcal{O}\left(\frac{NK}{BM}\right)$ and the assumption that $K \geq 16c^2 \log\left(\frac{2N}{\delta}\right)$, the sufficient condition for exact recovery is $M \geq \mathcal{O}\left(\frac{N}{B} K \log\left(\frac{N}{\delta}\right)\right)$. When \mathbf{F} is dense and uniform, the condition becomes $M \geq \mathcal{O}\left(K \log\left(\frac{N}{\delta}\right)\right)$. \square

TABLE I
SRMs EMPLOYED IN THE EXPERIMENT WITH SPARSE SIGNALS

Notation	\mathbf{R}	\mathbf{F}
WHT64-L	Local randomizer	64×64 block diagonal WHT
WHT64-G	Global randomizer	64×64 block diagonal WHT
WHT256-L	Local randomizer	256×256 block diagonal WHT
WHT256-G	Global randomizer	256×256 block diagonal WHT

V. NUMERICAL EXPERIMENTS

A. Simulation with Sparse Signals

In this section, we evaluate the sensing performance of several structurally random matrices and compare it with that of completely random projection. We also explore the connection among sensing performance (probability of exact recovery), streaming capacity (block size of \mathbf{F}) and structure of the sparsifying basis Ψ (e.g. sparsity and heterogeneity).

In the first simulation, the input signal \mathbf{x} of length $N = 256$ is sparse in the DCT domain, i.e. $\mathbf{x} = \Psi\boldsymbol{\alpha}$, where the sparsifying basis Ψ is the 256×256 IDCT matrix. Its transform coefficient vector $\boldsymbol{\alpha}$ has K nonzero entries whose magnitudes are Gaussian distributed and locations are at uniformly random, where $K \in \{10, 20, 30, 40, 50, 60\}$. With the signal \mathbf{x} , we generate a measurement vector of length $M = 128$: $\mathbf{y} = \Phi\mathbf{x}$, where Φ is some structurally random matrix or a completely Gaussian random matrix. SRMs under consideration are summarized in Table I.

The Orthogonal Matching Pursuit algorithm [6], is used to recover the signal from its measurements \mathbf{y} . For each value of sparsity $K \in \{10, 20, 30, 40, 50, 60\}$, we repeat the experiment 500 times and count the probability of exact recovery. The performance curve is plotted in Fig. 2(a). Numerical values on the x -axis denote signal sparsity K while those on the y -axis denote the probability of exact recovery. We then repeat similar experiments when an input signal is sparse in some sparse and non-uniform basis Ψ . Fig. 2(b) and Fig. 2(c) illustrate the performance curves when Ψ is the Daubechies-8 wavelet basis and the identity matrix, respectively.

These experiments verify that when performance of the SRM is comparable to that of a completely random matrix when the transform matrix \mathbf{F} is dense (all of its entries are non-zero) or when the sparsifying matrix Ψ of the input signal is dense (e.g. DCT). This implies that if we know the signal is sparse in a dense domain Ψ , we can sense the signal directly in its original domain (i.e., $\mathbf{F} = \mathbf{I}$) without performance loss. In addition, if we have no prior knowledge of a sparsifying transform, employing a SRM with the dense matrix \mathbf{F} guarantees optimal performance.

However, when both sparsifying matrix and SRM are sparse, sensing performance might drop quickly as illustrated in Fig. 2(c), revealing a trade-off between sensing performance and streaming capacity. In this case, Fig. 2(b) shows that the SRM with the global randomizer seems to work much better than the SRM with the local randomizer.

TABLE II
SRMS EMPLOYED IN THE EXPERIMENT WITH COMPRESSIBLE SIGNALS

Notation	R	F
DCT32-G	Global randomizer	32×32 block diagonal DCT
WHT32-G	Global randomizer	32×32 block diagonal WHT
DCT512-L	Local randomizer	512×512 block diagonal DCT
WHT512-L	Local randomizer	512×512 block diagonal WHT

B. Simulation with Compressible Signals

In this simulation, signals of interest are natural images of size 512×512 such as the 512×512 Lena, Barbara and Boat images. The sparsifying basis Ψ used for these natural images is the well-known Daubechies 9/7 wavelet transform. All images are implicitly regarded as 1-D signals of length 512^2 . The GPSR software in [3] is used for signal reconstruction.

For such a large scale simulation, it takes a huge amount of system resources to implement the sensing method of a completely random matrix. Thus, for the purpose of benchmark, we adopt a more practical scheme of partial FFT in the wavelet domain (WPFFT). The WPFFT is to sense wavelet coefficients in the wavelet domain using the method of partial FFT. Theoretically, WPFFT has optimal performance as the Fourier matrix is completely incoherent with the identity matrix. The WPFFT is a method of sensing a signal in the transform domain that also requires substantial amount of system resources. SRMs under consideration are summarized in Table II.

For the purpose of comparison, we also implement two popular sensing methods: partial FFT in the time domain (PFFT)[1] and the Scrambled/Permuted FFT (SFFT) in [25], [26] that corresponds to a dense SRM using the global randomizer.

The performance curves of these sensing ensembles are plotted in Fig. 3(a), Fig. 3(b) and Fig. 3(c), which correspond to the input signal Lena, Barbara and Boat images, respectively. Numerical value on the x -axis represents sampling rate, which is the number of measurements over the total number of samples. Value on y -axis is the quality of reconstruction (PSNR in dB). Lastly, Fig. 4 shows the visually reconstructed 512×512 Lena image from 25% of measurements using WPFFT, WHT32-G and WHT512-L ensembles

As clearly seen in Fig. 3, the PFFT is not an efficient sensing matrix for smooth signals like images because Fourier matrix and wavelet basis are highly coherent. On the other hand, the SRM method, which can roughly be viewed as the PFFT preceded by the pre-randomization process, is very efficient. In particular, with a dense SRM like SFFT, the performance difference between the SRM method and the benchmark one, WPFFT, is less than 1 dB. In addition, performance of DCT512-L and WHT512-L that are fully streaming capable SRM, degrades about 1.5 dB, which is a reasonable sacrifice as the buffer size required is less than 0.2 percent of the total length of the original signal. Less degradation is obtainable when the buffer size is increased. Also, in all cases, there is no

observable difference of performance between DCT and normalized WHT transforms. It implies that orthonormal matrices whose entries have the same order of absolute magnitude generate comparable performance. In addition, highly sparse SRM using the global randomizer such as DCT32-G and WHT32-G has experimental performance comparable to that of the dense SRM. Note that these SRM are highly sparse because their density are only 2^{-13} . This observation again verifies that a SRM using the global randomizer might, in general, outperform a SRM using the local randomizer. We leave the theoretical justification of this observation for our future research.

VI. DISCUSSION AND CONCLUSION

A. Complexity Discussion

We compare the computation and memory complexity between the proposed SRM and other random sensing matrices such as Gaussian or Bernoulli i.i.d. matrices. In implementation, the i.i.d Bernoulli matrix is obviously preferred than i.i.d Gaussian one as the former has integer entries 1, -1 and requires only 1 bit to represent each entry. A $M \times N$ i.i.d. Bernoulli sensing matrix requires MN bits for storing the matrix and MN additions and multiplications for sensing operation. A $M \times N$ structurally random matrix only requires $2N + N \log N$ bits for storage and $N + N \log N$ additions and multiplications for sensing operation. With SRM, its computational complexity and memory space required is independent with the number of measurements M . Note that with SRM, we do not need to store matrices \mathbf{D} , \mathbf{F} , \mathbf{R} explicitly. We only need to store the diagonals of \mathbf{D} and of \mathbf{R} and the fast transform \mathbf{F} , resulting in significant saving of both memory space and computational complexity.

Sparse signal recovery algorithms often require to compute \mathbf{A} and \mathbf{A}^T in each iteration for reconstructing the original sparse signal \mathbf{x} from the compressed measurement vector \mathbf{y} , where $\mathbf{A} = \mathbf{\Phi}\mathbf{\Psi}$. Speed of these reconstruction algorithm often depends critically on whether matrix-vector multiplications $\mathbf{A}\mathbf{u}$ and $\mathbf{A}^T\mathbf{u}$ can be computed quickly [3]. For the sake of simplicity, let's now assume $\mathbf{\Psi}$ is identity matrix. $\mathbf{A}\mathbf{u} = \mathbf{\Phi}\mathbf{u}$ requires $MN = \mathcal{O}(KN \log N)$ additions and multiplications for a random sensing matrix $\mathbf{\Phi}$ and $\mathcal{O}(N \log N)$ additions and multiplications for a SRM. This implies that SRM can speed up the reconstruction algorithm with at least K folds. With compressible signals (e.g., images), the number of measurements acquired tends to be proportional with the signal dimension, for example, $M = N/4$, then computational complexity reduction if using SRM is $\frac{N}{4 \log N}$.

Table III summarizes practical advantages of employing a SRM over a random sensing matrix.

B. Relationship with Other Related Works

When \mathbf{R} is a local randomizer, the SRM matrix is a little reminiscent to the so-called Fast Johnson-Lindenstrauss Transform (FJLT) [27]. However, the SRM is much easier and less expensive to implement due to the simplicity of the matrix \mathbf{D} . In FJLT, this matrix is a completely random matrix with sparse distribution. It is unknown if there exists an efficient implementation of such a sparse random matrix. As a result, SRMs are more appropriate for practical applications because of their simple implementation and its optimal performance guarantee.

In [25], [26], the Scrambled/Permuted FFT is experimentally proposed as a heuristic low-complexity sensing method that is efficient for sensing of a large signal. To the best of our knowledge, however, there has not been

TABLE III
PRACTICAL FEATURE COMPARISON

Features	SRMs	Completely Random Matrices
No. of measurements for exact recovery	$\mathcal{O}(K \log N)$	$\mathcal{O}(K \log N)$
Sensing complexity	$N \log N$	$\mathcal{O}(KN \log N)$
Reconstruction complexity at each iteration	$\mathcal{O}(N \log N)$	$\mathcal{O}(KN \log N)$
Implementation in hardware and optics	Very easy	Difficult
Fast computability	Yes	No
Block-based processing	Yes	No

any theoretical analysis for Scrambled FFT. It turns out the Scrambled FFT is superseded by our unified SRM framework.

Random Convolution convolving the input signal with a random pulse followed by randomly subsampling measurements is proposed in [19] as a promising sensing method for large scale, real signals. Although there are a few other methods that exploit the same idea of convolving a signal with a random pulse, for examples: Random Filter in [17] and Toeplitz structured sensing matrix in [18], only random convolution method can be shown to approach optimal sensing performance. The main difference among the random convolution method and other similar methods is while other methods such as Random Filter and Toeplitz-based CS methods subsample measurements structurally, the random convolution method applies a technique of randomly subsampling measurements that is also employed in the proposed SRM framework. In addition, in random convolution, randomness is introduced in the Fourier domain by randomizing phases of Fourier coefficients. These techniques help to decouple stochastic dependence among measurements and thus, enabling us to establish stronger claims of sensing performance.

Although sharing a few common features, the proposed SRM framework is distinct from all aforementioned methods, including random convolution. One of major differences is that in SRM, signal pre-randomization is performed directly in its original domain (via the global randomizer or the local randomizer), rather than in Fourier domain as in the random convolution method. As a result, the sensing system becomes much simpler and less expensive to implement (without performance loss). In addition, it extends the random convolution method by verifying that not only Fourier transform but also a wide variety of other popular fast transforms, especially ones that are easier to implement such as WHT, can be used to obtain optimal performance. Last but not least, the SRM framework presents a systematic method to design optimal and flexible sensing matrices with practical features.

APPENDIX I

Central Limit Theorem. Let Z_1, Z_2, \dots, Z_N be mutually independent random variables. Assume $E(Z_k) = 0$ and denote $\sigma^2 = \sum_{k=1}^N \text{Var}(Z_k)$. If for a given $\epsilon \geq 0$ and N sufficiently large, the following inequalities hold:

$$\text{Var}(Z_k) < \epsilon \sigma^2 \quad k = 1, 2, \dots, N$$

then distribution of the normalized sum $S = \sum_{k=1}^N Z_k$ converges to $\mathcal{N}(0, \sigma^2)$

Combinatorial Central Limit Theorem. Given two sequences $\{a_k\}_{k=1}^N$ and $\{b_k\}_{k=1}^N$. Assume the a_k are not all equal and b_k are also not all equal. Let $[\omega_1, \omega_2, \dots, \omega_N]$ be a uniform random permutation of $[1, 2, \dots, N]$. Denote $Z_k = a_{\omega_k}$ and

$$S = \sum_{k=1}^N Z_k b_k;$$

S is asymptotically normally distributed $\mathcal{N}(E(S), \text{Var}(S))$ if

$$\lim_{N \rightarrow \infty} N \frac{\max_{1 \leq k \leq N} (Z_k - \bar{Z})^2}{\sum_{k=1}^N (Z_k - \bar{Z})^2} \frac{\max_{1 \leq k \leq N} (b_k - \bar{b})^2}{\sum_{k=1}^N (b_k - \bar{b})^2} = 0;$$

where

$$\bar{b} = \frac{1}{N} \sum_{k=1}^N b_k \quad \text{and} \quad \bar{Z} = \frac{1}{N} \sum_{k=1}^N Z_k.$$

APPENDIX II

Hoeffding's Concentration Inequality. Suppose X_1, X_2, \dots, X_N are independent random variables and $a_k \leq X_k \leq b_k$ ($k = 1, 2, \dots, N$). Define a new random variable $S = \sum_{k=1}^N X_k$. Then for any $t > 0$

$$P(|S - E(S)| \geq t) \leq 2e^{-\frac{2t^2}{\sum_{k=1}^N (b_k - a_k)^2}}.$$

Ledoux's Concentration Inequality. Let $\{\eta_i\}_{1 \leq i \leq N}$ be a sequence of independent random variables such that $|\eta_i| \leq 1$ almost surely and $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N$ be vectors in Banach space. Define a new random variable: $S = \|\sum_{i=1}^N \eta_i \mathbf{v}_i\|$. Then for any $t > 0$,

$$P(S \geq E(S) + t) \leq 2 \exp\left(-\frac{t^2}{16\sigma^2}\right)$$

where σ^2 denote the variance of S and $\sigma^2 = \sup_{\|\mathbf{u}\| \leq 1} \sum_{i=1}^N |\langle \mathbf{u}, \mathbf{v}_i \rangle|^2$.

Talagrand's Concentration Inequality. Let Z_k be zero-mean i.i.d random variables and bounded $|Z_k| \leq \lambda$ and \mathbf{u}_k be column vectors of a matrix \mathbf{U} . Define a new random variable: $S = \|\sum_{i=1}^N Z_k \mathbf{u}_k\|$. Then for any $t > 0$:

$$P(S \geq E(S) + t) \leq 3 \exp\left(-\frac{t}{cB} \log\left(1 + \frac{Bt}{\sigma^2 + BE(S)}\right)\right)$$

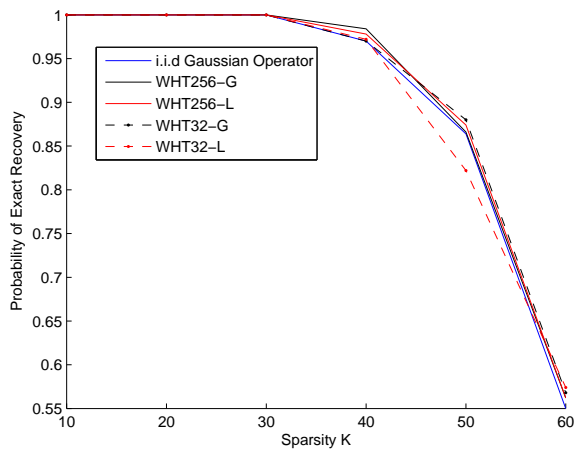
where c is some constant, variance $\sigma^2 = E(Z_k^2) \|\mathbf{U}\|^2$ and $B = \lambda \max_{1 \leq k \leq N} \|\mathbf{u}_k\|$.

Sourav's Concentration Inequality. Let $\{Z_{ij}\}_{1 \leq i, j \leq N}$ be a collection of numbers from $[0, 1]$. Let $[\omega_1, \omega_2, \dots, \omega_N]$ be a uniformly random permutation of $[1, 2, \dots, N]$. Define a new random variable: $S = \sum_{i=1}^N Z_{i\omega_i}$. Then for any $t \geq 0$

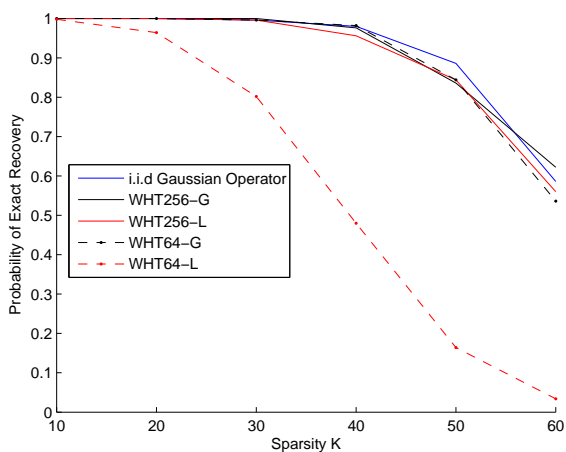
$$P(|S - E(S)| \geq t) \leq 2 \exp\left(-\frac{t^2}{4E(S) + 2t}\right).$$

REFERENCES

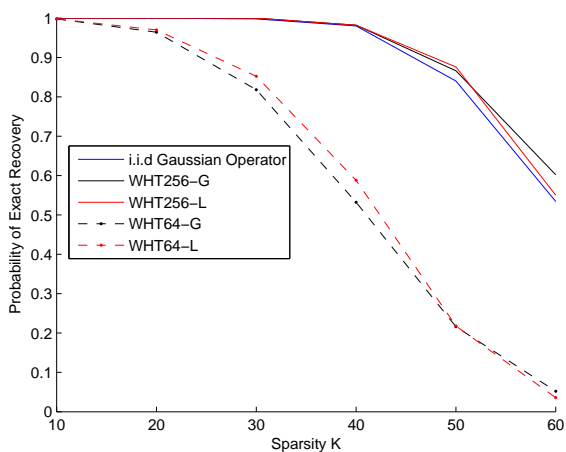
- [1] E. Candès, J. Romberg, and T. Tao, "Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information," *IEEE Trans. on Information Theory*, vol. 52, pp. 489 – 509, Feb. 2006.
- [2] D. L. Donoho, "Compressed sensing," *IEEE Trans. on Information Theory*, vol. 52, pp. 1289 – 1306, Apr. 2006.
- [3] M. A. T. Figueiredo, R. D. Nowak, and S. J. Wright, "Gradient projection for sparse reconstruction," *to appear in IEEE Journal of Selected Topics in Signal Processing*, 2007.
- [4] Elaine T. Hale, Wotao Yin, , and Yin Zhang, "A fixed-point continuation method for ℓ_1 -regularized minimization with applications to compressed sensing," *Technical Report*, Jul 2007.
- [5] Ewout V. D. Berg and Michael P. Friedlander, "Probing the pareto frontier for basis pursuit solutions," *Technical Report*, Jan 2008.
- [6] J. Tropp and A. Gilbert, "Signal recovery from random measurements via orthogonal matching pursuit," *IEEE Trans. Info. Theory*, vol. 53, pp. 4655–4666, Dec 2007.
- [7] D. Needell and J. A. Tropp, "Cosamp: Iterative signal recovery from incomplete and inaccurate samples," *Preprint*, Mar 2008.
- [8] W. Dai and O. Milenkovic, "Subspace pursuit for compressive sensing: Closing the gap between performance and complexity," *Preprint*, Mar 2008.
- [9] D. L. Donoho, Y. Tsaig, and Jean-Luc Starck, "Sparse solution of underdetermined linear equations by stagewise orthogonal matching pursuit," *Technical Report*, Mar. 2006.
- [10] Thong T. Do, Lu Gan, Nam Nguyen, and Trac D. Tran, "Sparsity adaptive matching pursuit algorithm for practical compressed sensing," *Asilomar Conference on Signals, Systems, and Computers, Pacific Grove, California*, Oct 2008.
- [11] D. L. Donoho and X. Huo, "Uncertainty principles and ideal atomic decomposition," *IEEE Trans. on Information Theory*, vol. 47(7), Nov. 2001.
- [12] E. Candès and T. Tao, "Near optimal signal recovery from random projections: Universal encoding strategies?," *IEEE Trans. on Information Theory*, vol. 52, pp. 5406 – 5425, Dec. 2006.
- [13] S. Mendelson, A. Pajor, and N. Tomczak-Jaegermann, "Uniform uncertainty principle for bernoulli and subgaussian ensembles.," *Preprint*, Aug. 2006.
- [14] E. Candès and J. Romberg, "Sparsity and incoherence in compressive sampling," *Inverse Problems*, vol. 23(3), pp. 969–985, 2007.
- [15] R. Coifman, F. Geshwind, and Y. Meyer, "Noiselets," *Appl. Comp. Harmonic Analysis*, vol. 10, pp. 27–44, 2001 2005.
- [16] E. Candès and T. Tao, "Decoding by linear programming," *IEEE Trans. on Information Theory*, vol. 51, pp. 4203–4215, Dec. 2005.
- [17] J. Tropp, M. Wakin, M. Duarte, D. Baron, and R. Baraniuk, "Random filters for compressive sampling and reconstruction," *Proc. IEEE ICASSP*, vol. 3, pp. 872–875, Toulouse, May 2006.
- [18] Waheed Bajwa, Jarvis Haupt, Gil Raz, Stephen Wright, and Robert Nowak, "Toeplitz-structured compressed sensing matrices," *IEEE Workshop on Statistical Signal Processing (SSP), Madison, Wisconsin.*, Aug 2007.
- [19] Justin Romberg, "Compressive sensing by random convolution.," *Preprint*, Jul 2008.
- [20] W. Hoeffding, "A combinatorial central limit theorem.," *The Annals of Mathematical Statistics*, vol. 22(4), pp. 558–566, Dec. 1951.
- [21] K. Schnass and P. Vandergheynst, "Average performance analysis for thresholding," *IEEE Signal Processing Letters*, vol. 14, Nov 2007.
- [22] M. Ledoux, "The concentration of measure phenomenon," *American Mathematical Society*, 2001.
- [23] Sourav Chatterjee, "Stein's method for concentration inequalities," *Probab. Theory Related Fields*, vol. 138, pp. 305–321, 2007.
- [24] M. Talagrand, "New concentration inequalities in product spaces," *Invent. Math.*, vol. 126, pp. 505–563, 1996.
- [25] E. Candès, J. Romberg, and T. Tao, "Stable signal recovery from incomplete and inaccurate measurements," *Communications on Pure and Applied Mathematics*, vol. 59, pp. 1207–1223, Aug. 2006.
- [26] M. F. Duarte, M. B. Wakin, and R. G. Baraniuk, "Fast reconstruction of piecewise smooth signals from incoherent projections," *SPARS'05*, Rennes, France, Nov 2005.
- [27] Nir Ailon and Bernard Chazelle, "Approximate nearest neighbors and the fast johnsonlindenstrauss transform," *Proceedings of the thirty-eighth annual ACM symposium on Theory of computing*, vol. 66, pp. 557 – 563, 2006.



(a)

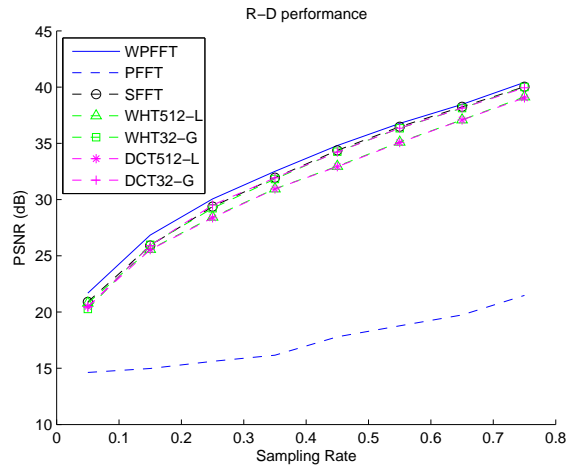


(b)

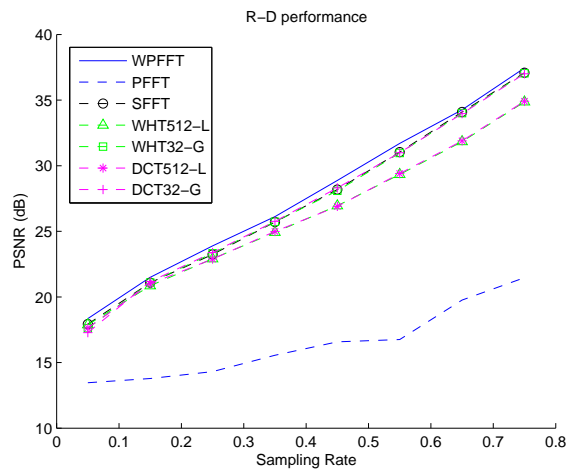


(c)

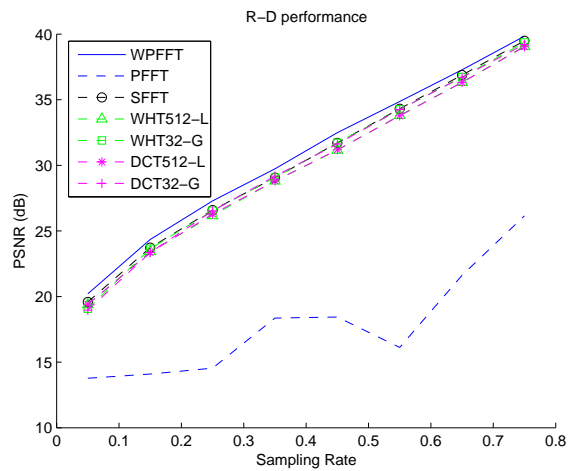
Fig. 2. Performance curves: probability of exact recovery vs. Sparsity K . (a) when Ψ is IDCT basis. (b) when Ψ is Daubechies-8 wavlet basis. (c) when Ψ is the identity basis



(a)



(b)



(c)

Fig. 3. Performance curves: Quality of signal reconstruction vs. sampling rate M/N . (a) the 512×512 Lena image. (b) the 512×512 Barbara image. (c) the 512×512 Boat image



(a)



(b)



(c)



(d)

Fig. 4. Reconstructed 512×512 *Lena* images from $M/N = 25\%$ sampling rate. (a) The original *Lena* image; (b) using the WPFFFT ensemble: 30.1dB; (c) using the WHT32-G ensemble: 29.3dB; (d) using the WHT512-L: 28.5dB