

Image representation and compression via sparse solutions of systems of linear equations

Alfredo Nava-Tudela^{a,b,*}, John J. Benedetto^b

^a*Institute for Physical Science and Technology, University of Maryland
College Park, MD 20742, USA*

^b*Norbert Wiener Center, Department of Mathematics, University of Maryland
College Park, MD 20742, USA*

Abstract

We are interested in finding sparse solutions to systems of linear equations $\mathbf{Ax} = \mathbf{b}$, where \mathbf{A} is underdetermined and full-rank. Our approach and theme is phenomenological and computational.

We use and test orthogonal matching pursuit (OMP) to study the compression properties of \mathbf{A} in the context of image processing. We follow the common technique of image blocking used in the JPEG and JPEG 2000 standards. We modify the stopping criteria of OMP, which results in better bit-rate vs image quality than with the usual ℓ^2 stopping rule. The image quality is measured by structural similarity (SSIM) and mean structural similarity (MSSIM) indices. These indices capture perceptual image quality, and compare them with the traditional peak signal-to-noise ratio.

We analyze various matrices whose column vectors come from the concatenation of waveforms based on the discrete cosine transform (DCT) and the Haar wavelet. We implement multiple vectorization algorithms and characterize their performance with respect to compression.

We compute the distortion properties of the image compression and representation methodology defined by the aforementioned matrices and blocking, respectively. In this setting, we propose a quantization/coding scheme, framed in the context of transform-based encoding, which allows us to compute bounds on the distortion measure. Combined with the statistics of

*Corresponding author

Email addresses: ant@umd.edu (Alfredo Nava-Tudela), jjb@umd.edu (John J. Benedetto)

natural images, this allows us to compute the normalized average distortion rate.

Keywords: Image processing, compression, sparsity, orthogonal matching pursuit, structural similarity index, peak signal-to-noise ratio, wavelet-spectral matrices.

1. Introduction

1.1. Background

In recent years, interest has grown in the study of sparse solutions of underdetermined systems of linear equations because of their many and potential applications [1]. In particular, these types of solutions can be used to describe images in a compact form, provided one is willing to accept an imperfect representation.

We are interested in studying this approach to image compression because of the ever-increasing volume of images in use by many multimedia channels. The basic idea is the following. Suppose that we have a full-rank matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$, where $n < m$, and that we want to find solutions to the equation,

$$\mathbf{Ax} = \mathbf{b}, \tag{1}$$

where \mathbf{b} is a given “signal.” Since the matrix \mathbf{A} is full-rank, and we have more unknowns than equations, we have an infinite number of solutions to Equation (1). What if from all possible solutions we could find \mathbf{x}_0 , the “sparsest” one, in the sense of having the least number of nonzero entries? Then, if the number of nonzero entries in \mathbf{x}_0 happens to be less than the number of nonzero entries in \mathbf{b} , we could store \mathbf{x}_0 instead of \mathbf{b} , achieving a representation of the original signal \mathbf{b} in a compressed way. This notion of *compression* is our point of view.

There are many questions that arise in this approach for image representation and compression. For example, is there a unique “sparsest” solution of Equation (1)? How does one find such a solution? What are the practical implications of this approach to image compression?

We note that this idea can be framed in the context of signal transform compression techniques. For example, the JPEG and JPEG 2000 standards have at their core transformations that result in different representations of the original image which can be truncated to achieve compression at the expense of some acceptable error [2, 3, 4].

1.2. Finding sparse solutions

The *orthogonal matching pursuit* (OMP) algorithm is one of the existing techniques to find sparse solutions of systems of linear equations, as in Equation (1). This is one of many greedy algorithms that attempts to solve the general problem,

$$(P_0^\epsilon) : \quad \min_{\mathbf{x}} \|\mathbf{x}\|_0 \quad \text{subject to} \quad \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2 < \epsilon. \quad (2)$$

Here, $\|\mathbf{x}\|_0 = \#\{j : |x_j| > 0\}$ is the “zero-norm” of vector \mathbf{x} , which counts the number of nonzero entries in \mathbf{x} . A greedy algorithm approach is necessary for the solution of the optimization problem defined in (2) because this is an NP-complete problem [5]. Moreover, it can be proven that under certain circumstances there is a unique sparsest solution to (P_0^ϵ) ; and, under those same circumstances, OMP is then guaranteed to find it [1].

1.3. Image representation

To make a practical implementation of the compression idea described in Section 1.1, we followed the approach used in the JPEG and JPEG 2000 standards [2, 3]. Also, in order to test our image representation and compression technique, we selected a set of four commonly used test images, three of them from [6].

All images in our image database are 512×512 , 8-bit depth, grayscale images. We proceeded to partition each image in subsets of 8×8 non-overlapping sub-images to process them individually. Partitioning a signal to work with more manageable pieces is a common technique [7]. Then, we vectorize each 8×8 sub-image into a vector $\mathbf{b} \in \mathbb{R}^{64}$ to be used as a right hand side in Equation (1). There are many ways to do this vectorization and we investigate three of them; see Sections 3.2.2 and 3.6.

To complete the setup, we need to choose a matrix \mathbf{A} . We begin with $\mathbf{A} = [\text{DCT}_1 \text{ Haar}_1] \in \mathbb{R}^{64 \times 128}$, where DCT_1 is a basis of one-dimensional discrete cosine transform waveforms and Haar_1 is a basis of Haar wavelets, resp.; see Section 3.3.2 for details. That is, we concatenate two bases of \mathbb{R}^{64} , since $\mathbf{b} \in \mathbb{R}^{64}$. We also consider bases for \mathbb{R}^{64} built from tensor products of the one-dimensional waveforms constructed analogously to the columns of DCT_1 and Haar_1 . This allows us to capture the two-dimensional nature of an image. These specific matrices $\mathbf{A} = \mathbf{A}_j = [\text{DCT}_{2,j} \text{ Haar}_{2,j}]$, $j = 1, 2, 3$, are defined in Section 3.3.3.

1.4. Compression results

Our theme is phenomenological and computational, as opposed to being theoretical.

We first compare the compression properties of $\mathbf{A} = [\text{DCT}_1 \text{ Haar}_1]$ to those of $\mathbf{B} = [\text{DCT}_1]$ and $\mathbf{C} = [\text{Haar}_1]$. We found that combining bases results in a representation \mathbf{x}_0 for each sub-image that requires fewer nonzero entries than for either the cases \mathbf{B} or \mathbf{C} , individually. We obtained similar compression results for the two-dimensional bases. A comparison between the one-dimensional and two-dimensional concatenation of bases shows that there is a range of tolerance values, which we have specified, in which the one-dimensional waveforms, resp., two-dimensional, outperform the two-dimensional waveforms, resp., one-dimensional; see Section 3.7.

Further, we have measured the quality of the reconstruction from these compressed representations with the *peak signal-to-noise ratio* (PSNR) [3, 8], the *structural similarity index* (SSIM), and the *mean structural similarity index* (MSSIM) [9], all as functions of the tolerance ϵ chosen when solving (P_0^ϵ) . This led us to a novel modification of the termination criteria in the OMP algorithm in order to achieve a desired PSNR or MSSIM value for the resulting decompressed image; see Sections 3.4.2 and 3.8.

The full potential of this approach to image representation and compression cannot be assessed until a bit-stream \mathbf{c} is produced to be able to quantify the net compression ratio, among other criteria; see Sections 3.4.1 and 3.7. These results must be interleaved into the structure of general transmission/storage systems in the context of the information-theoretical paradigms raised by the work of Shannon [10].

Our final results take advantage of the structure of our image compression approach in order to compute distortion properties; see Section 4.3

We have set the stage for addressing rigorously the natural questions that arise given our phenomenological/computational theme.

2. Finding sparse solutions with orthogonal matching pursuit

We give a brief description of one of the many greedy algorithms (GAs) that attempt to solve (P_0^ϵ) . Starting from $\mathbf{x}^0 = \mathbf{0}$, a greedy strategy iteratively constructs a k -term approximation \mathbf{x}^k by maintaining a set of active columns—initially empty—and, at each stage, expanding that set by one additional column. The column chosen at each stage maximally reduces the residual ℓ^2 error in approximating \mathbf{b} from the current set of active columns.

After constructing an approximation including the new column, the residual error ℓ^2 is evaluated. When it falls below a specified threshold $\epsilon > 0$, the algorithm terminates.

Orthogonal Matching Pursuit (OMP):

Task: Find the solution to $(P_0^\epsilon) : \min_{\mathbf{x}} \|\mathbf{x}\|_0$ subject to $\|\mathbf{Ax} - \mathbf{b}\|_2 < \epsilon$.

Parameters: Given \mathbf{A} , \mathbf{b} , and ϵ .

Initialization: Initialize $k = 0$, and set the following:

- The initial solution $\mathbf{x}^0 = \mathbf{0}$;
- The initial residual $\mathbf{r}^0 = \mathbf{b} - \mathbf{Ax}^0 = \mathbf{b}$;
- The initial solution support $\mathcal{S}^0 = \text{Support}\{\mathbf{x}^0\} = \emptyset$.

Main Iteration: Increment k by 1 and perform the following steps:

- **Sweep** – Compute the errors $\epsilon(j) = \min_{z_j} \|z_j \mathbf{a}_j - \mathbf{r}^{k-1}\|_2^2$ for all j using the optimal choice $z_j^* = \mathbf{a}_j^T \mathbf{r}^{k-1} / \|\mathbf{a}_j\|_2^2$;
- **Update Support** – Find a minimizer j_0 of $\epsilon(j)$: $\forall j \notin \mathcal{S}^{k-1}$, $\epsilon(j_0) \leq \epsilon(j)$, and update $\mathcal{S}^k = \mathcal{S}^{k-1} \cup \{j_0\}$;
- **Update Provisional Solution** – Compute \mathbf{x}^k , the minimizer of $\|\mathbf{Ax} - \mathbf{b}\|_2^2$ subject to $\text{Support}\{\mathbf{x}\} = \mathcal{S}^k$;
- **Update Residual** – Compute $\mathbf{r}^k = \mathbf{b} - \mathbf{Ax}^k$;
- **Stopping Rule** – If $\|\mathbf{r}^k\|_2 < \epsilon$, stop; else, apply another iteration.

Output: The proposed solution is \mathbf{x}^k obtained after k iterations.

This algorithm is known in the signal processing literature as *orthogonal matching pursuit* (OMP) [11, 12, 13, 14, 1], and it is the algorithm we have implemented, validated, and used throughout to find sparse solutions of (P_0^ϵ) .

The notion of *mutual coherence* gives us a simple criterion by which we can test when a solution of Equation (1) is the unique sparsest solution. In what follows, we assume that $\mathbf{A} \in \mathbb{R}^{n \times m}$, $n < m$, and $\text{rank}(\mathbf{A}) = n$.

Definition 1. The *mutual coherence* of a given matrix \mathbf{A} is the largest absolute normalized inner product between different columns from \mathbf{A} . Thus,

denoting the k -th column in \mathbf{A} by \mathbf{a}_k , the mutual coherence of \mathbf{A} is given by

$$\mu(\mathbf{A}) = \max_{1 \leq k, j \leq m, k \neq j} \frac{|\mathbf{a}_k^T \mathbf{a}_j|}{\|\mathbf{a}_k\|_2 \cdot \|\mathbf{a}_j\|_2}. \quad (3)$$

Theorem 1. *If \mathbf{x} solves the system of linear equations $\mathbf{A}\mathbf{x} = \mathbf{b}$ and $\|\mathbf{x}\|_0 < \frac{1}{2}(1 + 1/\mu(\mathbf{A}))$, then \mathbf{x} is the sparsest solution. Hence, if $\mathbf{y} \neq \mathbf{x}$ also solves the system, then $\|\mathbf{x}\|_0 < \|\mathbf{y}\|_0$.*

This same criterion can be used to test when OMP will find the sparsest solution.

Theorem 2. *For a system of linear equations $\mathbf{A}\mathbf{x} = \mathbf{b}$, if a solution \mathbf{x} exists obeying $\|\mathbf{x}\|_0 < \frac{1}{2}(1 + 1/\mu(\mathbf{A}))$, then an OMP run with threshold parameter $\epsilon = 0$ is guaranteed to find \mathbf{x} exactly.*

The proofs of these theorems can be found in [1].

3. Image representation and compression

3.1. Image representation concepts

For our purposes an image is a two dimensional sequence of sample values,

$$\mathcal{I}[n_1, n_2], \quad 0 \leq n_1 < N_1, \quad 0 \leq n_2 < N_2,$$

having finite extents, N_1 and N_2 , in the vertical and horizontal directions, respectively. The term ‘‘pixel’’ is synonymous with an image sample value. The first coordinate, n_1 , is the row index and the second coordinate, n_2 , is the column index of the pixel. The ordering of the pixels follows the canonical ordering of a matrix’s rows and columns, e.g., [3].

The sample value, $\mathcal{I}[n_1, n_2]$, represents the intensity (brightness) of the image at location $[n_1, n_2]$. The sample values will usually be B -bit signed or unsigned integers. Thus, we have

$$\begin{aligned} \mathcal{I}[n_1, n_2] &\in \{0, 1, \dots, 2^B - 1\} \quad \text{for unsigned imagery,} \\ \mathcal{I}[n_1, n_2] &\in \{-2^{B-1}, -2^{B-1} + 1, \dots, 2^{B-1} - 1\} \quad \text{for signed imagery.} \end{aligned}$$

In many cases, the B -bit sample values are interpreted as uniformly quantized representations of real-valued quantities, $\mathcal{I}'[n_1, n_2]$, in the range 0 to 1 (unsigned) or $-\frac{1}{2}$ to $\frac{1}{2}$ (signed). Letting $\text{round}(\cdot)$ denote rounding to the

nearest integer, the relationship between the real-valued and integer sample values may be written as

$$\mathcal{I}[n_1, n_2] = \text{round}(2^B \mathcal{I}'[n_1, n_2]). \quad (4)$$

This accounts for the sampling quantization error which is introduced by rounding the physically measured brightness at location $[n_1, n_2]$ on a light sensor to one of the allowed pixel values, e.g., [3, 7].

We shall use this framework to represent grayscale images, where a pixel value of 0 will represent “black” and a value of $2^B - 1$ “white”. The value of B is called the *depth* of the image, and typical values for B are 8, 10, 12, and 16. Color images are represented by either three values per sample, $\mathcal{I}_R[n_1, n_2]$, $\mathcal{I}_G[n_1, n_2]$, and $\mathcal{I}_B[n_1, n_2]$ each for the red, green, and blue channels, resp.; or by four values per pixel, $\mathcal{I}_C[n_1, n_2]$, $\mathcal{I}_M[n_1, n_2]$, $\mathcal{I}_Y[n_1, n_2]$, and $\mathcal{I}_K[n_1, n_2]$, each for the cyan, magenta, yellow, and black channels commonly used in applications for color printing, respectively. We shall restrict ourselves to grayscale images given that it is always possible to apply a compression system separately to each component in turn, e.g., [3].

3.2. Image compression via sparsity

3.2.1. Setup

In this section, we give an overview of image compression via sparsity. The basic idea is that if $\mathbf{Ax} = \mathbf{b}$, \mathbf{b} is dense—that is, it has mostly nonzero entries—, and \mathbf{x} is sparse, we can achieve compression by storing wisely \mathbf{x} instead of \mathbf{b} .

Specifically, suppose we have a signal $\mathbf{b} \in \mathbb{R}^n$ that requires n numbers for its description. However, suppose that we can solve problem (P_0^ϵ) , whose solution \mathbf{x}_0^ϵ has k nonzeros, with $k \ll n$, then we would have obtained an approximation $\hat{\mathbf{b}} = \mathbf{Ax}_0^\epsilon$ to \mathbf{b} using k scalars, with an approximation error of at most ϵ . Thus, by increasing ϵ we can obtain better compression at the expense of a larger approximation error. We shall characterize this relationship between error and compression, or equivalently, error and bits per sample, in Section 3.7.

3.2.2. From an image to a vector to an image

Following the approach to image processing at the core of the JPEG image compression standard [2], we subdivide each image in our database into 8×8 non-overlapping squares that will be treated individually. Since we need

to generate a right hand side vector \mathbf{b} to implement our compression scheme via sparsity, cf. Section 3.2.1, a sub-image $\mathbf{Y} \in \mathbb{R}^{8 \times 8}$ of size 8×8 pixels needs to be vectorized into a vector $\mathbf{y} \in \mathbb{R}^{64}$ to play the role of \mathbf{b} . There are many ways to do this, and we tested three possible approaches.

The first consists of concatenating one after the other the columns of \mathbf{Y} to form \mathbf{y} , and we shall call this method c_1 . It can be thought of as a bijection $c_1 : \mathbb{R}^{8 \times 8} \rightarrow \mathbb{R}^{64}$ that maps $\mathbf{Y} \mapsto \mathbf{y}$, see Figure 1(a).

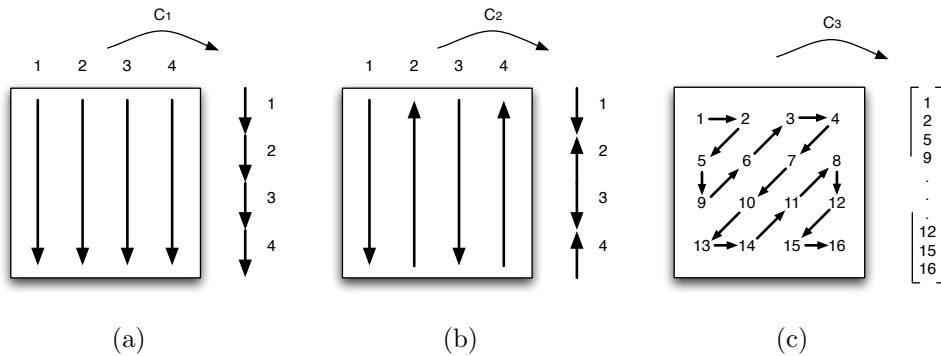


Figure 1: Three possible ways to vectorize a matrix. (a) Concatenate the columns of the matrix from top to bottom one after the other, from left to right; or (b) first flip every other column, and then concatenate the columns as before; or (c) traverse the elements of the matrix in a Cantor-Peano zigzag fashion as in this 4×4 example, but for the corresponding 8×8 matrix.

The second approach is to reverse the ordering of the entries of every even column of \mathbf{Y} and then concatenate the columns of the resulting matrix. We shall call this method c_2 . It is also a bijection $c_2 : \mathbb{R}^{8 \times 8} \rightarrow \mathbb{R}^{64}$ mapping $\mathbf{Y} \mapsto \mathbf{y}'$, see Figure 1(b).

Finally, a third method, $c_3 : \mathbb{R}^{8 \times 8} \rightarrow \mathbb{R}^{64}$, traverses an 8×8 sub-image \mathbf{Y} in a zigzag pattern from left to right, see Figure 1(c).

We must still designate a matrix $\mathbf{A} \in \mathbb{R}^{64 \times 128}$ to complete the setup. We shall address this issue in Section 3.3. For now, assume \mathbf{A} is given.

Once we have chosen c_i and \mathbf{A} , we proceed in the following way. Given a tolerance $\epsilon > 0$, and an image \mathcal{I} that has been partitioned into 8×8 non-overlapping sub-images \mathbf{Y}_l , where $l = 1, \dots, M$ and M is the number of sub-images partitioning \mathcal{I} , we obtain the approximation to $\mathbf{y}_l = c_i(\mathbf{Y}_l)$ derived from the OMP algorithm, i.e., from the sparse vector $\mathbf{x}_l = \text{OMP}(\mathbf{A}, \mathbf{y}_l, \epsilon)$, we compute $\tilde{\mathbf{y}}_l = \mathbf{A}\mathbf{x}_l$. Using $\tilde{\mathbf{y}}_l$ we can reconstruct a sub-image by setting $\tilde{\mathbf{Y}}_l = c_i^{-1}(\tilde{\mathbf{y}}_l)$.

Finally, we rebuild and approximate the original image \mathcal{I} by pasting together, in the right order and position, the set $\{\tilde{\mathbf{Y}}_l\}$ of sub-images and form the approximate image reconstruction $\tilde{\mathcal{I}}$ of \mathcal{I} . This new image $\tilde{\mathcal{I}}$ is an approximation, and not necessarily the original image \mathcal{I} , because in the process we have introduced an error by setting the tolerance $\epsilon > 0$, and not $\epsilon = 0$. On the other hand, we do have $\|\tilde{\mathbf{y}} - \mathbf{y}\|_2 < \epsilon$. Since $\|\mathbf{x}_l\|_0 \leq \|\mathbf{y}_l\|_0$, and more likely $\|\mathbf{x}_l\|_0 \ll \|\mathbf{y}_l\|_0$, storing wisely the set $\{\mathbf{x}_l\}_{l=1}^M$ will provide a compressed representation of the image \mathcal{I} .

The means to create efficiently a compressed representation of the image \mathcal{I} by using $\{\mathbf{x}_l\}_{l=1}^M$, the map $c_i : \mathbb{R}^{8 \times 8} \rightarrow \mathbb{R}^{64}$, and the matrix \mathbf{A} , and analyzing the effects of the choice of the tolerance ϵ on such a representation will all be addressed in subsequent sections.

3.3. Choosing a matrix \mathbf{A}

3.3.1. Setup

The choice of matrix \mathbf{A} is clearly central to our approach to compression. Because of the JPEG and JPEG 2000 standards that use at their core the discrete cosine transform (DCT) and a wavelet transform [2, 3], resp., we shall incorporate both transforms in our choices of \mathbf{A} .

3.3.2. One-dimensional basis elements

Given that the signal that we are going to process comes in the form of a vector \mathbf{b} , an inherently one-dimensional object, a first approach is to consider the one-dimensional DCT waveforms, and any one-dimensional wavelet basis for $L^2[0, 1]$. For the choice of the wavelet basis we opt for the Haar wavelet and its scaling function, the identity on $[0, 1]$.

More specifically, we know that the one-dimensional DCT-II transform [15, 16],

$$X_k = \sum_{n=0}^{N-1} x_n \cos \left[\frac{\pi}{N} \left(n + \frac{1}{2} \right) k \right], \quad k = 0, \dots, N-1, \quad (5)$$

has at its core a sampling of the function $f_{k,N}(x) = \cos \left(\pi \left(x + \frac{1}{2N} \right) k \right)$ on the regularly spaced set $\mathcal{S}(N) = \left\{ s_i \in [0, 1] : s_i = \frac{i}{N}, i = 0, \dots, N-1 \right\}$ of points. We define the vector $\mathbf{w}_{k,N} = \sqrt{2}^{\text{sgn}(k)} \cdot (f_{k,N}(s_0), \dots, f_{k,N}(s_{N-1}))^T \in \mathbb{R}^N$, which we generically call a *DCT waveform of wave number k and length N* . We shall use DCT waveforms with $N = 64$ for the one-dimensional

compression approach. This is because we subdivide each image in our database into collections of 8×8 non-overlapping sub-images, which are then transformed into vectors $\mathbf{b} \in \mathbb{R}^{64}$, and subsequently compressed, as described above; see Section 3.2.2. This collection of DCT waveforms is a basis for \mathbb{R}^{64} , and we arrange its elements column-wise in matrix form as $\text{DCT}_1 = (\mathbf{w}_{0,64} \dots \mathbf{w}_{63,64}) \in \mathbb{R}^{64 \times 64}$. Note that all column vectors of DCT_1 have the same ℓ^2 norm.

The corresponding basis for \mathbb{R}^{64} based on the Haar wavelet is built in the following way. Consider the Haar wavelet [17, 18],

$$\psi(x) = \begin{cases} 1 & \text{if } 0 \leq x < 1/2, \\ -1 & \text{if } 1/2 \leq x < 1, \\ 0 & \text{otherwise,} \end{cases} \quad (6)$$

and its scaling function,

$$\phi(x) = \begin{cases} 1 & \text{if } 0 \leq x < 1, \\ 0 & \text{otherwise.} \end{cases} \quad (7)$$

For each $n \geq 0$, define the set,

$$H_n = \{x \mapsto \psi_{n,k}(x) = 2^{n/2} \psi(2^n x - k) : 0 \leq k < 2^n\},$$

and define $H_{-1} = \{x \mapsto \phi(x)\}$. Note that $\#H_n = 2^n$ for $n \geq 0$ and $\#H_{-1} = 1$; and therefore $\#\bigcup_{j=-1}^n H_j = 1 + \sum_{j=0}^n 2^j = 2^{n+1}$. Since $64 = 2^6$, a value of $n = 5$ will produce 64 functions to choose from in $\mathcal{H}(n) = \bigcup_{j=-1}^n H_j$. For each function $h \in \mathcal{H}(5)$, define the vector $\mathbf{v}_{h,64} \in \mathbb{R}^{64}$ by sampling h on the set $\mathcal{S}(64)$, i.e., $\mathbf{v}_{h,64} = (h(s_0), \dots, h(s_{63}))^T$. Observe that $\|\mathbf{v}_{h,64}\|_2 = \|\mathbf{w}_{0,64}\|_2$ for all $h \in \mathcal{H}(5)$.

We can order the elements of $\mathcal{H}(5)$ in a natural way, viz., $h_0 = \phi$, $h_1 = \psi_{0,0}$, $h_2 = \psi_{1,0}$, $h_3 = \psi_{1,1}$, etc. It is easy to see that the set $\{\mathbf{v}_{h_j,64}\}_{j=0}^{63}$ of vectors is a basis for \mathbb{R}^{64} . As with the DCT waveforms, we arrange these vectors column-wise in matrix form as $\text{Haar}_1 = (\mathbf{v}_{h_0,64} \dots \mathbf{v}_{h_{63},64}) \in \mathbb{R}^{64 \times 64}$.

Then, for the one-dimensional approach, we define $\mathbf{A} = [\text{DCT}_1 \text{ Haar}_1] \in \mathbb{R}^{64 \times 128}$, the concatenation of both bases.

3.3.3. Two-dimensional basis elements

Another way to choose a basis for \mathbb{R}^{64} results from taking into account the intrinsic two-dimensional nature of an image and to define basis elements that reflect this fact. Specifically, consider the two-dimensional DCT-II,

$$X_{k_1, k_2} = \sum_{n_1=0}^{N_1-1} \sum_{n_2=0}^{N_2-1} x_{n_1, n_2} \cos \left[\frac{\pi}{N_1} \left(n_1 + \frac{1}{2} \right) k_1 \right] \cos \left[\frac{\pi}{N_2} \left(n_2 + \frac{1}{2} \right) k_2 \right], \quad (8)$$

used in the JPEG standard [2] and, as before, consider the family of functions indexed by k_1 and k_2 ,

$$g_{k_1, k_2, N_1, N_2}(x, y) = \cos \left[\pi \left(x + \frac{1}{2N_1} \right) k_1 \right] \cos \left[\pi \left(y + \frac{1}{2N_2} \right) k_2 \right],$$

sampled on all points $(x, y) \in \mathcal{S}(N_1) \times \mathcal{S}(N_2)$, where in our case, $N_1 = N_2 = 8$, and $k_1, k_2 \in \{0, \dots, 7\}$.

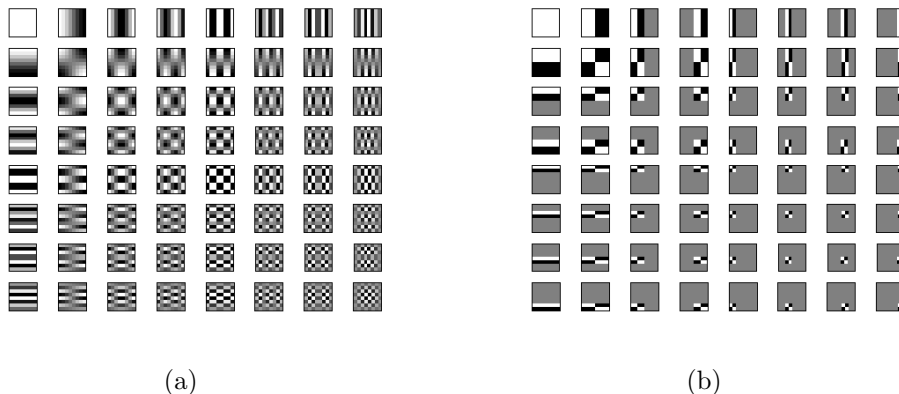


Figure 2: Natural representation for the 2-dimensional DCT (a) and Haar (b) bases for \mathbb{R}^{64} . White corresponds to the maximum value achieved by the basis element, black to the minimum. The intermediate shade of gray corresponds to 0.

Note that $g_{k_1, k_2, N_1, N_2}(x, y) = f_{k_1, N_1}(x) f_{k_2, N_2}(y)$, and therefore the image of $\mathcal{S}(8) \times \mathcal{S}(8)$ under $g_{k_1, k_2, 8, 8}$ can be naturally identified with the tensor product

$$\mathbf{w}_{k_1, 8} \otimes \mathbf{w}_{k_2, 8}, \quad k_1, k_2 = 0, \dots, 7, \quad (9)$$

modulo the constants $\sqrt{2}^{\text{sgn}(k_1)}$ and $\sqrt{2}^{\text{sgn}(k_2)}$, that make $\|\mathbf{w}_{k_1, 8}\|_2 = \|\mathbf{w}_{k_2, 8}\|_2$,

respectively. See Figure 2(a) for a graphical representation of all of these tensor products.

There is a total of 64 such tensor products, and for each of them we can obtain a vector $\tilde{\mathbf{w}}_{k_1, k_2, j} = c_j^{-1}(\mathbf{w}_{k_1, 8} \otimes \mathbf{w}_{k_2, 8})$; see Section 3.2.2. It is easy to see that the vector columns of the matrix,

$$\text{DCT}_{2,j} = (\tilde{\mathbf{w}}_{k_1, k_2, j}) \in \mathbb{R}^{64 \times 64}, \quad k_1, k_2 = 0, \dots, 7, \quad (10)$$

form a basis for \mathbb{R}^{64} . The ordering of the column vectors of $\text{DCT}_{2,j}$ follows the lexicographical order of the sequence of ordered pairs (k_1, k_2) for $k_1, k_2 = 0, \dots, 7$ when k_2 moves faster than k_1 , i.e., the ordering is $(0,0), (0,1), \dots, (0,7), (1, 0), \dots, (7,7)$.

In a similar fashion, we can build $\text{Haar}_{2,j}$. Specifically, first consider the set $\mathcal{H}(2)$, which contains 8 functions. Sampling $h_k \in \mathcal{H}(2)$ on $\mathcal{S}(8)$, we obtain the vector $\mathbf{v}_{h_k, 8} \in \mathbb{R}^8$. Given $k_1, k_2 \in \{0, \dots, 7\}$ we then compute the vector,

$$\tilde{\mathbf{v}}_{k_1, k_2, j} = c_j^{-1} \left(\mathbf{v}_{h_{k_1}, 8} \otimes \mathbf{v}_{h_{k_2}, 8} \right) \in \mathbb{R}^{64},$$

with which, for all $k_1, k_2 \in \{0, \dots, 7\}$ and following the same order for the ordered pairs (k_1, k_2) mentioned above, we can construct

$$\text{Haar}_{2,j} = (\tilde{\mathbf{v}}_{k_1, k_2, j}) \in \mathbb{R}^{64 \times 64}, \quad k_1, k_2 = 0, \dots, 7. \quad (11)$$

For a visual representation of the tensor products $\mathbf{v}_{h_{k_1}, 8} \otimes \mathbf{v}_{h_{k_2}, 8}$, see Figure 2(b).

Finally, we define $\mathbf{A} = \mathbf{A}_j = [\text{DCT}_{2,j} \text{ Haar}_{2,j}] \in \mathbf{R}^{64 \times 128}$, the concatenation of both bases.

3.4. Imagery metrics

3.4.1. Compression and bit-rate

Following the methodology described in Section 3.2.2, suppose that we have an image \mathcal{I} of size $N_1 \times N_2$ pixels ($N_1 = N_2 = 512$ in our case), and let $\{\mathbf{Y}_l\}_{l \in L}$ be a partition of \mathcal{I} into $\#L$ 8×8 sub-images. Let $\mathbf{y}_l = c_i(\mathbf{Y}_l)$ for some $i = 1, 2$, or 3 , where $c_i : \mathbb{R}^{8 \times 8} \rightarrow \mathbb{R}^{64}$, see Figure 1. Using OMP, with a full-rank matrix $\mathbf{A} \in \mathbb{R}^{64 \times 128}$ and a tolerance $\epsilon > 0$, obtain $\mathbf{x}_l = \text{OMP}(\mathbf{A}, \mathbf{y}_l, \epsilon)$, and compute the approximation to \mathbf{y}_l given by $\tilde{\mathbf{y}}_l = \mathbf{A}\mathbf{x}_l$.

We know that $\|\tilde{\mathbf{y}}_l - \mathbf{y}_l\|_2 < \epsilon$, and we can count how many nonzero entries there are in \mathbf{x}_l , viz., $\|\mathbf{x}_l\|_0$. With this, we can define the *normalized sparse*

bit-rate.

Definition 2 (Normalized Sparse Bit-Rate). Given a matrix \mathbf{A} and tolerance ϵ , the *normalized sparse bit-rate*, measured in bits per pixel (bpp), for the image \mathcal{I} is the number

$$nsbr(\mathcal{I}, \mathbf{A}, \epsilon) = \frac{\sum_l \|\mathbf{x}_l\|_0}{N_1 N_2}, \quad (12)$$

where the image \mathcal{I} is of size $N_1 \times N_2$ pixels ($N_1 = N_2 = 512$ in our case).

To interpret this definition, suppose a binary digit or “bit” represents a nonzero coordinate of the vector \mathbf{x}_l . Then we need at least $\|\mathbf{x}_l\|_0$ bits to store or transmit \mathbf{x}_l ; and so the total number of bits needed to represent image \mathcal{I} is at least $\sum_l \|\mathbf{x}_l\|_0$. Thus, the average number of bits per pixel (bpp) is obtained by dividing this quantity by $N_1 N_2$. This is the normalized sparse bit-rate.

Suppose \mathcal{I} is of depth B , and that \mathcal{I} can be represented by a string of bits (*bit-stream*) \mathbf{c} . An important objective of image compression is to assure that the length of \mathbf{c} , written $length(\mathbf{c})$, is as small as possible. In the absence of any compression, we require $N_1 N_2 B$ bits to represent the image sample values, e.g. [3]. Thus, the notion of *compression ratio* is defined as

$$cr(\mathcal{I}, \mathbf{c}) = \frac{N_1 N_2 B}{length(\mathbf{c})}, \quad (13)$$

and the *compressed bit-rate*, expressed in bpp, is defined as

$$br(\mathcal{I}, \mathbf{c}) = \frac{length(\mathbf{c})}{N_1 N_2}. \quad (14)$$

The compression ratio is a dimensionless quantity that tells how many times we have managed to reduce in size the original representation of the image, while the compressed bit-rate has bpp units, and it tells how many bits are used on average per sample by the compressed bit-stream \mathbf{c} to represent the original image \mathcal{I} .

We note from Equations (12) and (14) that it is likely that $nsbr(\mathcal{I}, \mathbf{A}, \epsilon) \leq br(\mathcal{I}, \mathbf{c})$ if the bit-stream \mathbf{c} is derived from the sparse representation induced by \mathbf{A} and ϵ . The rationale for this assertion is two-fold. First, it is unlikely that the coordinates in each of the resulting \mathbf{x}_l vectors will be realistically

represented by only one bit, and, second, because \mathbf{c} would somehow have to include a coding for the indices l for each \mathbf{x}_l , necessarily increasing the bit count some more. These particular issues, i.e., the number of bits to represent the entries of vectors \mathbf{x}_l , and the coding of the indices l for each of those vectors into a final compressed bit-stream \mathbf{c} , will be addressed in Section 4.

3.4.2. Error estimation criteria

Let $\mathbf{A} = [\text{DCT}_1 \text{ Haar}_1]$ and consider the vectorization function c_2 .

Given a 512×512 image \mathcal{I} in our database, we can proceed to compress it using the methodology described in Section 3.2.2. If \mathcal{I} is partitioned into 8×8 non-overlapping sub-images $\mathbf{Y}_l, l = 1, \dots, 4096$, we obtain for each of them a corresponding reconstructed sub image $\tilde{\mathbf{Y}}_l = c_2^{-1}(\tilde{\mathbf{y}}_l)$, and from those we reconstruct an approximation $\tilde{\mathcal{I}}$ to \mathcal{I} . Here, $\tilde{\mathbf{y}}_l = \mathbf{A}\mathbf{x}_l$, $\mathbf{x}_l = \text{OMP}(\mathbf{A}, \mathbf{y}_l, \epsilon)$, and $\mathbf{y}_l = c_2(\mathbf{Y}_l)$, as before. The compression would come from storing wisely $\{\mathbf{x}_l\}$. We summarize this procedure with the notation, $\tilde{\mathcal{I}} = \text{rec}(\mathcal{I}, \mathbf{A}, \epsilon)$.

In order to assess the quality of $\tilde{\mathcal{I}}$ when compared to \mathcal{I} , we introduce three error estimators.

Traditionally, the signal processing community has relied on the *peak signal-to-noise ratio*, or *PSNR* [3, 8].

Definition 3 (PSNR). The *peak signal-to-noise ratio* between two images $\tilde{\mathcal{I}}$ and \mathcal{I} is the quantity, measured in dB,

$$\text{PSNR}(\tilde{\mathcal{I}}, \mathcal{I}) = 20 \log_{10} \left(\frac{\max_{\mathcal{I}}}{\sqrt{\text{mse}(\tilde{\mathcal{I}}, \mathcal{I})}} \right),$$

where $\max_{\mathcal{I}}$ is the maximum possible value for any given pixel in \mathcal{I} , $\max_{\mathcal{I}} = 2^B - 1$ typically; and $\text{mse}(\tilde{\mathcal{I}}, \mathcal{I}) = \frac{1}{N_1 N_2} \sum_{i,j} (\tilde{\mathcal{I}}[i, j] - \mathcal{I}[i, j])^2$ is the mean square error between both images. Here N_1 and N_2 represent the dimensions of \mathcal{I} , and $\mathcal{I}[i, j]$ represents the value of the pixel at coordinates $[i, j]$ in image \mathcal{I} —similarly for $\tilde{\mathcal{I}}[i, j]$. In our case $N_1 = N_2 = 512$, and $\max_{\mathcal{I}} = 255$.

PSNR has the advantage that it is easy to compute and has widespread use, but it has been criticized for poorly correlating with perceived image quality [9, 19]. However, in recent years extensive work on other error estimators that take into account the human visual system have arisen. In

particular, we present and define the *structural similarity* and *mean structural similarity* indices [9].

Definition 4 (SSIM). Let $\tilde{\mathcal{I}}$ and \mathcal{I} be two images that have been decomposed in $L \times L$ non-overlapping sub images $\{\tilde{\mathbf{Y}}_l\}$ and $\{\mathbf{Y}_l\}$, respectively. Then the *structural similarity* index for two corresponding sub-image vectorizations, say $\tilde{\mathbf{y}}_l = c_2(\tilde{\mathbf{Y}}_l)$ and $\mathbf{y}_l = c_2(\mathbf{Y}_l)$, is defined as follows

$$\text{SSIM}(\tilde{\mathbf{y}}_l, \mathbf{y}_l) = \frac{(2\mu_{\tilde{\mathbf{y}}_l}\mu_{\mathbf{y}_l} + C_1)(2\sigma_{\tilde{\mathbf{y}}_l\mathbf{y}_l} + C_2)}{(\mu_{\tilde{\mathbf{y}}_l}^2 + \mu_{\mathbf{y}_l}^2 + C_1)(\sigma_{\tilde{\mathbf{y}}_l}^2 + \sigma_{\mathbf{y}_l}^2 + C_2)},$$

where $\mu_{\mathbf{y}_l}$ and $\sigma_{\mathbf{y}_l}$ represent the mean and standard deviation of \mathbf{y}_l , resp.; and similarly for $\tilde{\mathbf{y}}_l$. The term $\sigma_{\tilde{\mathbf{y}}_l\mathbf{y}_l}$ is the correlation between $\tilde{\mathbf{y}}_l$ and \mathbf{y}_l . The values C_1 and C_2 are two small constants.

For our purposes, we used the default values of $L = 11$, $C_1 = 0.01$, and $C_2 = 0.03$ used in [9] when assessing the SSIM of an image in our database and its reconstruction. We used a value of $L = 4$ when we modified OMP to use internally the SSIM as a stopping criteria. More on this later.

From the above definition, we can see that the SSIM index is a localized quality measure that can be represented on a plane that maps its values. It can take values from 0 to 1 and when it takes the value of 1 the two images are identical. In practice, we usually require a single overall quality of measure for the entire image. In that case we use the mean SSIM index to evaluate the overall image quality.

Definition 5 (MSSIM). Let $\tilde{\mathcal{I}}$ and \mathcal{I} be two images, where the former is the approximation and the later is the original. Then the *mean structural similarity index* is

$$\text{MSSIM}(\tilde{\mathcal{I}}, \mathcal{I}) = \frac{1}{M} \sum_{l=1}^M \text{SSIM}(\tilde{\mathbf{y}}_l, \mathbf{y}_l),$$

where $\tilde{\mathbf{y}}_l$ and \mathbf{y}_l are the image contents at the l -th local sub-image, and M is the number of local sub-images in the image.

Finally, we take a look at the relationship between the size of the sub-image and the tolerance, and how this affects the quality of the approximation. We analyze the idealized error distribution in which all pixels of

the approximation are c units apart from the original. Consider an $L \times L$ sub image that has been linearized to a vector \mathbf{y} of length L^2 . Assume that the OMP approximation within ϵ has distributed the error evenly, that is, if $\mathbf{x} = \text{OMP}(\mathbf{A}, \mathbf{y}, \epsilon)$ and $\tilde{\mathbf{y}} = \mathbf{A}\mathbf{x}$, then

$$\begin{aligned}
\|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2 < \epsilon &\Leftrightarrow \|\tilde{\mathbf{y}} - \mathbf{y}\|_2^2 < \epsilon^2, \\
&\Leftrightarrow \sum_{j=1}^{L^2} (\tilde{\mathbf{y}}(j) - \mathbf{y}(j))^2 < \epsilon^2, \\
&\Leftrightarrow L^2 c^2 < \epsilon^2, \\
&\Leftrightarrow c < \frac{\epsilon}{L}.
\end{aligned} \tag{15}$$

That is, if we want to be within c units from each pixel, we have to choose a tolerance ϵ such that $c = \epsilon/L$.

We note that the least-squares approximation at the core of OMP approximates the idealized error distribution. This can be seen in Figure 3 where the black dashed line represents this idealized error approximation. For tolerances $\epsilon > 40$, we can see that the PSNR for all images considered is above this idealized error distribution. This can be explained by noting that, for example, for $\epsilon = 2048$, we would have from Equation (15) that $c = 2048/8 = 256$, but the maximum pixel value is only 255. Therefore, unless the original image \mathcal{I} is just a white patch, the initial value of the OMP approximation being an all black image, there are matching pixels in the original and the approximation image $\tilde{\mathcal{I}} = \text{rec}(\mathcal{I}, \mathbf{A}, 2048)$ that are less than 256 units apart. This would necessarily, by Definition 3, imply $\text{PSNR}(\tilde{\mathcal{I}}, \mathcal{I}) > 0$, a value above the value of the PSNR for the idealized error distribution when $\epsilon = 2048$, which is a small negative value.

On the other hand, for small tolerances, about $\epsilon < 3$, we observe that the PSNR value for all images jumps again above the PSNR for the idealized error model. This is a happy case when roundoff error actually helps. What happens is that for such small tolerances, the roundoff to the closest integer for all entries in $\tilde{\mathbf{y}}_l = \mathbf{A}\mathbf{x}_l$ when we form the sub image approximation $\tilde{\mathbf{Y}}_l = c_2^{-1}(\tilde{\mathbf{y}}_l)$, coincides with the true value of the pixels in the original sub image \mathbf{Y}_l . Again, by Definition 3, this increases the value of $\text{PSNR}(\tilde{\mathcal{I}}, \mathcal{I})$ compared to the case where roundoff would not have taken place.

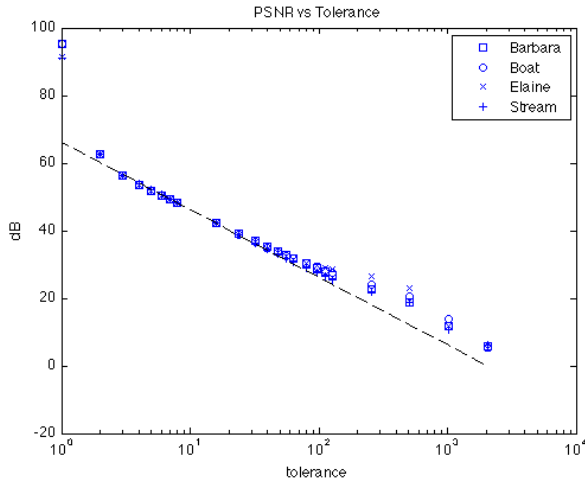


Figure 3: Peak Signal-to-Noise Ratio vs tolerance. We observe three typical behaviors for all images. For large values of the tolerance, about $\epsilon > 40$, the PSNR of all images is above the PSNR value for the idealized error distribution marked by the black dashed line. This behavior is also observed for very small values of the tolerance, about $\epsilon < 3$. For values between these two extreme behaviors, all images conform very closely to the idealized error distribution, a fact that is expected from the least-squares approximation at the core of the OMP algorithm.

3.5. Image database

To carry out our experiments and test image compression via sparsity, as well as the properties of the matrices described in Section 3.3 for that purpose, we selected 4 natural images, 3 of them from the University of Southern California’s Signal & Image Processing Institute (USC-SIPI) image database [6]. This database has been widely used for image processing benchmarking. The images are described in Table 1 and shown in Figure 4. All images are 512×512 , 8-bit grayscale images, which means they are composed of $512^2 = 262,144$ pixels that can take integer values from 0 (black) to 255 (white).

3.6. Effects of vectorization on image reconstruction

Given an image \mathcal{I} in our database, we can follow and apply to it the methodology described in Section 3.2.2, and obtain at the end of this process a reconstructed image $\tilde{\mathcal{I}}$ from it. In this section we explore the effects of the choice of map $c_i : \mathbb{R}^{8 \times 8} \rightarrow \mathbb{R}^{64}$ on the characteristics of image $\tilde{\mathcal{I}}$ for the different choices of matrix \mathbf{A} that we have selected to study.

USC-SIPI file	Name	Size
n/a	Barbara	512 × 512 pixels, 8-bit
boat.512	Boat	512 × 512 pixels, 8-bit
elaine.512	Elaine	512 × 512 pixels, 8-bit
5.2.10	Stream	512 × 512 pixels, 8-bit

Table 1: List of images used. Image *Barbara* is not in the USC-SIPI database but we included it since it is widely used by the image processing community.



(a) Barbara



(b) Boat



(c) Elaine



(d) Stream

Figure 4: Images used for our compression algorithms based on sparse image representation.

We summarize the empirical observations obtained from the experiments described below, and suggest a possible explanation for them without further proof.

3.6.1. Results for $\mathbf{A} = [DCT_1 \text{ Haar}_1]$

We set $\mathbf{A} = [DCT_1 \text{ Haar}_1]$, and choose a tolerance $\epsilon = 32$. Then, for each image \mathcal{I} in our database and each index $i = 1, 2, 3$ we choose the map $c_i : \mathbb{R}^{8 \times 8} \rightarrow \mathbb{R}^{64}$, and follow the methodology described in Section 3.2.2.

Image/Function	PSNR (dB)	Normalized sparse bit-rate (bpp)
Barbara		
c_1 :	36.8996	0.1833
c_2 :	36.9952	0.1863
c_3 :	36.8470	0.2338
Boat		
c_1 :	36.5791	0.1812
c_2 :	36.6020	0.1608
c_3 :	36.5615	0.2205
Elaine		
c_1 :	36.5003	0.1763
c_2 :	36.5155	0.1682
c_3 :	36.4877	0.1885
Stream		
c_1 :	36.4423	0.3161
c_2 :	36.4686	0.3050
c_3 :	36.4400	0.3504

Table 2: Performance results for $\mathbf{A} = [DCT_1 \text{ Haar}_1]$ for c_1 , c_2 , and c_3 . For each image in our test database, we vectorized each 8×8 sub-image using c_1 , c_2 , or c_3 . In all cases, the PSNR value was larger using c_2 ; and in all cases, except for image *Barbara*, the normalized sparse bit-rate was smaller. Both of these measures make c_2 a better choice than c_1 or c_3 . The results correspond to runs of OMP with an ℓ^2 stopping rule, and a tolerance $\epsilon = 32$.

The metrics that we use to measure the impact of the choice of map c_i are the *normalized sparse bit-rate*, and the *peak signal-to-noise ratio* (PSNR), see Definitions 2 and 3, respectively. A smaller value for the normalized sparse bit-rate is better than a larger one given that this implies fewer bits are necessary to represent the image. A larger value for the PSNR is better than

a smaller one as this means the fidelity of the representation is higher. Table 2 summarizes the results of the experiments. From these, we conclude that for $\mathbf{A} = [\text{DCT}_1 \text{ Haar}_1]$, the choice of c_2 over c_1 or c_3 produces better results. This could be because, if $\mathbf{Y} = (Y_{i,j})_{i,j=1,\dots,8}$ is a natural image, then on average $|Y_{8,j} - Y_{8,j+1}| < |Y_{8,j} - Y_{1,j+1}|$ for $j = 1, \dots, 7$, which makes $\mathbf{y}_{c_2} = c_2(\mathbf{Y})$ change more slowly than $\mathbf{y}_{c_1} = c_1(\mathbf{Y})$ or $\mathbf{y}_{c_3} = c_3(\mathbf{Y})$. By analogy to the behavior of the DFT, this must translate into fewer column vectors from \mathbf{A} to describe, within a certain error ϵ , the signal \mathbf{y}_{c_2} compared to the number of columns needed to approximate the signals \mathbf{y}_{c_1} or \mathbf{y}_{c_3} to the same error tolerance.

3.6.2. Results for $\mathbf{A}_j = [\text{DCT}_{2,j} \text{ Haar}_{2,j}]$

Proceeding similarly as in Section 3.6.1, we set $\mathbf{A}_j = [\text{DCT}_{2,j} \text{ Haar}_{2,j}]$ for $j = 1, 2$, and 3, and we perform the following experiment.

Define the vectorization function to match the same ordering that was used to form \mathbf{A}_j . This means we pick the vectorization function to be c_j . We compute $\mathbf{y}_l = c_j(\mathbf{Y}_l)$, where, as before, the sub-image \mathbf{Y}_l comes from the partition $\{\mathbf{Y}_l\}_{l \in L}$ of an image \mathcal{I} in our image database. Then, continuing with the compression methodology described in Section 3.2.2, we obtain $\mathbf{x}_l = \text{OMP}(\mathbf{A}, \mathbf{y}_l, \epsilon)$, setting $\epsilon = 32$ for this experiment. Finally, from the set of vectors $\{\mathbf{x}_l\}_{l \in L}$ we obtain $\tilde{\mathbf{y}}_l = \mathbf{A}\mathbf{x}_l$, and use $\{\tilde{\mathbf{y}}_l\}_{l \in L}$ to obtain the reconstructed image $\tilde{\mathcal{I}}$ of our original image \mathcal{I} . Again, as in Section 3.6.1, we assess the effects of the choice of the vectorization function c_i by the values of PSNR and normalized sparse bit-rate resulting from this representation of \mathcal{I} by $\tilde{\mathcal{I}}$. We give a summary of the results of this experiment in Table 3.

We point out that choosing c_i , with $i \neq j$, when $\mathbf{A}_j = [\text{DCT}_{2,j} \text{ Haar}_{2,j}]$, results in worse values of both PSNR and normalized sparse bit-rate than when $i = j$. We record only the results where $i = j$.

Moreover, any choice of $\mathbf{A}_j = [\text{DCT}_{2,j} \text{ Haar}_{2,j}]$ with a matching vectorization function c_j performs better than when $\mathbf{A} = [\text{DCT}_1 \text{ Haar}_1]$ for the normalized sparse bit-rate metric, and better for the PSNR metric except for the image *Stream*. Also, on average, the vectorization order imposed by c_3 is slightly better than those by either c_1 or c_2 , although the difference is practically imperceptible to the human eye. The normalized sparse bit-rate figures all coincide.

Image/Function	PSNR (dB)	Normalized sparse bit-rate (bpp)	Matrix
Barbara			
c_1 :	37.0442	0.1634	[DCT _{2,1} Haar _{2,1}]
c_2 :	37.0443	0.1634	[DCT _{2,2} Haar _{2,2}]
c_3 :	37.0443	0.1634	[DCT _{2,3} Haar _{2,3}]
Boat			
c_1 :	36.6122	0.1541	[DCT _{2,1} Haar _{2,1}]
c_2 :	36.6120	0.1541	[DCT _{2,2} Haar _{2,2}]
c_3 :	36.6120	0.1541	[DCT _{2,3} Haar _{2,3}]
Elaine			
c_1 :	36.5219	0.1609	[DCT _{2,1} Haar _{2,1}]
c_2 :	36.5219	0.1609	[DCT _{2,2} Haar _{2,2}]
c_3 :	36.5220	0.1609	[DCT _{2,3} Haar _{2,3}]
Stream			
c_1 :	36.4678	0.2957	[DCT _{2,1} Haar _{2,1}]
c_2 :	36.4676	0.2957	[DCT _{2,2} Haar _{2,2}]
c_3 :	36.4677	0.2957	[DCT _{2,3} Haar _{2,3}]

Table 3: Performance results for $\mathbf{A}_j = [\text{DCT}_{2,j} \text{ Haar}_{2,j}]$, $j = 1, 2, 3$, with corresponding vectorization functions c_1 , c_2 , and c_3 . In all cases, the PSNR and normalized sparse bit-rate values were almost identical. Matrix \mathbf{A}_3 performs slightly better on average. Mismatching function c_i with matrix $\mathbf{A}_j = [\text{DCT}_{2,j} \text{ Haar}_{2,j}]$, when $i \neq j$, results in degraded performance. The values correspond to runs of OMP with an ℓ^2 stopping rule, and a tolerance $\epsilon = 32$.

3.7. Normalized sparse bit-rate vs tolerance

Notwithstanding the remarks in Section 3.4.1, there is still value in using the normalized bit-stream measure to quantify and plot normalized sparse bit-rate vs tolerance graphs to gauge the compression properties of various compression matrices.

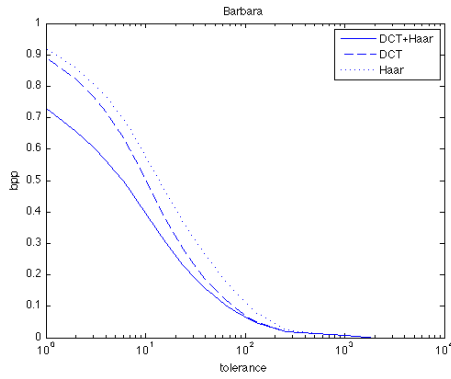
Given the results in Table 3, we look into the compression properties of matrices $\mathbf{A} = [\text{DCT}_1 \text{ Haar}_1]$, and $\mathbf{A}_3 = [\text{DCT}_{2,3} \text{ Haar}_{2,3}]$. We compare these properties for both matrices relative to each other, and to the compression properties of \mathbf{B} and \mathbf{C} , which are formed from the DCT_1 or the Haar_1 submatrices of matrix \mathbf{A} , respectively. We plot for all images in our database their respective normalized sparse bit-rate vs tolerance graphs. We took $\epsilon \in T = \{2^k\}_{k=0}^{11} \cup \{3, 5, 6, 7, 24, 40, 48, 56, 80, 96, 112\}$ and for each image \mathcal{I} in our image database we obtained the corresponding normalized sparse bit-rates $nsbr(\mathcal{I}, \mathbf{A}, \epsilon)$, $nsbr(\mathcal{I}, \mathbf{B}, \epsilon)$, and $nsbr(\mathcal{I}, \mathbf{C}, \epsilon)$ to obtain the plots in Figure 5. In Figure 6 we compare $\mathbf{A} = [\text{DCT}_1 \text{ Haar}_1]$ with $\mathbf{A}_3 = [\text{DCT}_{2,3} \text{ Haar}_{2,3}]$. We observe that up to a tolerance $\epsilon_{\mathcal{I}}$, dependent on image \mathcal{I} , \mathbf{A}_3 performs better for tolerance values $\epsilon \geq \epsilon_{\mathcal{I}}$. That is, the value of the normalized sparse bit-rate is smaller when performing compression utilizing \mathbf{A}_3 . For values of $\epsilon \leq \epsilon_{\mathcal{I}}$, compression with \mathbf{A} results in better normalized sparse bit-rate values. We shall see in Section 3.8, with the aid of Figure 3, that for values of $\epsilon = 32$, and smaller, the quality of the image reconstruction is satisfactory. We note from Figure 6 that, for all images in our database, $\epsilon_{\mathcal{I}} < 32$. This means that, for most practical cases, the use of $[\text{DCT}_{2,3} \text{ Haar}_{2,3}]$ results in slightly smaller normalized sparse bit-rate values than when using $[\text{DCT}_1 \text{ Haar}_1]$.

From the results shown in Figure 5, we can see that the DCT_1 basis elements perform better compression for any given tolerance than when using the Haar_1 basis elements, except for image *Stream*. In fact, when the tolerance ϵ is close to but less than 3, the Haar_1 basis elements result in a smaller normalized sparse bit-rate value. Moreover, and more importantly, combining both the DCT_1 and Haar_1 bases results in better compression than if either basis is used alone. The same is true for $\text{DCT}_{2,3}$ and $\text{Haar}_{2,3}$.

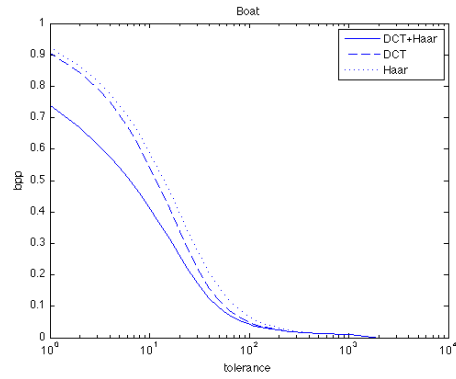
In this light, there are natural questions dealing with and relating image reconstruction quality, range of effective tolerances, and error estimators.

3.8. PSNR and MSSIM comparison

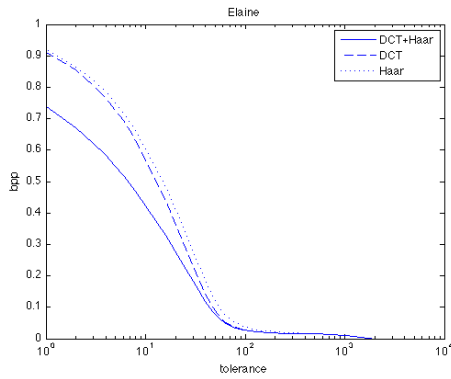
In Figure 7 we have plotted the normalized sparse bit-rate versus both error indices MSSIM and PSNR. The first thing that we observe is that the



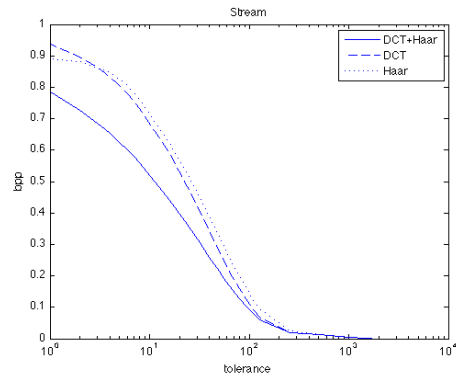
(a) Barbara



(b) Boat



(c) Elaine



(d) Stream

Figure 5: Normalized sparse bit-rate vs tolerance: One-dimensional basis elements. We observe that for all images the best normalized sparse bit-rate for a given tolerance is obtained for matrix $\mathbf{A} = [\text{DCT}_1 \text{ Haar}_1]$ which combines both the DCT_1 and Haar_1 bases for \mathbb{R}^{64} .

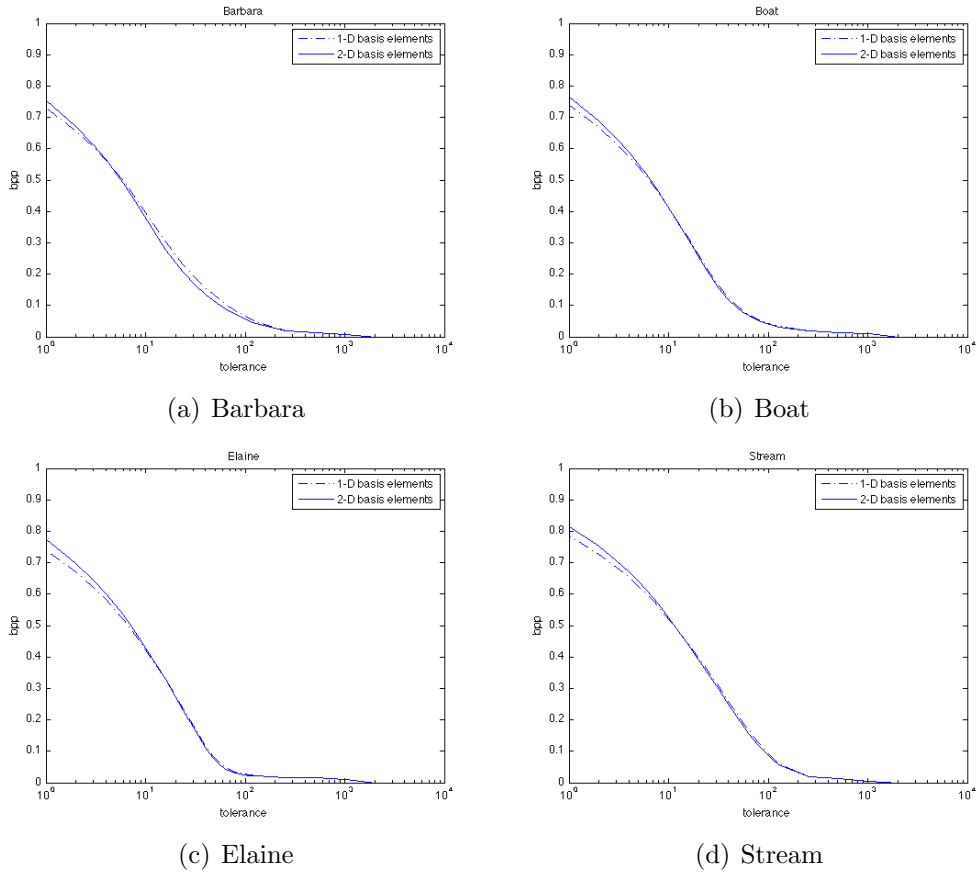


Figure 6: Normalized sparse bit-rate vs tolerance. We compare the performance of $\mathbf{A} = [\text{DCT}_1 \text{ Haar}_1]$ and $\mathbf{A}_3 = [\text{DCT}_{2,3} \text{ Haar}_{2,3}]$.

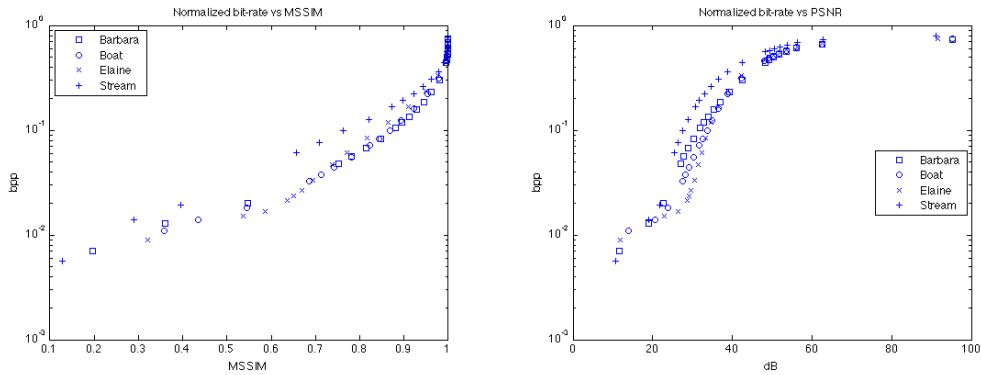


Figure 7: Normalized sparse bit-rate vs MSSIM, PSNR

sensitivity for PSNR varies more than the sensitivity for MSSIM over the range of tolerances chosen.

From Figure 8 we observe that, for the range of 20 to 40 dB in PSNR, the MSSIM index ranges from about 0.33 to 0.98. Since a value of 1 in MSSIM corresponds to two identical images, we can focus on values of PSNR no greater than 40 dB in our analysis. Also, in Figure 8, we corroborate the criticism that has been addressed to PSNR as a measure of image quality. For example, at 20 dB, the image *Stream* has an MSSIM value of 0.33, whereas the image *Elaine* has an MSSIM value of 0.48. Similarly, at 30 dB, the image *Elaine* has an MSSIM value of 0.69, whereas the image *Stream* has an MSSIM value of 0.86. It is not until 40 dB that we have a much smaller range of MSSIM values, viz., 0.96 (*Elaine*) to 0.98 (*Stream*). Therefore, if SSIM and MSSIM capture more accurately the human visual system’s perception of image quality, then the PSNR index is not a reliable estimator until values larger than or equal to 35 dB.

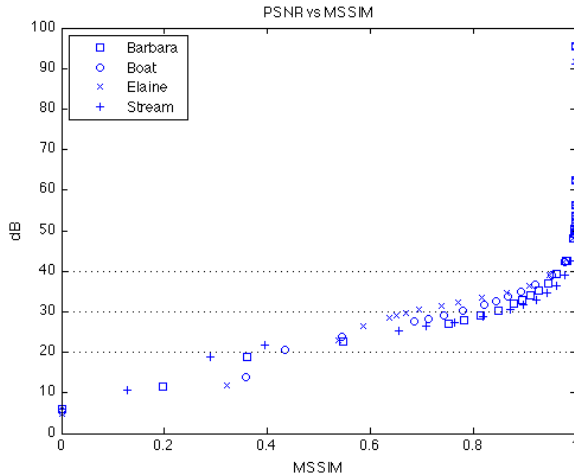


Figure 8: Peak Signal-to-Noise Ratio vs Mean Structural Similarity

Because of this observation about the PSNR index, we shall focus on the SSIM and MSSIM indices. We address the questions at the end of Section 3.7 and answer them with Figure 9. From this figure, if we were to consider desirable values of MSSIM to be greater than or equal to 0.9, we would see that this corresponds to a tolerance $\epsilon \leq 32$ for the image *Elaine* and $\epsilon \leq 48$ for the image *Stream*. All other tolerances for the other two images fall between these two values.

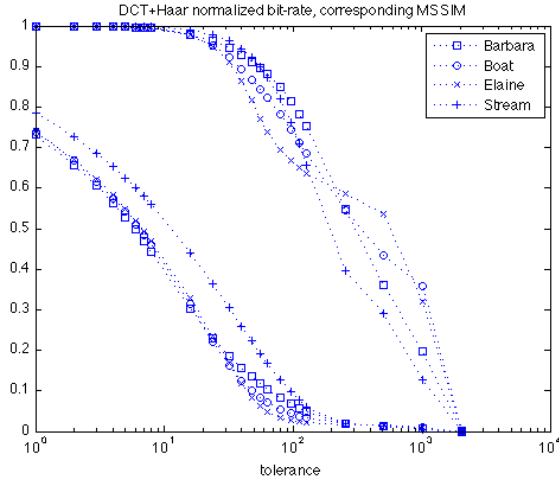


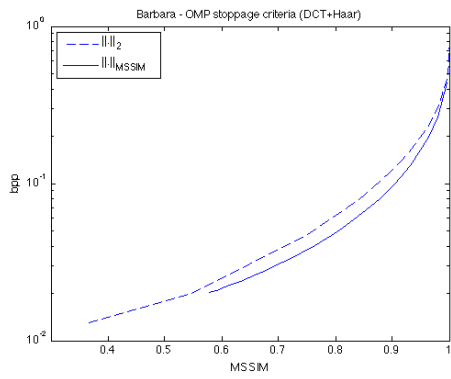
Figure 9: Normalized sparse bit-rate and corresponding MSSIM vs tolerance. In this graph we have plotted together the best normalized sparse bit-rate obtained by combining the DCT and Haar bases, and the corresponding value of the MSSIM index for a given tolerance. The normalized sparse bit-rate graphs are on the bottom left, and the MSSIM index values are above these.

This means that if we wanted all images to have an MSSIM index of 0.9 or larger, we would have to pick a tolerance no larger than $\epsilon = 32$. According to Equation (15), this tolerance corresponds to a distance on average of no more than $32/8 = 4$ units per pixel between the reconstructed image and the original. Under these circumstances we would achieve a normalized sparse bit-rate of 0.160 to 0.305 bits per pixel. It is natural to ask if there is a modification of OMP which guarantees a certain minimum MSSIM quality level. It turns out that such a modification is possible.

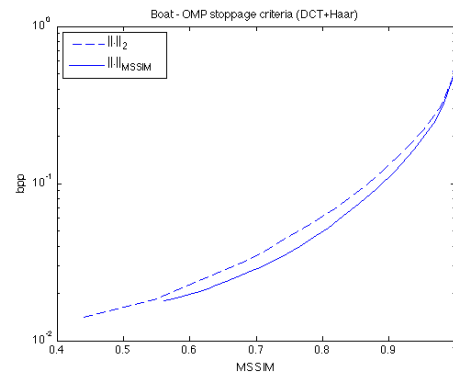
Consider the following change to the stopping rule, $\|\mathbf{Ax} - \mathbf{b}\|_2 < \epsilon$, for the OMP algorithm:

$$\|\mathbf{Ax} - \mathbf{b}\|_{MSSIM} \equiv \text{MSSIM}(c_2^{-1}(\mathbf{Ax}), c_2^{-1}(\mathbf{b})) > \delta_0,$$

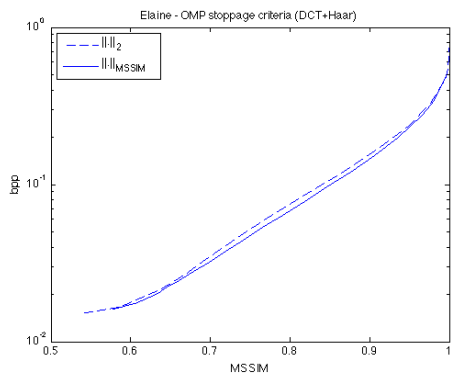
where δ_0 is a desired minimum MSSIM index value to be achieved in each individual sub-image of the reconstruction of \mathcal{I} . When we make this change, and recompute the normalized sparse bit-rate vs MSSIM graphs, we obtain the plots shown in Figure 10. In this figure, we observe that changing the stopping rule for OMP leads to an improvement in the normalized sparse bit-rate without sacrificing image quality. To see this from the opposite per-



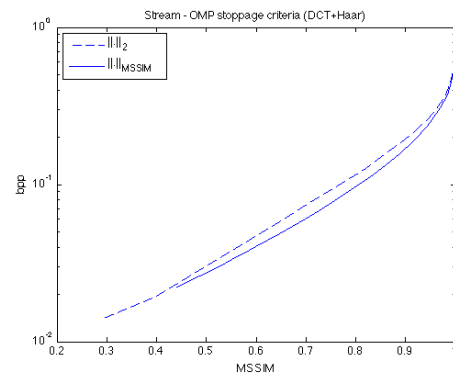
(a) Barbara



(b) Boat



(c) Elaine



(d) Stream

Figure 10: Normalized sparse bit-rate vs MSSIM. Comparison of the two different stopping rules for OMP.

spective, given a normalized sparse bit-rate, we can achieve a better MSSIM image quality index when we use the new stopping criterion. In fact, this change redistributes the work that OMP performs more evenly across the image.

Figures 11 through 18 each consist of two images, the reconstruction from the original (left) and the corresponding SSIM index map (right). The SSIM map represents the localized quality of the image reconstruction. Lighter values are values closer to 1 (“white” = 1), whereas darker values are values closer to 0 (“black” = 0). For each image \mathcal{I} we obtained a reconstruction $\tilde{\mathcal{I}}_1$ for $\epsilon = 32$ for the ℓ^2 stopping criterion, and a reconstruction $\tilde{\mathcal{I}}_2$ for the MSSIM stopping criterion choosing $\delta_0 < \text{MSSIM}(\tilde{\mathcal{I}}_1, \mathcal{I})$ in such a way that $\text{MSSIM}(\tilde{\mathcal{I}}_2, \mathcal{I}) \simeq \text{MSSIM}(\tilde{\mathcal{I}}_1, \mathcal{I})$.

3.9. Other matrices: rotations of the basis elements of $\text{DCT}_{2,3}$

Consider the matrix $\text{DCT}_{2,3}$ defined in Section 3.3.3, and recall that its columns are obtained from the tensor products of DCT waveforms, see Section 3.3.2, and Equations (9) and (10). We recall that the two-dimensional DCT-II transform plays a fundamental role in the definition of the column vectors of the matrix $\text{DCT}_{2,3}$; specifically, see Equation (8), i.e.,

$$X_{k_1, k_2} = \sum_{n_1=0}^{N_1-1} \sum_{n_2=0}^{N_2-1} x_{n_1, n_2} \cos \left[\frac{\pi}{N_1} \left(n_1 + \frac{1}{2} \right) k_1 \right] \cos \left[\frac{\pi}{N_2} \left(n_2 + \frac{1}{2} \right) k_2 \right].$$

As pointed out in Section 3.3.3, at the core of this transform we find the family of functions, indexed by k_1 and k_2 ,

$$g_{k_1, k_2, N_1, N_2}(x, y) = \cos \left[\pi \left(x + \frac{1}{2N_1} \right) k_1 \right] \cos \left[\pi \left(y + \frac{1}{2N_2} \right) k_2 \right],$$

sampled on all points $(x, y) \in \mathcal{S}(N_1) \times \mathcal{S}(N_2)$, where in our case, $N_1 = N_2 = 8$, $k_1, k_2 \in \{0, \dots, 7\}$, and $\mathcal{S}(N) = \{s_i \in [0, 1) : s_i = \frac{i}{N}, i = 0, \dots, N-1\}$.

Let $\rho(\theta)$ be the counter-clockwise rotation by angle θ in the \mathbb{R}^2 plane, and identify it with its matrix representation, i.e.,

$$\rho(\theta) = \begin{pmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{pmatrix}.$$

Then, we have the following construction for a new matrix \mathbf{A} . Let $\theta_{n,N} = \frac{\pi}{2} \frac{n}{N}$,

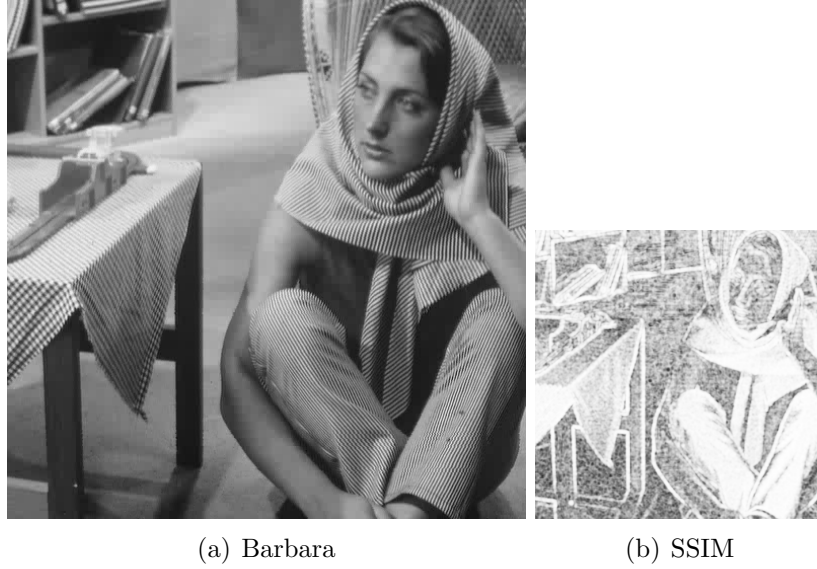


Figure 11: Barbara: $\epsilon = 32$, PSNR = 36.9952 dB, MSSIM = 0.9444, normalized sparse bit-rate = 0.1863 bpp, stopping rule: $\|\cdot\|_2$.

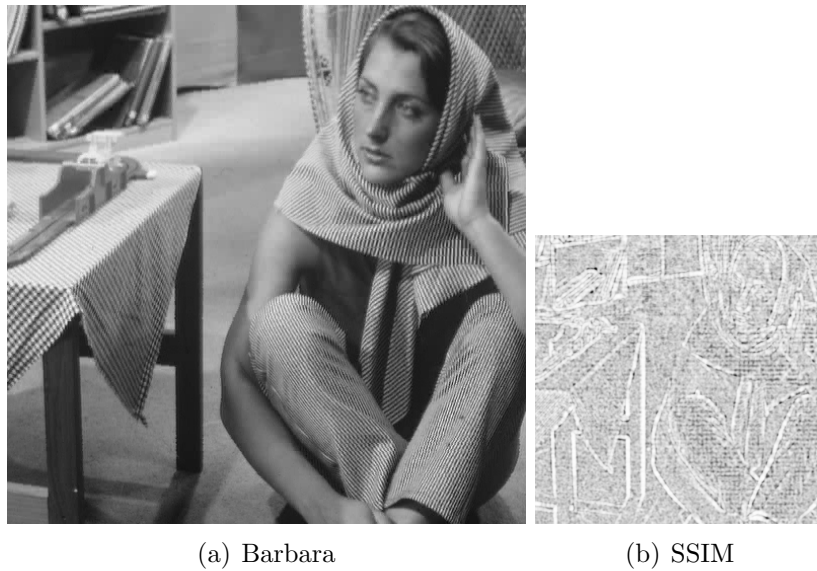


Figure 12: Barbara: $\delta_0 = 0.94$, PSNR = 32.1482 dB, MSSIM = 0.9462, normalized sparse bit-rate = 0.1539 bpp, stopping rule: $\|\cdot\|_{MSSIM}$.

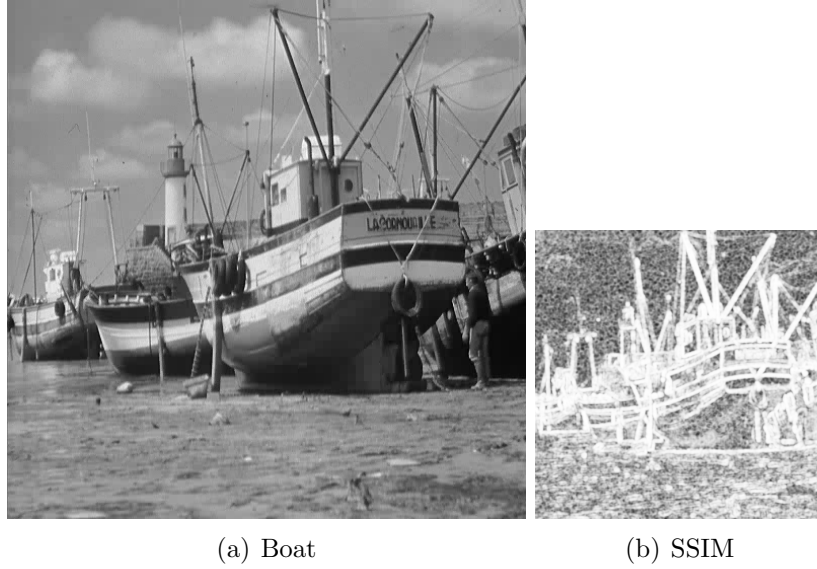


Figure 13: Boat: $\epsilon = 32$, PSNR = 36.6020 dB, MSSIM = 0.9210, normalized sparse bit-rate = 0.1608 bpp, stopping rule: $\|\cdot\|_2$.

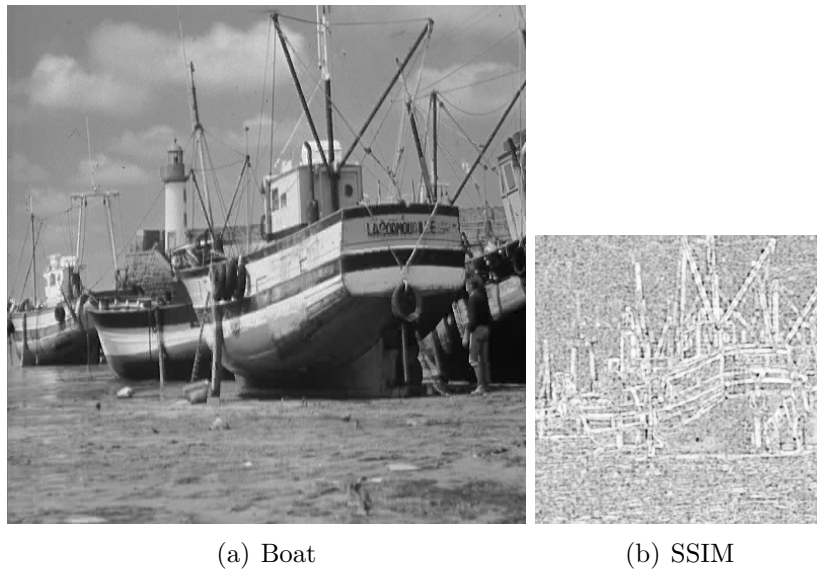


Figure 14: Boat: $\delta_0 = 0.92$, PSNR = 34.1405 dB, MSSIM = 0.9351, normalized sparse bit-rate = 0.1595 bpp, stopping rule: $\|\cdot\|_{MSSIM}$.

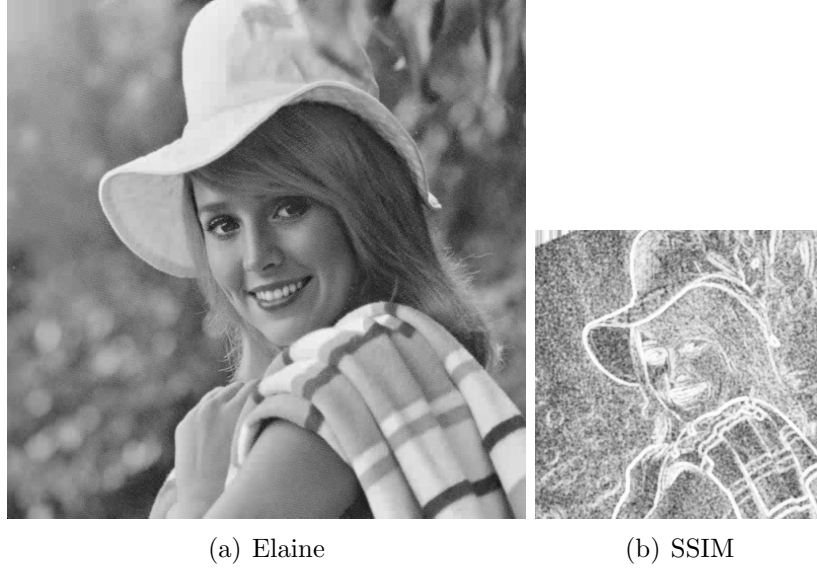


Figure 15: Elaine: $\epsilon = 32$, PSNR = 36.5155 dB, MSSIM = 0.9096, normalized sparse bit-rate = 0.1682 bpp, stopping rule: $\|\cdot\|_2$.

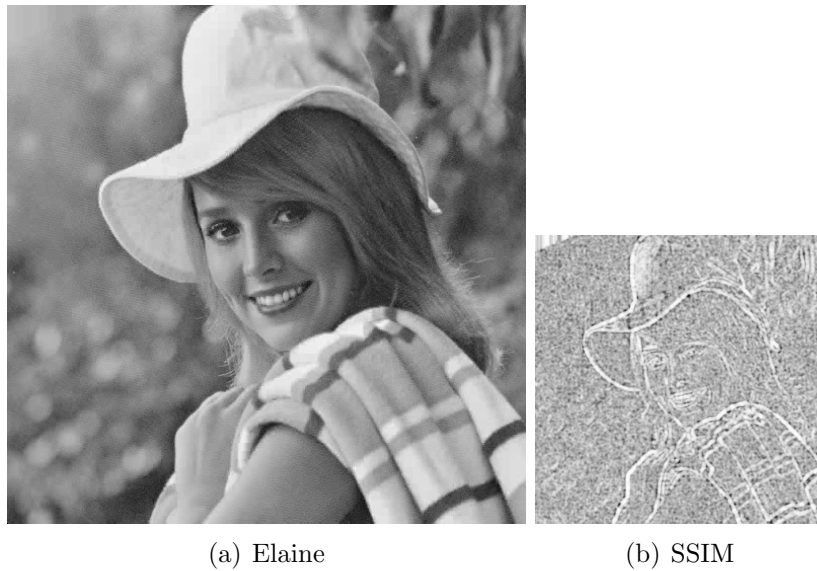


Figure 16: Elaine: $\delta_0 = 0.90$, PSNR = 35.6288 dB, MSSIM = 0.9168, normalized sparse bit-rate = 0.1686 bpp, stopping rule: $\|\cdot\|_{MSSIM}$.

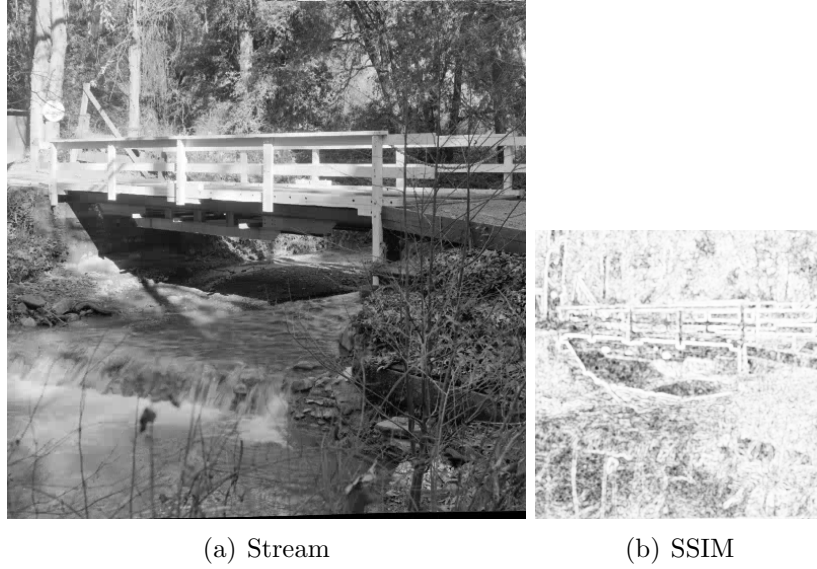


Figure 17: Stream: $\epsilon = 32$, PSNR = 36.4686 dB, MSSIM = 0.9622, normalized sparse bit-rate = 0.3050 bpp, stopping rule: $\|\cdot\|_2$.

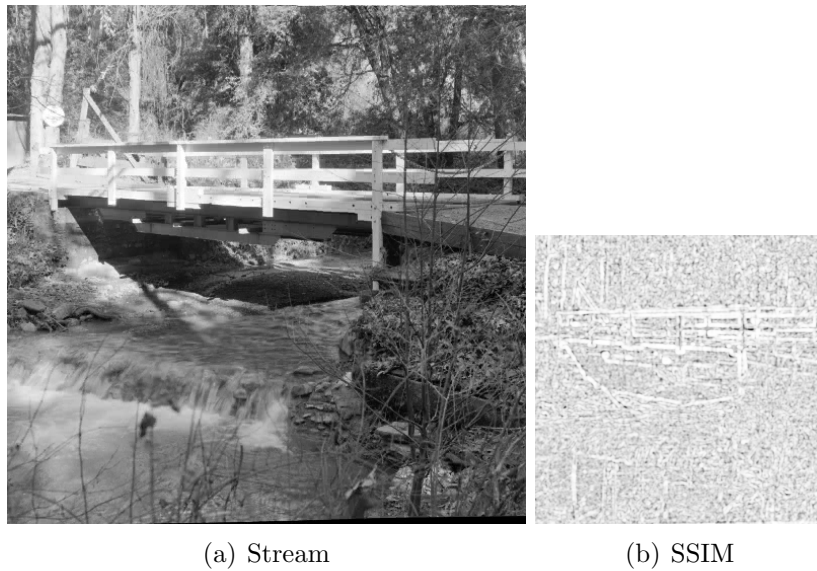


Figure 18: Stream: $\delta_0 = 0.95$, PSNR = 34.1398 dB, MSSIM = 0.9634, normalized sparse bit-rate = 0.2853 bpp, stopping rule: $\|\cdot\|_{MSSIM}$.

with $n = 0, \dots, N - 1$, be N equidistant points in $[0, \frac{\pi}{2})$. For each ordered pair $(k_1, k_2) \in \{0, \dots, 7\} \times \{0, \dots, 7\}$, and vector $\mathbf{v}_{i,j} = (\frac{i}{8}, \frac{j}{8})^T \in \mathcal{S}(8) \times \mathcal{S}(8)$, we form the matrix,

$$\mathbf{W}_{k_1, k_2}(\theta_{n, N}) = \begin{pmatrix} g_{k_1, k_2, 8, 8}(\rho(\theta_{n, N})\mathbf{v}_{0,0}) & \cdots & g_{k_1, k_2, 8, 8}(\rho(\theta_{n, N})\mathbf{v}_{0,7}) \\ \vdots & & \vdots \\ g_{k_1, k_2, 8, 8}(\rho(\theta_{n, N})\mathbf{v}_{7,0}) & \cdots & g_{k_1, k_2, 8, 8}(\rho(\theta_{n, N})\mathbf{v}_{7,7}) \end{pmatrix} \in \mathbb{R}^{8 \times 8}.$$

Using this matrix, form the column vector,

$$\mathbf{w}_{k_1, k_2, j}(\theta_{n, N}) = c_j(\mathbf{W}_{k_1, k_2}(\theta_{n, N})),$$

where the vectorization functions c_j were defined in Section 3.2.2, see Figure 1.

Choosing $j = 3$, we can compare $\mathbf{w}_{k_1, k_2, 3}(\theta_{n, N})$ with $\tilde{\mathbf{w}}_{k_1, k_2, 3}$ in Equation (10). Basically, when $n = 0$, we have $\mathbf{w}_{k_1, k_2, 3}(\theta_{0, N}) = \mathbf{w}_{k_1, k_2, 3}(0) = \tilde{\mathbf{w}}_{k_1, k_2, 3}$. We recall that $\{\tilde{\mathbf{w}}_{k_1, k_2, 3} : k_1, k_2 = 0, \dots, 7\}$ forms the columns of the $\text{DCT}_{2,3}$ submatrix in $\mathbf{A}_3 = [\text{DCT}_{2,3} \text{ Haar}_{2,3}]$, hence the name for this section.

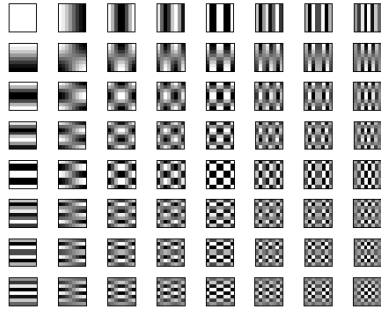
Now, traversing k_1 and k_2 in lexicographical order for the sequence of ordered pairs (k_1, k_2) with $k_1, k_2 = 0, \dots, 7$ and k_2 moving faster than k_1 , i.e., $(0,0), (0,1), \dots, (0,7), (1,0), \dots, (7,7)$, we can form a new matrix,

$$\text{DCT}_{2,3}(n, N) = (\mathbf{w}_{0,0,3}(\theta_{n, N}) \dots \mathbf{w}_{7,7,3}(\theta_{n, N})).$$

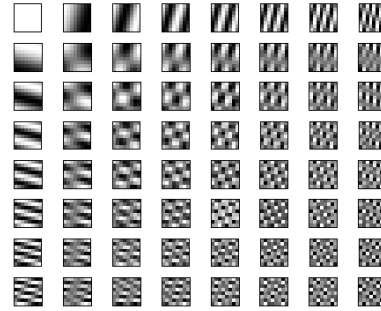
Suppose that we choose $N = 6$, and that we concatenate the six matrices $\{\text{DCT}_{2,3}(n, 6)\}_{n=0}^5$, to form the matrix,

$$\mathbf{A}_\rho(6) = [\text{DCT}_{2,3}(0, 6) \dots \text{DCT}_{2,3}(5, 6)]. \quad (16)$$

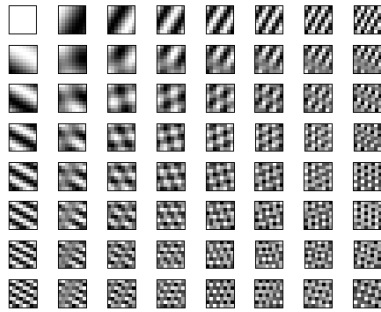
Thus, $\mathbf{A}_\rho(6)$ is a matrix that uses rotations of some of the column vectors of $[\text{DCT}_{2,3} \text{ Haar}_{2,3}]$ that could be used in problem (P_0^c) . We have run experiments using $\mathbf{A}_\rho(6)$ and a tolerance of $\epsilon = 32$ for comparison with the results obtained with $[\text{DCT}_1 \text{ Haar}_1]$ and $\mathbf{A}_j, j = 1, 2, 3$. The results are shown in Table 4. For a visual representation of the basis elements of $\mathbf{A}_\rho(6)$, see Figure 19.



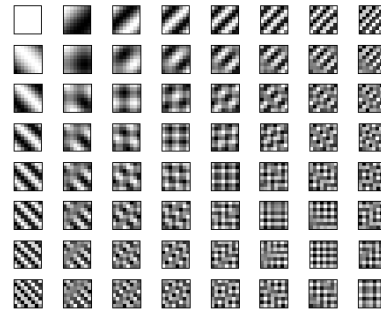
(a)



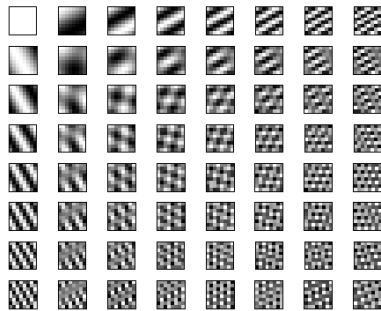
(b)



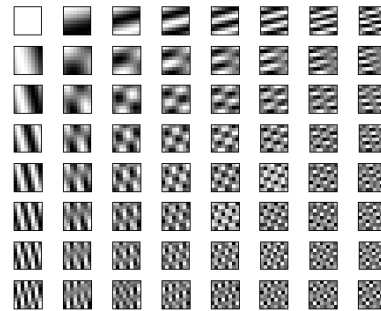
(c)



(d)



(e)



(f)

Figure 19: Full natural 2D representation for (a) $\text{DCT}_{2,3}(0,6)$, (b) $\text{DCT}_{2,3}(1,6)$, (c) $\text{DCT}_{2,3}(2,6)$, (d) $\text{DCT}_{2,3}(3,6)$, (e) $\text{DCT}_{2,3}(4,6)$, and (f) $\text{DCT}_{2,3}(5,6)$, bases for \mathbb{R}^{64} . White corresponds to the maximum value achieved by the basis element, black to the minimum. The intermediate shade of gray corresponds to 0.

Image	PSNR (dB)	Normalized sparse bit-rate (bpp)	MSSIM
Barbara	37.1180	0.1325	0.9464
Boat	36.6579	0.1321	0.9234
Elaine	36.5671	0.1359	0.9118
Stream	36.5089	0.2482	0.9636

Table 4: Performance results for $\mathbf{A} = [\text{DCT}_{2,3}(0,6) \dots \text{DCT}_{2,3}(5,6)]$. In all cases both the PSNR and normalized sparse bit-rate values were better than the values obtained for $\mathbf{A} = [\text{DCT}_{2,3} \text{ Haar}_{2,3}]$, i.e., larger and smaller, respectively. See Table 3 for comparison. The values correspond to runs of OMP with an ℓ^2 stopping rule, and a tolerance $\epsilon = 32$.

4. Quantization

4.1. Background

As a precursor to Shannon, Hartley wrote the following equation to quantify “information” in a discrete setting:

$$H = n \log s,$$

where H is the amount of information, n is the number of symbols transmitted, and s is the size of a given alphabet from which the symbols are drawn [20]. Shannon extended this notion by identifying the amount of information with entropy, see [10]. Specifically, in the case of a discrete information source, Shannon represented it as a Markov process, and asked if one could “define a quantity which will measure, in some sense, how much information is ‘produced’ by such a process, or better, at what rate information is produced?” In fact, he defined this quantity H in terms of entropy as

$$H = - \sum_{i=1}^n p_i \log_2 p_i, \quad (17)$$

where we suppose that we have a set of n possible events whose probabilities of occurrence are p_1, p_2, \dots, p_n .

To interpret Equation (17) we assume that we are given a random variable X on the finite set $\{1, 2, \dots, n\}$ with probability distribution p . The elements $X(1) = x_1, X(2) = x_2, \dots, X(n) = x_n$ are distinct and $p(x_1), p(x_2), \dots, p(x_n)$ are nonnegative real numbers with $p(x_1) + p(x_2) + \dots + p(x_n) = 1$. We write

$p_i = p(x_i)$ as a shorthand for $\text{prob}(X = x_i)$. The smaller the probability $p(x_i)$, the more uncertain we are that an observation of X will result in x_i . Thus, we can regard $1/p(x_i)$ as a measure of the uncertainty of x_i . The smaller the probability, the larger the uncertainty; see [10, 21, 22]. Shannon thought of uncertainty as information. In fact, if an event has probability 1, there is no information gained in asking the outcome of such an event given that the answer will always be the same.

Consequently, if we define the *uncertainty* of x_i to be $-\log_2 p(x_i)$, measured in *bits*, the *entropy* of the random variable X is defined to be the expected value,

$$H(X) = - \sum_{i=1}^n p(x_i) \log_2 p(x_i),$$

of the uncertainty of X , i.e., the entropy of X will measure the information gained from observing X .

Further, Shannon defined the capacity C of a discrete noiseless channel as

$$C = \lim_{T \rightarrow \infty} \frac{\log N(T)}{T},$$

where $N(T)$ is the number of allowed signals of duration T .

Theorem 3 (Shannon, Fundamental Theorem for a Noiseless Channel [10]). *Let a source have entropy H (bits per symbol) and let a channel have a capacity C (bits per second). Then it is possible to encode the output of the source in such a way as to transmit at the average rate $\frac{C}{H} - \epsilon$ symbols per second over the channel where ϵ is arbitrarily small. It is not possible to transmit at an average rate greater than $\frac{C}{H}$.*

This existential result has been a driving force for developing constructive quantization and coding theory through the years, e.g., [23] for an extraordinary engineering perspective and highlighting the role of redundancy in source signals, cf., the *theory of frames* [24, Chap. 3 and 7] and [25, 26, 27].

4.2. Quantization (coding), rate, and distortion

A *scalar quantizer* can be defined as consisting of a set \mathcal{S} of intervals or *cells* $S_i \subset \mathbb{R}, i \in I$, that form a partition of the real line, where the index set I is ordinarily a collection of consecutive integers beginning with 0 or 1, together with a set \mathcal{C} of *reproduction values* or *levels* $y_i \in \mathbb{R}, i \in I$, so that

the overall quantizer q is defined by $q(x) = y_i$ for $x \in S_i$, expressed concisely as

$$q(x) = \sum_{i \in I} y_i \mathbb{1}_{S_i}(x), \quad (18)$$

where the indicator function $\mathbb{1}_S(x)$ is 1 if $x \in S$ and 0 otherwise [23, 28].

More generally, a class of *memoryless quantizers* can be described as follows. A quantizer of dimension $k \in \mathbb{N}$ takes as input a vector $\mathbf{x} = (x_1, \dots, x_k)^T \in \mathcal{A} \subseteq \mathbb{R}^k$. Memoryless refers to a quantizer which operates independently on successive vectors. The set \mathcal{A} is called the *alphabet* or *support* of the source distribution. If $k = 1$ the quantizer is *scalar*, and, otherwise, it is *vector*. The quantizer then consists of three components: a *lossy encoder* $\alpha : \mathcal{A} \rightarrow I$, where the index set I is an arbitrary countable set; a *reproduction decoder* $\beta : I \rightarrow \hat{\mathcal{A}}$, where $\hat{\mathcal{A}} \subset \mathbb{R}^k$ is the *reproduction alphabet*; and a *lossless encoder* $\gamma : I \rightarrow J$, an invertible mapping (with probability 1) into a collection J of variable-length binary vectors that satisfies the *prefix condition*, that is, no vector in J can be the prefix of any other vector in the collection [23].

Alternatively, a lossy encoder is specified by a partition $\mathcal{S} = \{S_i \subset \mathbb{R} : i \in I\}$ of \mathcal{A} ; a reproduction decoder is specified by a *codebook* $\mathcal{C} = \{\beta(i) \in \hat{\mathcal{A}} : i \in I\}$ of *points*, *codevectors*, or *reproduction codewords*, also known as the *reproduction codebook*; and the lossless encoder γ can be described by its *binary codebook* $J = \{\gamma(i) : i \in I\}$ containing *binary* or *channel codewords*. The *quantizer rule* is the function $q(\mathbf{x}) = \beta(\alpha(\mathbf{x}))$ or, equivalently, $q(\mathbf{x}) = \beta(i)$ whenever $\mathbf{x} \in S_i$ [23].

The instantaneous rate of a quantizer applied to a particular input is the normalized length $r(\mathbf{x}) = \frac{1}{k}l(\gamma(\alpha(\mathbf{x})))$ of the channel codeword, the number of bits per source symbol that must be sent to describe the reproduction. If all binary codewords have the same length, it is referred to as a *fixed-length* or *fixed-rate* quantizer.

To measure the quality of the reproduction, we assume the existence of a nonnegative distortion measure $d(\mathbf{x}, \hat{\mathbf{x}})$ which assigns a distortion or cost to the reproduction of input \mathbf{x} by $\hat{\mathbf{x}}$. Ideally, one would like a distortion measure that is easy to compute, useful in analysis, and perceptually meaningful in the sense that small (large) distortion means good (poor) perceived quality. No single distortion measure accomplishes all three goals [23]. However, $d(\mathbf{x}, \hat{\mathbf{x}}) = \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2$ satisfies the first two.

We also assume that $d(\mathbf{x}, \hat{\mathbf{x}}) = 0$ if and only if $\mathbf{x} = \hat{\mathbf{x}}$. In this light we

say that a code is *lossless* if $d(\mathbf{x}, \beta(\alpha(\mathbf{x}))) = 0$ for all inputs \mathbf{x} , and *lossy* otherwise.

Finally, the overall performance of a quantizer applied to a source is characterized by the *normalized rate*,

$$\begin{aligned} R(\alpha, \gamma) &= E[r(X)] = \frac{1}{k} E[l(\gamma(\alpha(X)))] \\ &= \frac{1}{k} \sum_i l(\gamma(i)) \int_{S_i} f(\mathbf{x}) d\mathbf{x}, \end{aligned} \quad (19)$$

and the *normalized average distortion*,

$$\begin{aligned} D(\alpha, \beta) &= \frac{1}{k} E[d(X, \beta(\alpha(X)))] \\ &= \frac{1}{k} \sum_i \int_{S_i} d(\mathbf{x}, \mathbf{y}_i) f(\mathbf{x}) d\mathbf{x}. \end{aligned} \quad (20)$$

Here, we assume that the quantizer operates on a k -dimensional random vector $X = (X_1, \dots, X_k)$ that is described by a probability density function $f(x)$. Every quantizer (α, γ, β) is thus described by a rate-distortion pair $(R(\alpha, \gamma), D(\alpha, \beta))$. The goal of a compression system design is to optimize the rate-distortion trade-off [23].

In light of the results by Shannon in Section 4.1, compression system design will also have to take into account the characteristics of the communication channel in managing the rate-distortion trade-off. Also, from the definitions above, it is clear that knowledge of the probability density function of the source messages is relevant, see Figures 21 and 22.

4.3. Image quantization and encoding

With the perspective from Sections 4.1 and 4.2, we return to the topic of image quantization and encoding with regard to compression. The image standards, JPEG and JPEG 2000, are framed in the transform coding paradigm, and contain two steps beyond their respective discrete cosine transform (DCT) and the Cohen-Daubechies-Feauveau 5/3 (lossless) and 9/7 (lossy) biorthogonal wavelet transforms. Both have *quantization* (α) and *encoding* (γ) steps, with their respective “inverses”, to complete their definitions [2, 3, 4].

The general schematic for transform coding is described in Figure 20.

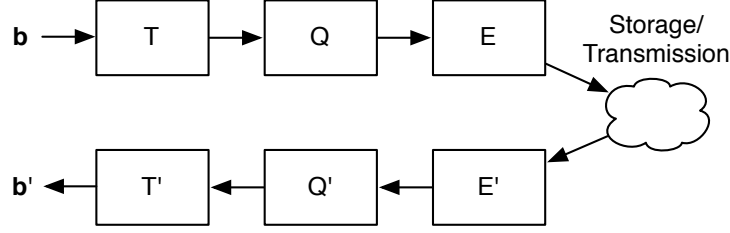


Figure 20: Schematic diagram of a general transform coding system.

Example 1. The transform T is meant to exploit the redundancy in the source signal \mathbf{b} and decorrelate it. It has an inverse T^{-1} or, minimally, a left inverse T' such that $T'T\mathbf{b} = \mathbf{b}$. In our approach we have $\mathbf{A} = T'$, and T is defined via OMP by $T\mathbf{b} = \text{OMP}(\mathbf{A}, \mathbf{b}, \epsilon_0)$. Thus, we have $\|T'T\mathbf{b} - \mathbf{b}\|_2 < \epsilon$. Therefore, our compression scheme is lossy. Q is a non-invertible scalar quantizer that will be applied to the coefficients of the vector $T\mathbf{b}$, and Q' is its reproduction function. Finally, we have an invertible lossless encoder E , defined as γ in the previous section, with $E' = E^{-1}$. The composition QT is equivalent to the lossy encoder α , and the composition $T'Q'$ corresponds to the reproduction decoder β from Section 4.2.

In order to describe the overall performance of our quantizer

$$(\alpha, \gamma, \beta) = (QT, E, T'Q'),$$

we must characterize the rate-distortion pair $(R(QT, E), D(QT, T'Q'))$.

Proposition 4. Let $n < m$ and let $\mathbf{A} = (\mathbf{a}_j) \in \mathbb{R}^{n \times m}$ be a full-rank matrix with each $\|\mathbf{a}_j\|_2 = c$. Given $a > 0$ and $\mathbf{y} \in \mathbb{R}^n$. Suppose that $\mathbf{x} \in \mathbb{R}^m$ has the property that $a\mathbf{A}\mathbf{x} = \mathbf{y}$ and that $\epsilon \in \mathbb{R}^m$ satisfies $\|\epsilon\|_0 \leq \|\mathbf{x}\|_0$. If $\tilde{\mathbf{x}} = \mathbf{x} + \epsilon$ and $\tilde{\mathbf{y}} = a\mathbf{A}\tilde{\mathbf{x}}$, then

$$\|\tilde{\mathbf{y}} - \mathbf{y}\|_2 \leq ac\|\epsilon\|_\infty\|\mathbf{x}\|_0. \quad (21)$$

Proof. Let $\epsilon = (\epsilon_1, \dots, \epsilon_m)^T \in \mathbb{R}^m$ with $\|\epsilon\|_0 \leq \|\mathbf{x}\|_0$ and let $\tilde{\mathbf{x}} = \mathbf{x} + \epsilon$. Then,

we compute

$$\begin{aligned}
\|\tilde{\mathbf{y}} - \mathbf{y}\|_2 &= \|a\mathbf{A}\tilde{\mathbf{x}} - a\mathbf{A}\mathbf{x}\|_2 = a\|\mathbf{A}(\mathbf{x} + \epsilon) - \mathbf{A}\mathbf{x}\|_2 \\
&= a\|\mathbf{A}\epsilon\|_2 = a\left\|\sum_{j=1}^m \mathbf{a}_j\epsilon_j\right\|_2 = a\left\|\sum_{\epsilon_j \neq 0} \mathbf{a}_j\epsilon_j\right\|_2 \\
&\leq a\sum_{\epsilon_j \neq 0} \|\mathbf{a}_j\epsilon_j\|_2 = a\sum_{\epsilon_j \neq 0} \|\mathbf{a}_j\|_2 |\epsilon_j| = a\sum_{\epsilon_j \neq 0} c|\epsilon_j| \\
&\leq a\sum_{\epsilon_j \neq 0} c\|\epsilon\|_\infty = ac\|\epsilon\|_\infty \|\epsilon\|_0 \leq ac\|\epsilon\|_\infty \|\mathbf{x}\|_0
\end{aligned}$$

□

Remark 1. a. Moreover, the value of $\|\epsilon\|_0$ is linked to the sparsity of \mathbf{x} , because in our case the error ϵ comes from scalar quantizing the entries of \mathbf{x} . That is, if $\tilde{\mathbf{x}} = \text{round}(\mathbf{x})$, where $\tilde{\mathbf{x}}$ is the vector whose entries are exactly those of \mathbf{x} but rounded to the closest integer, then necessarily

$$\|\epsilon\|_0 = \|\tilde{\mathbf{x}} - \mathbf{x}\|_0 \leq \|\mathbf{x}\|_0. \quad (22)$$

Hence, in the case where a scalar quantization scheme satisfies inequality (22), Proposition 4 gives

$$\|\tilde{\mathbf{y}} - \mathbf{y}\|_2 \leq ac\|\epsilon\|_\infty \|\mathbf{x}\|_0, \quad (23)$$

with $\mathbf{A} = (\mathbf{a}_j)$, $\tilde{\mathbf{y}} = a\mathbf{A}\tilde{\mathbf{x}}$, $\mathbf{y} = a\mathbf{A}\mathbf{x}$, and $c = \|\mathbf{a}_j\|_2$ for all j . Observe that, in particular, when $\tilde{\mathbf{x}} = \text{round}(\mathbf{x})$, we have $\|\epsilon\|_\infty = 1/2$.

From Equation (23) we see that the error in the reconstruction due to scalar quantization is linked to the size of c , $\|\epsilon\|_\infty$, and the sparsity of \mathbf{x} . We are tempted to make c as small as possible, or modify the quantization scheme to make $\|\epsilon\|_\infty$ smaller to reduce the reconstruction error due to scalar quantization.

b. If \mathbf{x}_0 is the sparsest vector that solves $\mathbf{A}\mathbf{x} = \mathbf{y}$, then $a^{-1}\mathbf{x}_0$ is the sparsest vector that solves $a\mathbf{A}\mathbf{x} = \mathbf{y}$, and $\|\mathbf{x}_0\|_0 = \|a\mathbf{x}_0\|_0$. We conclude that the norm of $a^{-1}\mathbf{x}_0$ has an inversely proportional relationship with the size of a . Therefore, if we are to use a finite predetermined number of bits to represent $a^{-1}\mathbf{x}_0$, the solution of $a\mathbf{A}\mathbf{x} = \mathbf{y}$, we necessarily have a constraint on a .

c. We know that the magnitude of the coordinates of \mathbf{x} are bounded by a

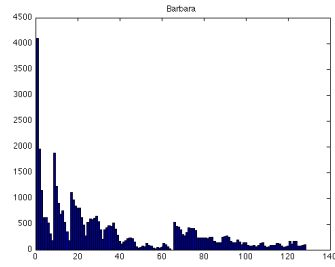
multiple of $\|\mathbf{y}\|_2$, see [29]. This has an impact on how many bits are needed to represent \mathbf{x} . Therefore, when choosing a in Proposition 4 we have to take into consideration the maximum value that the value of $\|\mathbf{y}\|_2$ will impose on the magnitude of the coordinates of \mathbf{x} , see Examples 2 and 3.

Finally, let us recall that our image compression approach comes from the OMP solution \mathbf{x}_0 to problem $(P_0^{\epsilon_0})$ for a given matrix $\mathbf{A} = (\mathbf{a}_j)$ whose column vectors satisfy $\|\mathbf{a}_j\|_2 = c$, a vector \mathbf{b} , and a tolerance $\epsilon_0 > 0$ (such that $\|\mathbf{A}\mathbf{x}_0 - \mathbf{b}\|_2 < \epsilon_0$.) Then, choosing $a > 0$, and following the description of T at the beginning of this section, if we set $\mathbf{x}_0 = T\mathbf{b} = \text{OMP}(a\mathbf{A}, \mathbf{b}, \epsilon_0)$ for a given signal \mathbf{b} and a tolerance $\epsilon_0 > 0$, $a\mathbf{A} = T'$, Q a scalar quantizer that satisfies inequality $\|\epsilon\|_0 = \|Q(\mathbf{x}) - \mathbf{x}\|_0 \leq \|\mathbf{x}\|_0$, and Q' its corresponding reproduction function, the triangle inequality and inequality (21) give

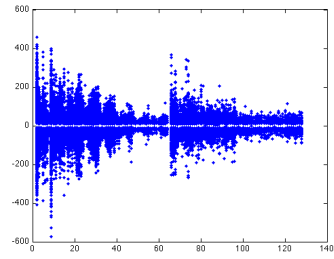
$$\begin{aligned} d(\beta(\alpha(\mathbf{b})), \mathbf{b}) &= \|T'Q'QT\mathbf{b} - \mathbf{b}\|_2 \\ &= \|T'Q'QT\mathbf{b} - T'T\mathbf{b} + T'T\mathbf{b} - \mathbf{b}\|_2 \\ &= \|a\mathbf{A}\tilde{\mathbf{x}}_0 - a\mathbf{A}\mathbf{x}_0 + a\mathbf{A}\mathbf{x}_0 - \mathbf{b}\|_2 \\ &\leq \|a\mathbf{A}\tilde{\mathbf{x}}_0 - a\mathbf{A}\mathbf{x}_0\|_2 + \|a\mathbf{A}\mathbf{x}_0 - \mathbf{b}\|_2 \\ &= ac\|\delta\|_\infty\|\mathbf{x}_0\|_0 + \epsilon_0, \end{aligned}$$

where $\delta = \tilde{\mathbf{x}}_0 - \mathbf{x}_0$. This inequality would give us a footing in the computation of the normalized average distortion $D(\alpha, \beta)$, see Equation (20).

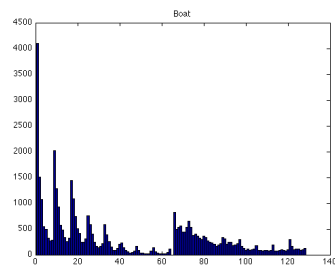
Example 2. From the definition of $D(\alpha, \beta)$, it is clear that we need to know something about the probability density function of the input sources, i.e., the statistics of the 8×8 vectorized sub-images into which each image is partitioned, if we are to compute D . In place of such knowledge, we can observe the distribution of the coefficients for each of the vectors resulting from the analysis of the images in our database and their corresponding histograms. This is what Figure 21 shows. For each such image \mathcal{I} , we used the matrix $\mathbf{A}_3 = [\text{DCT}_{2,3} \text{ Haar}_{2,3}] = (\mathbf{a}_i)$ with a tolerance of $\epsilon = 32$ to compute its statistics. On the x-axis of each subfigure in Figure 21 we have matched column \mathbf{a}_i at position i to the integer i . Hence, positions 1 to 64 correspond to the DCT waveforms, and positions 65 to 128 to the Haar waveforms. All subfigures on the left are the histograms for the frequency with which each column vector of \mathbf{A}_3 is chosen. For example, since there are 4096 sub-images of size 8×8 in a 512×512 image, the column vector \mathbf{a}_1 will be chosen 4096 times since it corresponds to the constant vector, which computes



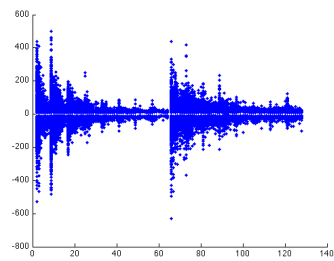
(a)



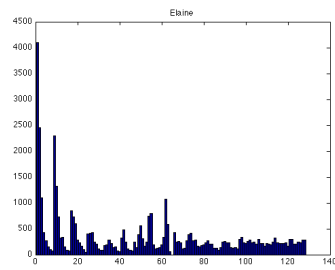
(b)



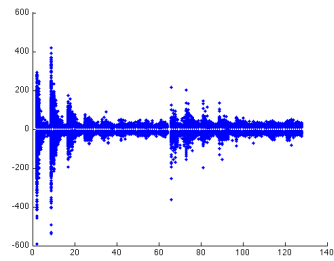
(c)



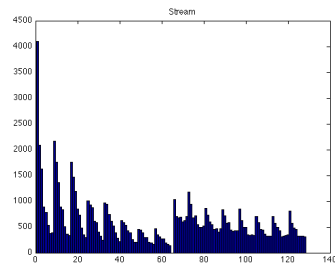
(d)



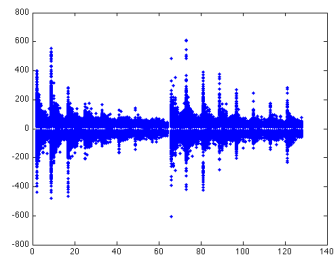
(e)



(f)



(g)



(h)

Figure 21: Histograms for images *Barbara*, *Boat*, *Elaine*, and *Stream*. We set the tolerance $\epsilon = 32$, and used $\mathbf{A}_3 = [\text{DCT}_{2,3} \text{ Haar}_{2,3}]$.

the mean or DC component for each sub-image. All subfigures to the right correspond to partial representations of the distribution of the coefficients that multiply each and every column vector whenever such a vector is chosen in the representation/approximation of an input \mathbf{b} . For example, suppose that column \mathbf{a}_{74} was multiplied by a coefficient $a_{74} = 3.2310$ to obtain the representation of some input $\mathbf{b} = a_{74}\mathbf{a}_{74} + \mathbf{r}$ within a tolerance of $\epsilon = 32$. Then we would have plotted point $(74, 3.2310)$ in its corresponding subfigure to the right. We have not plotted the coefficients for \mathbf{a}_1 since they correspond to the DC components of the sub-images of \mathcal{I} , which vary between 0 and 255. We note that all images in our database have a similar structure.

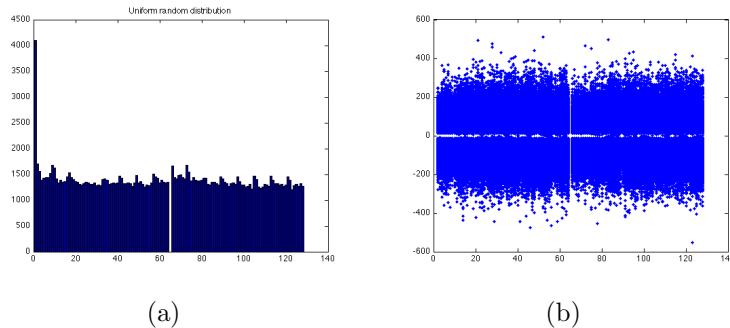
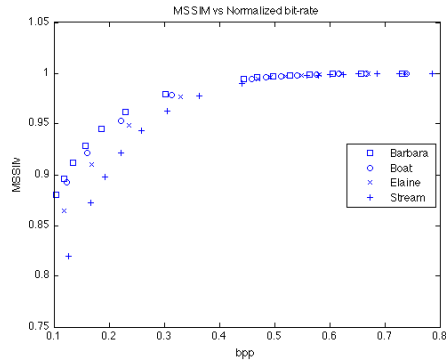


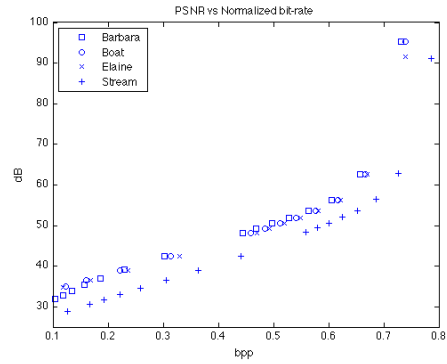
Figure 22: Histograms for uniform random input. $\mathbf{A}_3 = [\text{DCT}_{2,3} \text{ Haar}_{2,3}]$, with $\epsilon = 32$.

Example 3. For comparison purposes, we obtained the histogram and the distribution of coefficients for a randomly generated image with a uniform distribution on $[0 \ 255]$, see Figure 22. The first thing we note is that unlike the natural images in our database, all column vectors, except \mathbf{a}_1 and \mathbf{a}_{65} , which correspond to the constant vectors (one for the DCT and one for the Haar waveforms), are chosen about the same number of times regardless of their position. Further, the distribution of the values of the coefficients is uniform. It is also clear that, in order to be reconstructed, this is the image that requires the most nonzero coefficients within the tolerance $\epsilon = 32$. This is consistent with the definition of information by Shannon: the more the uncertainty in a source, the more information it carries.

Example 4. We have plotted the distortion as measured by the MSSIM index and the PSNR versus the idealized normalized sparse bit-rate, see Figure

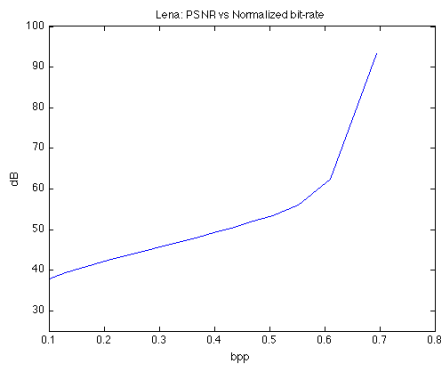


(a)

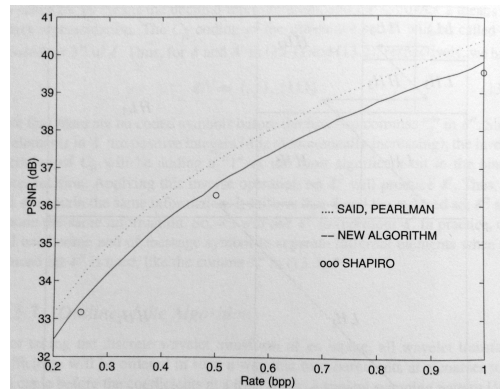


(b)

Figure 23: (a) MSSIM vs Normalized sparse bit-rate, (b) PSNR vs Normalized sparse bit-rate.



(a)



(b)

Figure 24: PSNR vs bit-rate: (a) Normalized sparse bit-rate results for $\mathbf{A} = [\text{DCT}_1 \text{ Haar}_1]$ prior to any γ coding, and (b) bit-rate coding performances published in [18] for image *Lena*: Said and Pearlman’s SPIHT algorithm [30], Embedded Coding and the Wavelet-Difference-Reduction compression algorithm (“new algorithm”) [31], and Shapiro’s EZW algorithm [32].

23. This bit-rate is unattainable in practice but nonetheless gives us an idea of an upper bound in the rate-distortion trade-off, and lets us compare how much room we have to select a lossless encoder γ to complete the implementation of a quantizer using our sparse image representation approach. Figure 23 shows the MSSIM and PSNR versus normalized sparse bit-rate trade-off. Figure 24(a) shows the PSNR versus normalized sparse bit-rate trade-off for image *Lena*, which we computed to compare with Figure 24(b) which shows results for that image for three published fully implemented quantizers. We observe that there is enough room to pick an encoder γ that could compete with these implementations.

Regarding the computation of the rate $R(\alpha, \gamma)$ for our image quantizer, we would have to choose a lossless encoder γ , which we have not done here.

5. Acknowledgement

The first name author gratefully acknowledges the support of the Institute for Physical Science and Technology of the University of Maryland, College Park. The second named author gratefully acknowledges the support of MURI-ARO Grant W911NF-09-1-0383 and NGA Grant 1582-08-1-0009. We are both appreciative of expert observations by Professor Ramani Duraiswami.

References

- [1] A. M. Bruckstein, D. L. Donoho, M. Elad, From sparse solutions of systems of equations to sparse modeling of signals and images, *SIAM Review* 51 (2009) 34–81.
- [2] G. K. Wallace, The JPEG still picture compression standard, *Communications of the ACM* 34 (1991) 30–44.
- [3] D. S. Taubman, M. W. Marcellin, *JPEG 2000: Image Compression Fundamentals, Standards and Practice*, Kluwer Academic Publishers, Norwell, MA, second edition, 2002.
- [4] C. Christopoulos, A. Skodras, T. Ebrahimi, The JPEG 2000 still image coding system: an overview, *IEEE Transactions on Consumer Electronics* 46 (2000) 1103–1127.

- [5] B. K. Natarajan, Sparse approximate solutions to linear systems, *SIAM Journal on Computing* 24 (1995) 227–234.
- [6] Viterbi School of Engineering, University of Southern California, The USC-SIPI Image Database, <http://sipi.usc.edu/database/>, 2012.
- [7] M. V. Wickerhauser, *Adapted Wavelet Analysis from Theory to Software*, A K Peters, Ltd., 1996.
- [8] Wikipedia - The Free Encyclopedia, Peak Signal-to-Noise Ratio, http://en.wikipedia.org/wiki/Peak_signal-to-noise_ratio, 2012.
- [9] Z. Wang, A. C. Bovik, H. R. Sheikh, E. P. Simoncelli, Image quality assessment: From error measurement to structural similarity, *IEEE Transactions on Image Processing* 13 (2004) 1–14.
- [10] C. E. Shannon, A mathematical theory of communication, *The Bell System Technical Journal* 27 (1948) 379–423, 623–656.
- [11] P. J. Huber, Projection pursuit, *Annals of Statistics* 13 (1985) 435–475.
- [12] D. Donoho, I. Johnstone, P. Rousseeuw, W. Stahel, Discussion: Projection pursuit, *Annals of Statistics* 13 (1985) 496–500.
- [13] Y. Pati, R. Rezaifar, P. Krishnaprasad, Orthogonal matching pursuit: recursive function approximation with application to wavelet decomposition, in: *27th Asilomar Conference on Signals, Systems and Computers*, 1993, pp. 40–44.
- [14] S. G. Mallat, Z. Zhang, Matching pursuits with time-frequency dictionaries, *IEEE Transactions on Signal Processing* 41 (1993) 3397–3415.
- [15] W. L. Briggs, V. E. Henson, *The DFT, an Owner’s Manual for the Discrete Fourier Transform*, SIAM, Philadelphia, PA, 1995.
- [16] K. R. Rao, P. Yip, *Discrete Cosine Transform: Algorithms, Advantages, Applications*, Academic Press Professional, Inc., San Diego, CA, USA, 1990.
- [17] S. G. Mallat, *A Wavelet Tour of Signal Processing*, Academic Press, San Diego, CA, 1998.

- [18] H. L. Resnikoff, R. O. Wells, Jr., *Wavelet Analysis. The Scalable Structure of Information*, Springer-Verlag, New York, NY, 1998. Corrected 2nd printing.
- [19] A. B. Watson (Ed.), *Digital Images and Human Vision*, MIT Press, Cambridge, MA, 1993.
- [20] J. Gleick, *The Information: a History, a Theory, a Flood*, Pantheon Books, New York, NY, 2011.
- [21] W. C. Huffman, V. Pless (Eds.), *Handbook of Coding Theory*, volume 1, Elsevier Science B. V., Amsterdam, The Netherlands, 1998.
- [22] W. C. Huffman, V. Pless, *Fundamentals of Error-Correcting Codes*, Cambridge University Press, New York, NY, 2010.
- [23] R. M. Gray, D. L. Neuhoff, Quantization, *IEEE Transactions on Information Theory* 44 (1998) 2325–2383.
- [24] J. J. Benedetto, M. W. Frazier, *Wavelets: Mathematics and Applications*, CRC Press, Boca Raton, FL, 1994.
- [25] J. J. Benedetto, A. M. Powell, O. Yilmaz, Sigma-delta ($\Sigma\Delta$) quantization and finite frames, *IEEE Transactions on Information Theory* 52 (2006) 1990–2005.
- [26] O. Christensen, *An Introduction to Frames and Riesz Bases*, Springer-Birkhäuser, NY, 2003.
- [27] P. G. Casazza, J. Kovačević, Uniform tight frames with erasures, *Advances in Computational Mathematics* 18 (2003) 387–430.
- [28] J. C. Candy, G. C. Temes (Eds.), *Oversampling Delta-Sigma Data Converters*, IEEE Press, NY, 1992.
- [29] J. J. Benedetto, A. Nava-Tudela, *Frame estimates for OMP*, 2012. Preprint.
- [30] A. Said, W. A. Pearlman, A new, fast, and efficient image codec based on set partitioning in hierarchical trees, *IEEE Transactions on Circuits and Systems for Video Technology* 6 (1996) 243–250.

- [31] J. Tian, R. O. Wells, Jr., A lossy image codec based on index coding, in: IEEE Data Compression Conference, DCC '96, p. 456.
- [32] J. M. Shapiro, Embedded image coding using zerotrees of wavelet coefficients, IEEE Transactions on Signal Processing 41 (1993) 3445–3462.