

BAYESIAN WAVELET-DOMAIN IMAGE MODELING USING HIDDEN MARKOV TREES

Justin K. Romberg, Hyeokho Choi and Richard G. Baraniuk

Department of Electrical and Computer Engineering, Rice University, Houston, TX 77251–1892, USA

ABSTRACT

Wavelet-domain hidden Markov models have proven to be useful tools for statistical signal and image processing. The hidden Markov tree (HMT) model captures the key features of the joint statistics of the wavelet coefficients of real-world data. One potential drawback to the HMT framework is the need for computationally expensive iterative training (using the EM algorithm, for example). In this paper, we propose two reduced-parameter HMT models that capture the general structure of a broad class of grayscale images. The image HMT (iHMT) model leverages the fact that for a large class of images the structure of the HMT is self-similar across scale. This allows us to reduce the complexity of the iHMT to just nine easily trained parameters (independent of the size of the image and the number of wavelet scales). In the universal HMT (uHMT) we take a Bayesian approach and fix these nine parameters. The uHMT requires no training of any kind. While simple, we show using a series of image estimation/denoising experiments that these two new models retain nearly all of the key structures modeled by the full HMT. Based on these new models, we develop a shift-invariant wavelet denoising scheme that outperforms all algorithms in the current literature.

1. INTRODUCTION

Statistical image processing problems, such as estimation, detection, and classification, rely on knowledge of the joint probability density function (pdf), $f(\mathbf{x})$, of the image \mathbf{x} . Since $f(\mathbf{x})$ is usually not known or is too complex to specify exactly, models that accurately approximate $f(\mathbf{x})$ are critical to image processing algorithms.

There have been several approaches to modeling the local joint statistics of image pixels in the spatial domain, the Markov random field model [1] being the most prevalent. However, spatial-domain models are limited in their ability to describe large-scale behavior. Markov random field models can be improved by incorporating a larger neighborhood of pixels, but this rapidly increases their complexity.

Transform-domain models are based on the idea that often a linear, invertible transform will “restructure” the

This work was supported by the National Science Foundation, grant no. MIP-9457438, DARPA, grant no. DARPA/AFOSR F49620-97-1-0513, ONR, grant no. N00014-99-1-0813, Texas Instruments, and the Rice Consortium for Computational Seismic/Signal Interpretation.

Email: jrom@rice.edu, choi@ece.rice.edu, richb@rice.edu
Web: www.dsp.rice.edu

image, leaving transform coefficients whose structure is simpler to model. Real-world images are well characterized by their *singularity* (edge and ridge) structure. The wavelet transform captures this singularity structure, and provides a natural and powerful domain for image modeling [2]. We aim, therefore, to model the joint pdf of the wavelet coefficients \mathbf{w} of an image, and design wavelet-domain processing algorithms based on this model of $f(\mathbf{w})$.

2. IMAGES IN THE WAVELET DOMAIN

The wavelet transform is an atomic decomposition of an image with basis functions that are shifted and dilated versions of an oscillating mother wavelet [2]. The *primary properties* of wavelet transforms make wavelet-domain statistical image processing attractive [2, 3]:

- P1. Locality:** Each wavelet coefficient represents the image content localized in spatial location and frequency.
- P2. Multiresolution:** The wavelet transform represents the image at a nested set of scales.
- P3. Edge Detection:** Wavelets act as local edge detectors. The edges in the image are represented by large wavelet coefficients at the corresponding spatial locations.
- P4. Decorrelation:** The wavelet coefficients of real-world images tend to be approximately decorrelated.
- P5. Energy Compaction:** The wavelet transforms of real-world images tend to be sparse. A wavelet coefficient is large only if edges are present within the support of the wavelet.

Properties **P1** and **P2** lead to a natural arrangement of the wavelet coefficients in a quadtree structure with three subbands representing the horizontal, vertical, and diagonal edges in the image (see Fig. 1). The Compaction property (**P5**) follows from the fact that the edges constitute only a very small portion of a typical image; consequently, we can closely approximate an image by just a few (large) wavelet coefficients. Furthermore, the Decorrelation property (**P4**) indicates that the dependencies between wavelet coefficients are predominantly local. The primary properties give wavelet transforms significant structure, which we codify in the following *secondary properties*:

- S1. NonGaussianity:** The wavelet coefficients have peaky, heavy-tailed marginal distributions.
- S2. Persistency:** Large/small values of wavelet coefficients tend to propagate through the scales of the quadtrees.

NonGaussianity follows immediately from Energy Compaction (**P5**). Persistency follows from the Edge Detection (**P3**) and Multiresolution (**P2**) properties.

In general, the wavelet coefficients \mathbf{w} are indexed by two integers: one for the scale (dilation), and one for the shift. In this paper, we will adopt an abstract indexing system and use only one integer whose value ranges from 1 to N^2 for an $N \times N$ image.

3. HIDDEN MARKOV TREE MODELS

The secondary properties of wavelet transforms give rise to joint wavelet statistics that are succinctly captured by the *wavelet-domain hidden Markov tree (HMT) model*, introduced by Crouse et al. (see [4] for a more detailed discussion).

The HMT models the nonGaussian marginal pdf $f(w_i)$ (**S1**) as a Gaussian mixture whose components are labeled by a *hidden state* $S_i \in S, L$. The S_i dictate from which of the two components in the mixture model w_i is drawn, and thus characterize (in the statistical sense) the magnitude of w_i . State S corresponds to a zero-mean, low-variance Gaussian, while state L corresponds to a zero-mean, high-variance Gaussian. If we let

$$g(x; \mu, \sigma^2) := \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(x - \mu)^2}{2\sigma^2} \right\} \quad (1)$$

denote the Gaussian pdf, then we can write

$$f(w_i | S_i = S) := g(w_i; 0, \sigma_{S,i}^2), \quad (2)$$

$$f(w_i | S_i = L) := g(w_i; 0, \sigma_{L,i}^2) \quad (3)$$

with $\sigma_L^2 > \sigma_S^2$. The marginal pdf $f(w_i)$ is the convex combination of the conditional densities

$$f(w_i) = p_i^S g(w_i; 0, \sigma_{S,i}^2) + p_i^L g(w_i; 0, \sigma_{L,i}^2), \quad (4)$$

with $p_i^S = 1 - p_i^L$. The p_i^S and p_i^L can be interpreted as the probability that w_i is small or large (in the statistical sense), respectively.

The persistence of wavelet coefficient magnitudes across scale (**S2**) is modeled by linking the hidden states in a Markov tree. The resulting dependency graph has a quadtree topology that mirrors the quadtree topology of the wavelet coefficients, see Fig. 1(b). Each subband is represented with its own quadtree; this assumes that the subbands are statistically independent.

Each parent→child state-to-state link has a corresponding transition matrix that quantifies statistically the degree of persistence of large/small coefficients:

$$A_i := \begin{bmatrix} p_i^{S \rightarrow S} & p_i^{S \rightarrow L} \\ p_i^{L \rightarrow S} & p_i^{L \rightarrow L} \end{bmatrix} \quad (5)$$

with $p_i^{S \rightarrow L} = 1 - p_i^{S \rightarrow S}$ and $p_i^{L \rightarrow S} = 1 - p_i^{L \rightarrow L}$.

Denote the parameters needed to specify a HMT model by the vector Θ . Members of Θ are the mixture variances for each state, $\sigma_{S,i}$ and $\sigma_{L,i}$, the transition probabilities $p_i^{S \rightarrow S}$ and $p_i^{L \rightarrow L}$, and a mass function for the hidden state of the root node, p_0^L . These parameters can be fit to a given set of training data using the Expectation-Maximization (EM)

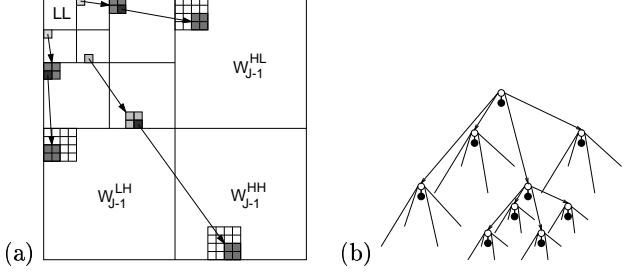


Figure 1: (a) Quadtree organization of the wavelet coefficients. The four children wavelet coefficients divide the spatial localization of the parent coefficient. (b) 2-D HMT model. Each black node is a wavelet coefficient; each white node is the corresponding hidden state. Links represent dependencies between states.

algorithm [4]. The training yields an approximate maximum likelihood estimate of the model parameters given the training data, yielding a good approximation of the joint density function $f(\mathbf{w})$ of the wavelet coefficients and thus $f(\mathbf{x})$.

In general, the HMT model for an $N \times N$ image has approximately $4n$ parameters, with $n := N^2$. In some applications, this large number of parameters could make the HMT model cumbersome. To accurately specify $4n$ parameters for an n -pixel image requires significant a priori information about the image. If this information is unavailable, we run the risk of over-fitting the model. Crouse et al. [4] reduces the total number of HMT parameters to approximately $4L$, with L the number of wavelet scales (typically 4–10), by tying within scale. While a significant reduction, a large quantity of a priori image information is still required to specify the parameters without over-fitting.

Often, the a priori image information takes the form of training data. Training algorithms such as the EM algorithm, especially for large data sets or data that have been severely corrupted by noise, can be computationally prohibitive. This makes the wavelet HMT model impractical for applications requiring computationally efficient processing. Furthermore, in many applications training data is unavailable. In such cases, an empirical Bayesian approach could be taken and a model fit to the data at hand. This is an effective approach if processing time is not an issue (see denoising examples in Fig. 2). However, if the observed data is severely corrupted (by noise, for example), then training may not be robust, and the model parameters will not characterize the joint image pdf accurately.

To address these problems, we will reduce the number of parameters in the HMT model. In doing this, the HMT model will become less accurate; two images that have different parameterizations with the general form of the HMT may have the same parameterization in the reduced-parameter model. What we gain is a reduction in complexity; less a priori information will be needed to specify the model parameters, and training will become more robust.

4. REDUCED-PARAMETER HMT IMAGE MODELS

Crouse et al. assumed that every image has a different HMT model, with the $4L$ parameters being specified by training on an observation [4]. In this section we take a different approach. We specify a new HMT model, called the *iHMT*, with a drastically reduced set of parameters (only 9), that incorporates properties common to all images in a class. The parameterization of the iHMT is based on the fact that for real-world images, the structure of the HMT is *self-similar* across scale [6].

Furthermore, we have found that many real-world images have similar iHMT parameters. By fixing one set of parameters, called the *universal HMT* (*uHMT*), we can take a strictly Bayesian approach to the estimation problem, eliminating the need for training altogether.

4.1. Tertiary Properties of the Wavelet Coefficients

The wavelet transforms of real-world images exhibit additional strong statistical properties in addition to the primary (**P1–P5**) and the secondary (**S1,S2**) properties. In designing our reduced-parameter HMT models, we will leverage the following *tertiary* properties of wavelet transform:

- T1. Exponential decay across scale:** The magnitudes of the wavelet coefficients of real-world images tend to decay exponentially across scale.
- T2. Stronger persistence at finer scales:** The persistence of large/small wavelet coefficient magnitudes becomes stronger at finer scales.

The exponential decay property (**T1**) stems from the overall smoothness and self-similarity of images. Roughly speaking, a typical real-world image consists of smooth regions separated by a finite number of discontinuities. This results in a $1/f$ -type spectral behavior, which leads to the exponential decay of the wavelet coefficients across scale [2].

We can obtain intuition behind property **T2** by considering the simple yet powerful image model of Cohen and D'Ales [5]. They model an image as piecewise smooth with a finite number of discontinuities. Consider a 1-D slice from such an image. Clearly it is also piecewise smooth with a finite number (say M) of discontinuities.

Since there are a finite number of discontinuities and the spatial resolution of the wavelet coefficients becomes finer as j increases (**P2**), there is some j_{crit} such that for all $j \geq j_{\text{crit}}$, each wavelet basis function has at most one discontinuity inside its spatial support. We call this condition *isolation of the edges*. Given no a priori information about the locations of the discontinuities, the fact that the spatial resolutions of the wavelet coefficients become finer exponentially implies that the probability that every edge is isolated goes to 1 exponentially.

By **P4**, for fine scales such that $j \gg j_{\text{crit}}$ there will be on the order of M wavelet coefficients that are “large” when compared to other coefficients at the same scale (exactly M if we are using the Haar wavelet). Each of these large coefficients will also have a large child, since the children wavelet basis functions simply divide up the spatial support of the parent. Each of the small coefficients’ children will

have small children, since there is no chance for any of them to encounter an edge.

In 2-D, the situation is similar except that instead of a discontinuities at points, we now have discontinuities along curves. At j_{crit} , all wavelet basis functions that have spatial support intersecting this curve will be “large.” Again, each of these coefficients will also have at least one large child, while the small coefficients will spawn only small children.

4.2. The iHMT model

Based on the tertiary properties of the wavelet transforms of real-world images, we can specify the HMT model parameters in a hyper-parametric form. The coefficient decay and the change in coefficient persistence are easily modeled by imposing structure on how the mixture variances and state transition probabilities change across scale. Because the tertiary properties are common to many real-world images, the resulting model describes the common overall behavior of real-world images in the wavelet domain.

We can easily model the exponential decay of wavelet coefficients (**T1**) through the mixture variances of the HMT model. Since the HMT mixture variances characterize the magnitudes of the wavelet coefficients, we will require that they decay exponentially across scale:

$$\sigma_{S,j}^2 = C_{\sigma_S} 2^{-j^{2\alpha_S}}, \quad (6)$$

$$\sigma_{L,j}^2 = C_{\sigma_L} 2^{-j^{2\alpha_L}}. \quad (7)$$

To have $\sigma_{S,j}^2 < \sigma_{L,j}^2$ for all scales, we require $\alpha_S \geq \alpha_L$. The result is an HMT for images with $1/f$ power spectra.

We will model the change in the degree of coefficient magnitude persistency by considering the way that the state transition probabilities change across scale.

Again, consider a 1-D signal consisting of smooth regions having M jump discontinuities. The isolation of edges at fine scales controls the persistency and novelty probabilities (and hence the form of the transition matrix) in the HMT. If each of the M edges in the 1-D slice is isolated then there is no opportunity for a novel large coefficient to come from a small parent; the only way a coefficient can be large is if its parent is large. Thus, $p_j^{S \rightarrow L} \rightarrow 0$ exponentially as $j \rightarrow \infty$. In other words, $p_j^{S \rightarrow S} \rightarrow 1$, since once a wavelet basis function lies over a smooth region, all of its children also lie over that smooth region. If a basis function lies over an edge, one and only one of its children will lie over the edge. This is an exact statement for the Haar basis functions, and a close approximation for longer wavelets. Therefore, the large wavelet coefficient gives rise to one large and one small wavelet coefficient and $p_j^{L \rightarrow L} \rightarrow \frac{1}{2}$. For a more in-depth analysis, see [6].

The edge isolation probability going to 1 exponentially means that the asymptotic values for persistency and novelty parameters are also approached exponentially. This gives a state transition matrix (see (5)) specified by four parameters:

$$A_j = \begin{bmatrix} 1 - C_{SS} 2^{-\gamma_{Sj}} & C_{SS} 2^{-\gamma_{Sj}} \\ \frac{1}{2} - C_{LL} 2^{-\gamma_{Lj}} & \frac{1}{2} + C_{LL} 2^{-\gamma_{Lj}} \end{bmatrix}. \quad (8)$$

The only parameter in the HMT not yet accounted for is the probability mass function on the hidden state value

of the root coefficients (just one number in our case, $p_{j_0}^L$, since the hidden state can only take two different values). Taking this parameter as is, we have reduced the number of parameters that specify the iHMT model to nine:

$$\Theta_i = \{\alpha_S, \alpha_L, C_{\sigma_S}, C_{\sigma_L}, \gamma_S, \gamma_L, C_{SS}, C_{LL}, p_{j_0}^L\}. \quad (9)$$

4.3. A “universal” iHMT: The uHMT

Now that we have an image model specified by a small set of parameters Θ_i , we must find a way of determining them. The first possibility would be to derive a constrained EM algorithm to give pseudo-MLE estimates of Θ_i given training data. Deriving the steps for this algorithm is difficult, and there is no guarantee that the training would be faster than in the unconstrained case.

Another possibility is to fix the parameters directly. This yields an iHMT model for a class of images, with each member in the class being treated as statistically equivalent. Although we clearly lose accuracy by viewing all images of interest as statistically equivalent, we totally eliminate the need for training. This saves us a tremendous amount of computation. For example, on a 512×512 image the EM algorithm can take anywhere from minutes to hours to converge on a typical workstation.

To see how much variation in iHMT parameters there is across grayscale, photograph-like images, we trained HMT models for a set of normalized images and examined their parameters. The variance and persistence decays were measured by fitting a line to the log of the variance vs. scale for each state. The decays were very similar for all of the images. Since the images were normalized, the range over which the variances decayed was similar as well. These observations lead us to believe that a specific, “universal” set of iHMT parameters can reasonably characterize photograph-like images. We call the HMT with this set of parameters the *uHMT* model.

The simplicity of the uHMT model also allows us to apply it in situations where the cost of a standard HMT would be prohibitive. For instance, we have developed a fast $O(n \log n)$ shift-invariant estimation scheme (discussed briefly in Section 5 and in detail in [6]) based on the uHMT parameters that delivers state-of-the-art performance (see Fig. 2).

5. APPLICATION TO IMAGE DENOISING

To demonstrate the effectiveness of the uHMT for modeling an image’s wavelet coefficients, we estimate an image submerged in additive white Gaussian noise. Translated into the wavelet domain, the problem is as follows:

$$\text{given } \mathbf{y} = \mathbf{w} + \mathbf{n}, \text{ estimate } \mathbf{w}, \quad (10)$$

where \mathbf{n} is a Gaussian random field whose components are independent and identically distributed with zero mean and known variance σ_n^2 .

Since we are viewing \mathbf{w} as a realization of a random field whose joint pdf is modeled by the HMT, we take a Bayesian approach to the estimation problem. The conditional density $f(\mathbf{y}|\mathbf{w})$ is given by the problem; it is an independent, Gaussian random field with mean \mathbf{w} . Using

the HMT model for $f(\mathbf{w})$, we can solve the Bayes equation for the posterior $f(\mathbf{w}|\mathbf{y})$.

To obtain the model parameters, Crouse et al. takes an empirical Bayesian approach [4]. The HMT parameters used to model $f(\mathbf{w}|\Theta)$ are first estimated from the observed noisy data \mathbf{y} and then “plugged-in” to the Bayes equation (after accounting for the noise).

For the Bayes estimator, we calculate the conditional mean of the posterior $f(\mathbf{w}|\mathbf{y}, \Theta)$ using the pointwise transformation

$$\hat{w}_i = E[w_i|\mathbf{y}, \Theta] = \sum_q p(S_i = q|\mathbf{y}, \Theta) \frac{\sigma_{q;i}^2}{\sigma_n^2 + \sigma_{q;i}^2} y_i \quad (11)$$

to obtain the minimum mean-square estimate (MMSE) of \mathbf{w} . Results using the empirical Bayesian HMT estimator, shown in Fig. 2(d) and Table 1, are competitive in both visual quality and PSNR to redundant wavelet shrinkage.

With the uHMT parameters, we have a prior on \mathbf{w} and the estimation problem can be approached from a purely Bayesian standpoint. Since we have eliminated training, the estimation algorithm is truly $O(n)$ and takes only a few seconds to run on a workstation.

To test this new Bayesian estimator, we denoised a set of images using the uHMT with parameters: $\alpha_L = \alpha_S = 5/4$, $C_{\sigma_S} = 2^7$, $C_{\sigma_L} = 2^{13}$, $\gamma_S = \gamma_L = 1$, $C_{SS} = C_{LL} = 32/5$, and $p_0^L = 1/2$. The results, given in Table 1 and Fig. 2(e), are almost identical to the more complicated empirical Bayes HMT approach, suggesting that we have lost almost nothing by totally eliminating training.

Image estimates obtained using an orthogonal wavelet transform frequently exhibit visual artifacts, usually in the form of ringing around edges. These artifacts can be combatted by averaging together estimates obtained from all different shifts of the image [7]. The resulting shift-invariant estimate is given by

$$\hat{x} = \text{Average}(\mathbf{S}_{-k,-m}(\mathbf{D}(\mathbf{S}_{k,m}(y))))_{0 \leq k, m \leq N-1} \quad (12)$$

where $\mathbf{S}_{k,m}(y) = y(s - k, t - m)$ is the 2-D shift operator and \mathbf{D} denotes the estimator (11). Implementing (12) directly would have computational complexity $O(n^2)$ and would thus be infeasible for large images. To streamline the algorithm, we must exploit the redundancies in the wavelet representations between different shifts of the image.

In the wavelet domain, each shift of the image corresponds to a different tree of wavelet coefficients. The wavelet coefficient trees for different shifts overlap, with common coefficients occupying entire subtrees. Averaging estimates for different shifts amounts to averaging the $p(S_i = q|\mathbf{y}, \Theta)$ for each tree in which w_i appears, and then using the result in (11) (we assume that Θ is the same for each shift of the image). The way in which the wavelet coefficient trees of different shifts overlap allows an $O(n \log n)$ shift-invariant denoising algorithm [6]. The results of Table 1 and Fig. 2(f) indicate that this denoising algorithm defines the new state-of-the-art: in general, we gain a 1–1.5 dB gain over thresholding with the redundant wavelet transform [7, 8].

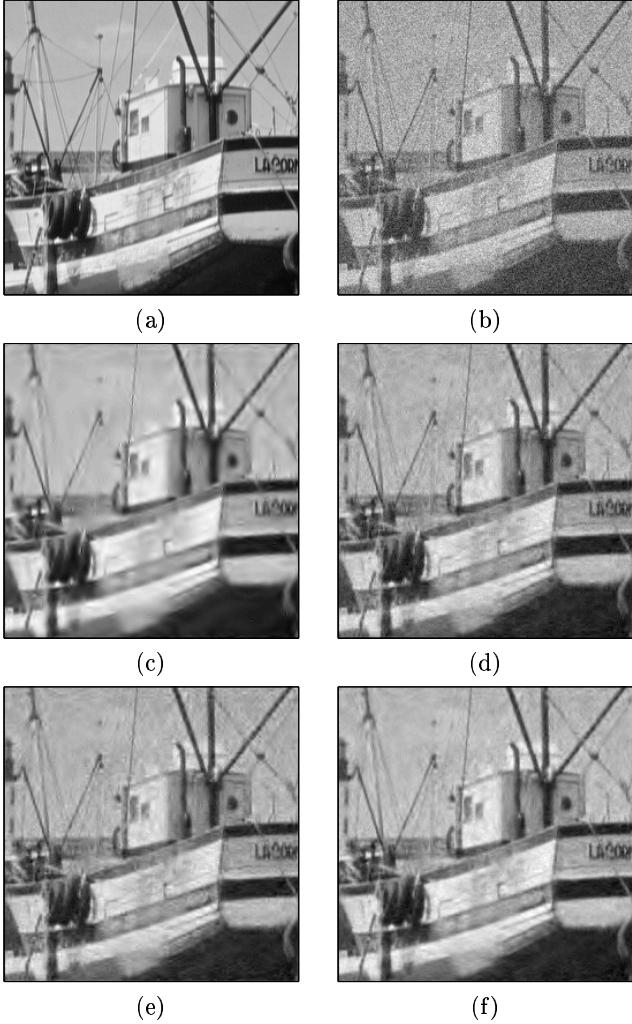


Figure 2: (a) Original 256×256 “Boats” image; (b) Noisy boats image, with $\sigma_n = 0.1$, $PSNR = 20$ dB. Boats image denoised using (c) redundant hard-thresholding using empirical best threshold [8], $PSNR = 26.3$ dB; (d) empirical Bayesian HMT estimator [4] $PSNR = 26.5$ dB; (e) uHMT Bayesian estimator, $PSNR = 26.4$ dB; (f) shift-invariant uHMT estimator, $PSNR = 27.4$ dB.

6. CONCLUSIONS

Hidden Markov Trees capture the primary aspects of image structure in the wavelet domain. In this paper, we have shown that additional image structure can be exploited by constraining the HMT parameters to have a certain form. The resulting model, the iHMT, has only 9 parameters.

A set of “universal” parameters arises naturally from the form of the iHMT. These nine numbers completely specify a model for a large class of real-world images, eliminating any need for training in the estimation algorithm without compromising denoising performance. Having the model fully specified facilitates the implementation of a shift-invariant estimation algorithm which offers state-of-the-art performance.

Table 1: Image estimation results for 256×256 images corrupted with additive white Gaussian noise of $\sigma_n = 0.1$. Entries are the peak signal to noise ratio (PSNR), $PSNR := -20 \log_{10}(\|\hat{x} - x\|_2/N)$. Pixel intensity values were normalized between 0 and 1. All results use the Daubechies-8 wavelet. “R-HMT” is the shift-invariant estimator; “uHMT” uses the “universal” parameters presented in Section 5; “E-HMT” uses the empirical Bayesian estimator of [4]; “R-Thr” uses a hard thresholded redundant wavelet transform using the thresholds in [8]

Image	R-HMT	uHMT	E-HMT	R-Thr
Baby	29.6	28.9	29.2	29.5
Birthday	26.4	25.8	25.8	25.3
Boats	27.4	26.4	26.5	26.3
Bridge	25.3	24.6	25.0	23.7
Buck	29.6	28.4	28.6	29.7
Building	26.6	25.9	26.3	25.8
Camera	27.0	26.2	26.4	26.3
Clown	27.8	26.8	26.8	26.5
Fruit	29.7	28.5	28.6	29.0
Kgirl	29.3	28.3	28.3	28.4
Lenna	27.6	26.7	26.7	26.3

7. REFERENCES

- [1] S. Geman and D. Geman, “Stochastic relaxation, Gibbs distribution, and the Bayesian restoration of images,” *IEEE Trans. on Pattern Anal. Machine Intell.*, vol. 6, pp. 721–741, 1984.
- [2] S. Mallat, *A Wavelet Tour of Signal Processing*. San Diego: Academic Press, 1998.
- [3] M. Vetterli and J. Kovačević, *Wavelets and Subband Coding*. Englewood Cliffs NJ: Prentice Hall, 1995.
- [4] M. S. Crouse, R. D. Nowak, and R. G. Baraniuk, “Wavelet-based statistical signal processing using hidden Markov models,” *IEEE Trans. Signal Proc.*, vol. 46, pp. 886–902, April 1998.
- [5] A. Cohen and J. P. D’Ales, “Nonlinear approximation of random functions,” *SIAM J. Appl. Math.*, vol. 57, April 1997.
- [6] J. K. Romberg, H. Choi, and R. G. Baraniuk, “Bayesian tree-structured image modeling using wavelet-domain hidden Markov models.” Preprint, available at www.dsp.rice.edu.
- [7] R. Coifman and D. Donoho, “Translation-invariant de-noising,” in *Wavelets and Statistics*, Lecture Notes in Statistics, Springer-Verlag, 1995.
- [8] M. Lang, H. Guo, J. E. Odegard, C. S. Burrus, and R. O. Wells, “Nonlinear processing of a shift invariant DWT for noise reduction,” in *Proceedings of SPIE*, 1995.