

Spatially Adaptive Wavelet Thresholding with Context Modeling for Image Denoising

S. Grace Chang, *Student Member, IEEE*, Bin Yu, *Senior Member, IEEE*, and Martin Vetterli, *Fellow, IEEE*

Abstract—The method of wavelet thresholding for removing noise, or denoising, has been researched extensively due to its effectiveness and simplicity. Much of the literature has focused on developing the best uniform threshold or best basis selection. However, not much has been done to make the threshold values adaptive to the spatially changing statistics of images. Such adaptivity can improve the wavelet thresholding performance because it allows additional local information of the image (such as the identification of smooth or edge regions) to be incorporated into the algorithm. This work proposes a spatially adaptive wavelet thresholding method based on context modeling, a common technique used in image compression to adapt the coder to changing image characteristics. Each wavelet coefficient is modeled as a random variable of a generalized Gaussian distribution with an unknown parameter. Context modeling is used to estimate the parameter for each coefficient, which is then used to adapt the thresholding strategy. This spatially adaptive thresholding is extended to the overcomplete wavelet expansion, which yields better results than the orthogonal transform. Experimental results show that spatially adaptive wavelet thresholding yields significantly superior image quality and lower MSE than the best uniform thresholding with the original image assumed known.

Index Terms—Adaptive method, context modeling, image denoising, image restoration, wavelet thresholding.

I. INTRODUCTION

IN THIS paper, we address the classical problem of removing additive noise from a corrupted image, or *denoising*. In recent years there has been a plethora of work on using *wavelet thresholding* [6] for denoising in both the signal processing and statistics community, due to its effectiveness and simplicity. In its most basic form, this technique denoises in the orthogonal wavelet domain, where each coefficient is *thresholded* by comparing against a threshold; if the coefficient is smaller than the threshold, it is set to zero, otherwise it is kept or modified. The intuition is that because the wavelet transform is good at energy

compaction, small coefficients are more likely due to noise, and large coefficients due to important signal features (such as edges). The threshold thus acts as an oracle deciding whether or not to keep the coefficients. Most of the literature thus far has concentrated on developing threshold selection methods, with the threshold being uniform or at best one threshold for each subband. Very little has been done on developing thresholds that are adaptive to different spatial characteristics. Other works investigate the choice of wavelet basis or expansion for the thresholding framework. One particularly interesting result is that (uniform) thresholding in a shift-invariant expansion (dubbed *translation-invariant (TI) denoising* by Coifman and Donoho [4]) eliminates some of the unpleasant artifacts introduced by the modification of the orthogonal wavelet expansion coefficients. In this paper, we use the wisdom that thresholding in a shift-invariant, overcomplete representation outperforms the orthogonal basis, and also investigate an issue that has not been explored, namely, the spatial adaptivity of the threshold value.

To motivate spatially adaptive thresholding, consider the example in Fig. 1, where a square pulse function has been corrupted by additive noise, and the goal is to recover the original function. The wavelet coefficients of the original and the noisy function are displayed in Fig. 1(a) and (b), respectively. The noisy coefficients are (soft) thresholded by a single threshold in Fig. 1(c), and one can see that, especially in the finest scale, there are some coefficients corresponding to noise which have not been set to zero, and that some of these noisy coefficients are larger in magnitude than those coefficients corresponding to the signal. Thus, with a uniform threshold, it may not be feasible to have both the benefits of keeping the important signal features and killing the noisy coefficients. On the other hand, one can reap both benefits with adaptive thresholds by choosing the threshold value to be very small in the regions of the peaks due to the step function, and large otherwise [see Fig. 1(d)].¹ Fig. 1(e) compares the reconstruction from both methods, and adaptive thresholding yields a more accurate reconstruction and retains better the sharp edges. Thus the question becomes, how one distinguishes between the coefficients that are mainly due to signal and those mainly due to noise. Also, how should the thresholds be adjusted? These are the questions that we will answer in this paper with our proposed algorithm.

Most natural images have changing characteristics, since they typically consist of regions of smoothness and sharp transitions. These regions of varying characteristics can be well differentiated in the wavelet domain, as can be seen in the wavelet decom-

Manuscript received August 10, 1998; revised March 21, 2000. This work was supported in part by the NSF Graduate Fellowship and the University of California Dissertation Fellowship to S. G. Chang; ARO Grant DAAH04-94-G-0232 and NSF Grant DMS-9322817 to B. Yu; and NSF Grant MIP-93-213002 and Swiss NSF Grant 20-52347.97 to M. Vetterli. Part of this work was presented at the IEEE International Conference on Image Processing, Chicago, IL, October 1998. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Scott T. Acton.

S. G. Chang was with the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley, CA 94720 USA. She is now with Hewlett-Packard Company, Grenoble, France (e-mail: grchang@yahoo.com).

B. Yu is with the Department of Statistics, University of California, Berkeley, CA 94720 USA (e-mail: binyu@stat.berkeley.edu)

M. Vetterli is with the Laboratory of Audiovisual Communications, Swiss Federal Institute of Technology (EPFL), Lausanne, Switzerland and also with the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley, CA 94720 USA.

Publisher Item Identifier S 1057-7149(00)06911-6.

¹For the sake of illustrating the effectiveness of varying thresholds, the regions of the true peaks are assumed to be known in this example.

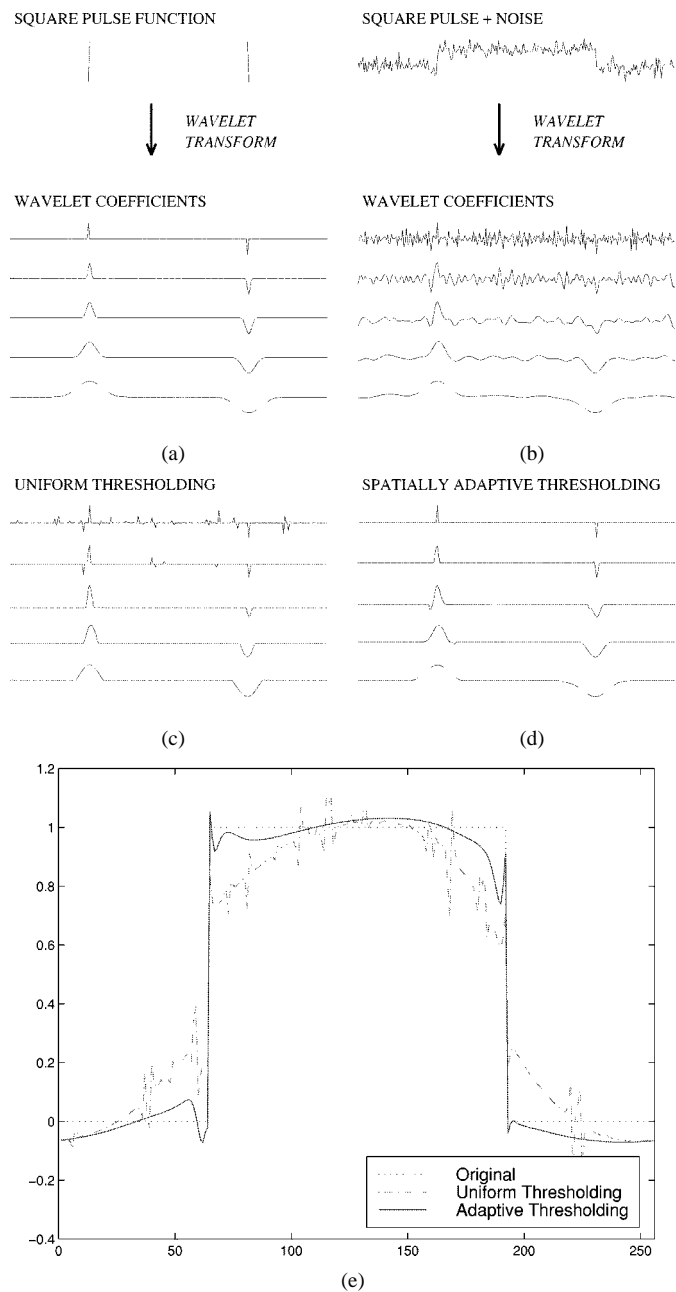


Fig. 1. Motivation for adaptive thresholds: (a) and (b) show a square pulse function and its corrupted version, respectively, along with their wavelet decomposition of four scales. The wavelet coefficients are thresholded by a uniform threshold in (c) and spatially adaptive thresholds in (d). The original and the reconstructions from (c) and (d) are shown in (e).

position of the *Lena* image in Fig. 2. One observes areas of high and low energy (or large and small coefficient magnitude), represented by white and black pixels, respectively. High energy areas correspond to signal features of sharp variation such as edges and textures; low energy areas correspond to smooth regions. When noise is added, it tends to increase the magnitude of the wavelet coefficient, on average. Specifically, in smooth regions, one expects the coefficients to be dominated by noise, thus most of these coefficients should be removed, especially since noise is highly visible here. In regions of sharp transition, the coefficients have a lot of energy due to the signal, and some due to noise (which is not as visible in these regions), thus they

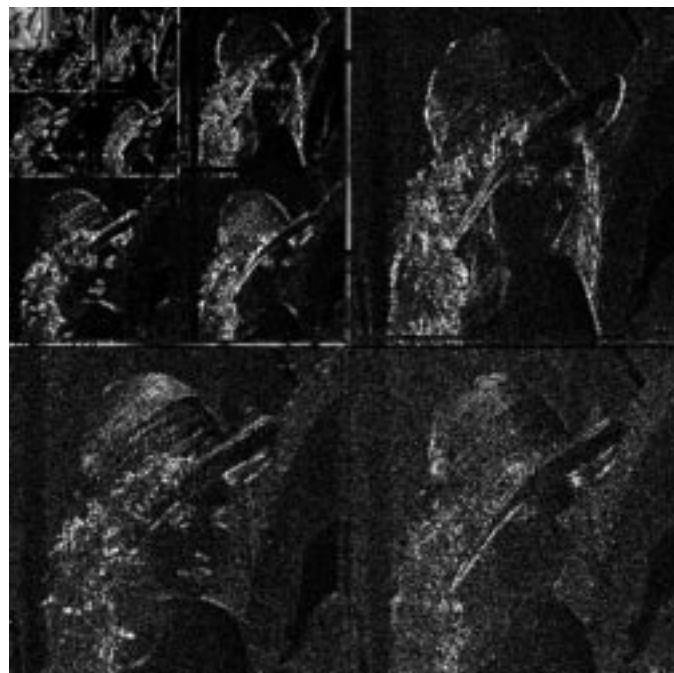


Fig. 2. Four level wavelet decomposition of *Lena*. White pixels indicate large magnitude coefficients, and black signifies small magnitude.

should be kept, or at most modified only a little, to ensure that most of the signal details are retained. Thus, the idea is to distinguish between the low and high energy regions, and modify the coefficients using a spatially adaptive thresholding strategy.

To accomplish spatial adaptive thresholding, we model each wavelet coefficient as a realization from a given probability distribution, whose parameter is to be estimated. This parameter in turn is used to find the appropriate threshold. It has long been accepted in the subband coding community that for a large class of images, the coefficients in each subband form a distribution well described by a generalized Gaussian distribution (GGD) [15]. Instead of using one parameter for each subband level, several wavelet-based image coders have achieved better performances by modeling the wavelet coefficients as a mixture of GGD with unknown slowly spatially varying parameters [8], [16]. The estimation of the parameter for a given coefficient is conditioned on a function of its neighboring coefficients, a method called *context-modeling* frequently used in compression for differentiating pixels of varied characteristics and adapting the coder. Context modeling allows one to group pixels of similar nature but not necessarily spatially adjacent, and to gather statistical information from these pixels. Now, given that one can estimate the parameter for each coefficient, the next step is to use them to calculate the threshold. Our work in [3] found that when the signal coefficients are modeled as GGD random variables and the noise as Gaussian, the threshold $T_B = \sigma_n^2 / \sigma_X$ is a good approximation to the optimal threshold which minimizes the mean squared error of the soft-thresholding estimator, where σ_n^2 is the noise power, and σ_X is the standard deviation of the signal. The simplicity of this threshold makes it easy to achieve spatial adaptivity—one uses context modeling to quantify the local characteristic in σ_X , which in turn yields a threshold T_B adaptive on a pixel-by-pixel manner.

Our proposed algorithm is based on using adaptive thresholding in the overcomplete wavelet expansion (specifically, nonsubsampling expansion). It outperforms both using only adaptive thresholding in the orthogonal expansion or using only uniform thresholding in the overcomplete expansion like the TI denoising. That is, *by combining spatially adaptive thresholding and overcomplete expansion, we achieve results which are significantly superior than either method alone*. First, the adaptive threshold selection is effective at removing noise in smooth regions, while not disturbing too much the edge and texture regions. Second, thresholding in the orthogonal expansion has been observed to produce Gibbs-like edge artifacts. Thresholding in the overcomplete expansion has the interpretation of averaging circularly-shifted, denoised versions of the signal, thus providing an additional smoothing which attenuates these unpleasant artifacts [4].

The organization of this paper is as follows. Section II introduces the threshold selection method when there is only one class of Generalized Gaussian distributed random variables corrupted by independent additive Gaussian noise. Because this threshold selection is based on *iid* noise assumption, the discussion will first be set in the orthogonal wavelet transform. The method of context modeling is then introduced to allow coefficients be modeled as random variables of different parameters, and these parameters are used to make the threshold spatially adaptive. The explanation of our proposed algorithm is completed by extending this adaptive method in the orthogonal expansion to the overcomplete expansion. In Section III, we will compare the spatially adaptive results with those from the best uniform thresholding strategy (in the mean squared error sense, and based on knowing the original image), in both the orthogonal and overcomplete expansion, and also with a state-of-the-art denoising method [10]. Results will show that the combination of using spatially adaptive thresholding and overcomplete expansion yields significantly better results in both visual quality and mean squared error.

II. ADAPTIVE ALGORITHM

The adaptive algorithm will be developed in the following manner. First, the concept of wavelet thresholding in the orthogonal wavelet transform is introduced, and coefficients in each subband are modeled as realizations of one class of GGD (with unknown parameter), corrupted by additive white Gaussian noise. We then describe the threshold selection method developed in [3] which achieves near-optimal soft-thresholding under these distributions. To make this thresholding approach spatially adaptive, *each coefficient* (rather than each subband) is modeled as a GGD random variable with a different unknown parameter estimated on a pixel level using context modeling. This spatial mixture of distributions allows the image characteristics to be quantified locally via the distribution parameters, which are then used to adjust the threshold for each coefficient. Last, since the aforementioned algorithm is developed in the orthogonal expansion where the coefficients are uncorrelated, the algorithm will need to be modified to extend to the overcomplete expansion where coefficients are

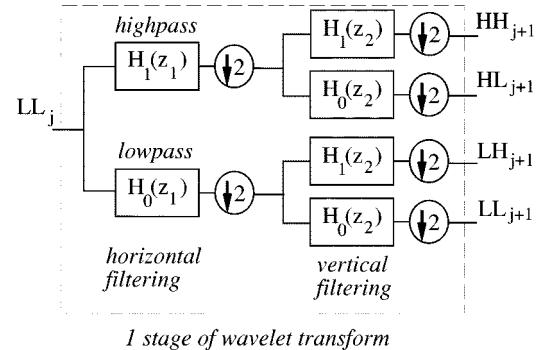


Fig. 3. One stage of the 2-D wavelet transform.

correlated. Several alternative approaches will be discussed in terms of their applicability to our adaptive algorithm.

A. Coefficient Modeling and Threshold Selection

Let the original image be $\{f[i, j], i, j = 1, \dots, N\}$, where N is some integer power of 2. The image has been corrupted by additive noise and one observes

$$g[i, j] = f[i, j] + \varepsilon[i, j], \quad j = 1, \dots, N \quad (1)$$

where $\{\varepsilon[i, j]\}$ are independent and identically distributed (*iid*) as normal $N(0, \sigma_n^2)$ and independent of $\{f[i, j]\}$. The goal is to remove the noise, or “denoise” $\{g[i, j]\}$, and to obtain an estimate $\{\hat{f}[i, j]\}$ of $\{f[i, j]\}$ which minimizes the mean squared error (MSE)

$$\text{MSE}(\hat{f}) = \frac{1}{N^2} \sum_{i, j=1}^N (\hat{f}[i, j] - f[i, j])^2. \quad (2)$$

To accomplish wavelet thresholding for denoising, the observations $\{g[i, j]\}$ are first transformed into the wavelet domain. The necessary notations for the wavelet transform will be introduced here, and the readers are referred to references such as [9], [14] for more details. The two-dimensional (2-D) discrete orthogonal wavelet transform (DWT) can be implemented as a critically sampled octave-band filter bank, where separable filtering is used (Fig. 3). Let $\mathbf{g} = \{g[i, j]\}_{i, j}$, $\mathbf{f} = \{f[i, j]\}_{i, j}$, $\boldsymbol{\varepsilon} = \{\varepsilon[i, j]\}_{i, j}$, that is, the boldfaced letters will denote the matrix representation of the signals under consideration. Let $\mathbf{Y} = \mathcal{W}\mathbf{g}$ denote the matrix of wavelet coefficients of \mathbf{g} , where \mathcal{W} is the 2-D dyadic orthogonal wavelet transform operator, and similarly $\mathbf{X} = \mathcal{W}\mathbf{f}$ and $\mathbf{V} = \mathcal{W}\boldsymbol{\varepsilon}$. It is often convenient to cluster these coefficients into groups or *subbands* of different scales and orientations as in Fig. 4, where, for example, the label HL_1 refers to those coefficients at the first scale of decomposition which are the output of the highpass filter in the horizontal direction and the lowpass filter in the vertical direction. The subbands $HH_k, HL_k, LH_k, k = 1, 2, \dots, J$ are called the *details*, where k is the *scale*, with J being the largest (or coarsest) scale in the decomposition, and a subband at scale k has size $N/2^k \times N/2^k$. The subband LL_J is the *low resolution residual*, and J is typically chosen large enough such that $N/2^J \ll N$ and $N/2^J > 1$. Let $Y^{(s, o)}[i, j], i, j =$

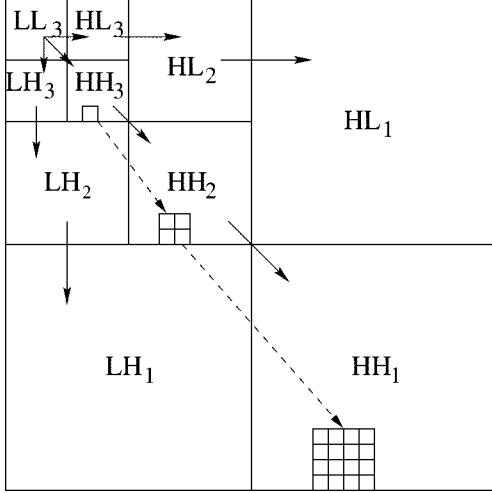


Fig. 4. Subbands of the orthogonal DWT. The arrows show the relationship from parent to child.

$1, \dots, N/2^s$, denote the wavelet coefficients of $\{g[i, j]\}$ at a particular scale s and orientation o , where $s = 1, 2, \dots, J$ and $o \in \{HL, LH, HH, LL\}$.

The method of wavelet thresholding denoising filters each coefficient Y_{ij} from the detail subbands with a threshold function (to be explained shortly) to obtain \hat{X}_{ij} . The denoised estimate is then $\hat{\mathbf{f}} = \mathcal{W}^{-1}\hat{\mathbf{X}}$, where \mathcal{W}^{-1} is the inverse wavelet transform operator.

It has been observed that for a large class of images, the coefficients from each subband (except LL) form a symmetric distribution that is sharply peaked at zero [9], [15], well described by the zero-mean GGD

$$GG_{\beta, \sigma_X}(x) = C(\beta, \sigma_X) e^{-(\alpha(\beta, \sigma_X)|x|)^\beta}, \quad -\infty < x < \infty, \sigma_X > 0, \beta > 0 \quad (3)$$

where

$$\alpha(\beta, \sigma_X) = \sigma_X^{-1} \left[\frac{\Gamma\left(\frac{3}{\beta}\right)}{\Gamma\left(\frac{1}{\beta}\right)} \right]^{1/2}, \quad C(\beta, \sigma_X) = \frac{\beta \alpha(\beta, \sigma_X)}{2\Gamma\left(\frac{1}{\beta}\right)} \quad (4)$$

and $\Gamma(t) = \int_0^\infty e^{-u} u^{t-1} du$ is the gamma function. The parameter σ_X is the standard deviation and determines the spread of the density function, and the parameter β is called the *shape parameter*. Well-known special cases of the GGD density function include the Gaussian distribution with $\beta = 2$, and the Laplacian with $\beta = 1$. Let $\{X^{(s,o)}[i, j]\}$ and $\{V^{(s,o)}[i, j]\}$ denote the wavelet coefficients of the original signal $\{f[i, j]\}$ and the noise $\{\varepsilon[i, j]\}$, respectively. For each subband, the signal coefficients $\{X^{(s,o)}[i, j]\}$ are modeled as independent samples of distribution $p_X(x) = GG_{\beta, \sigma_X}(x)$, and since the wavelet transform is orthogonal, the noise coefficients are independent samples of the Gaussian distribution $p_V(v) = \phi(v, \sigma_n^2) = 1/\sqrt{2\pi\sigma_n^2} \exp\{-v^2/2\sigma_n^2\}$. Let the estimator be restricted to the *soft-threshold* estimator of the form $\hat{X}^{(s,o)}[i, j] = \eta_T(Y^{(s,o)}[i, j])$, where

$\eta_T(x) = \text{sgn}(x) \cdot \max(|x| - T, 0)$ is the soft-thresholding function, and T is called the *threshold*. The optimal threshold T^* is defined to be the argument which minimizes the expected squared error

$$T^* = \arg \min_T E_{Y|X, X} (\eta_T(Y) - X)^2 \quad (5)$$

where $Y|X \sim \phi(y - x, \sigma_n^2)$ and $X \sim GG_{\beta, \sigma_X}(x)$. The soft-thresholding function is chosen over another popular choice, the *hard-thresholding* function, $\psi_T(x) = x \cdot \mathbf{1}_{\{|x| > T\}}$, because it yields denoised images with better visual quality. In [3], it was found that at least for the range of $\beta \in [0.5, 4]$, T^* can be well approximated by

$$T_B = \frac{\sigma_n^2}{\sigma_X} \quad (6)$$

with at most 5% difference from the minimum expected squared error. Because the threshold T_B depends only on the standard deviation σ_X and not on the shape parameter β , it may not yield a good approximation for other values of β than the range tested here, and the threshold may need to be modified to incorporate β . However, in practice it has been observed in several work [13], as well as in our own experience, that β typically falls within the range $[0.5, 1]$, well within the range of β tested here, thus the simple form of the threshold T_B is appropriate for our purpose. The curve of expected squared error is very flat near the optimal threshold T^* , implying that the error is not very sensitive to a slight perturbation near T^* .

The threshold $T_B = \sigma_n^2/\sigma_X$ is not only nearly optimal but also has an intuitive appeal. For such a choice, the normalized threshold T_B/σ_n is inversely proportional to σ_X , the standard deviation of X , and proportional to σ_n , the noise standard deviation. When $\sigma_n/\sigma_X \ll 1$, the signal is much stronger than the noise, thus T_B/σ_n is chosen to be small in order to preserve most of the signal and remove some of the noise; vice versa, when $\sigma_n/\sigma_X \gg 1$, the noise dominates and the normalized threshold is chosen to be large to remove the noise which has overwhelmed the signal.

The proposed threshold T_B can easily be adjusted to the signal and noise energy as reflected in σ_X and σ_n . By estimating the parameter σ_X for each subband, we have an *uniform* threshold T_B adaptive on a subband-level. Better denoising performance can be achieved by using spatially adaptive thresholds, whose derivation will be described in the following section.

B. Context Modeling for Spatial Adaptivity

To achieve a spatially adaptive thresholding strategy, the wavelet coefficients are modeled as components in a discrete random field, with a collection of independent zero-mean GGD random variables whose parameters β and σ_X are spatially varying. As discussed previously, mainly the parameter σ_X is of interest since the threshold $T_B = \sigma_n^2/\sigma_X$ depends on it, and β is assumed to be in the range for which this threshold is appropriate. The parameter σ_X needs to be estimated for each coefficient to make the threshold T_B spatially adaptive. This can be accomplished by *context modeling*, an idea used frequently in image compression for adapting the coder to

changing image characteristics. That is, the statistical model for a given coefficient is conditioned on a function of its neighbors. Several model-based coders have utilized information from causal quantized neighbors to determine the context model and the model parameters for each coefficient (for example, [8], [16]). The most relevant work is the wavelet-based compression scheme in [16], where context modeling was used to classify coefficients into several classes of Laplacian distributions with different values of σ_X . The conditioning was based on the weighted average of the magnitude of quantized coefficients in a causal neighborhood, and each class was formed by clustering coefficients whose associated weighted averages fall within a specified range. The distribution parameter is estimated from the coefficients for each class, which is then used to adapt the coder. Since the parameter and the description of each class need to be sent as overhead, only four classes were used in [16].

From [16], we adopt the idea of clustering pixels with similar context for parameter estimation. For the denoising problem, it is not necessary to explicitly cluster the pixels into a discrete number of classes in order to conserve bits. Thus, one has the luxury of estimating the parameters for each coefficient via, say, a moving window, resulting in virtually an infinite mixture of distributions. The estimation method used here will be for the GGD parameters, and our context model generalizes that of [16] to be more flexible and advantageous for the denoising task.

Consider one particular subband with M^2 coefficients, $\{Y^{(s,o)}[i,j]\}$. To simplify notation, the superscript (s,o) is dropped and will be used only when necessary for clarity. Each coefficient $Y[i,j]$ is modeled as a random variable whose variance can be estimated as follows. Consider a neighborhood of $Y[i,j]$, and the absolute value of its p elements are placed in a $p \times 1$ vector \mathbf{u}_{ij} . One possible choice is the eight nearest neighbors of $Y[i,j]$ in the same subband, plus its parent coefficient $Y^{(s+1,o)}[\lceil i/2 \rceil, \lceil j/2 \rceil]$ (see Fig. 4 for the definition of parent-child relationship). To characterize the activity level of the current pixel, we calculate the context as a weighted average of the absolute value of the neighbors

$$Z[i,j] = \mathbf{w}^t \mathbf{u}_{ij}. \quad (7)$$

The weight \mathbf{w} is found by using the least squares estimate, that is,

$$\mathbf{w}_{LS} = \arg \min_{\mathbf{w}} \sum_{i,j} (|Y[i,j]| - \mathbf{w}^t \mathbf{u}_{ij})^2 \quad (8)$$

$$= (\mathbf{U}^t \mathbf{U})^{-1} \mathbf{U}^t |\mathbf{Y}| \quad (9)$$

where \mathbf{U} is a $M^2 \times p$ matrix with each row being \mathbf{u}_{ij}^t , for all i,j , and \mathbf{Y} is the $M^2 \times 1$ vector containing all coefficients $Y[i,j]$. Notice that the absolute values of the neighbors rather than their original values are used in the averaging. This is because orthogonal wavelet coefficients are uncorrelated, and thus an average of the neighbors does not yield much information about the coefficient of interest. However, the absolute value or the squared values of neighboring coefficients are correlated [12], and therefore their averages are useful in collecting information about other coefficients in the near vicinity.

The variance of the random variable $Y[i,j]$ is estimated from other coefficients whose context variable

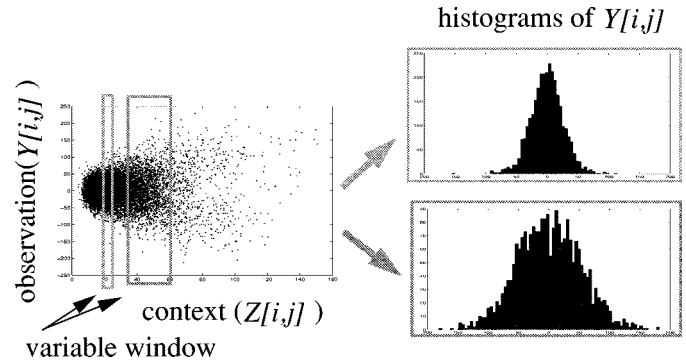


Fig. 5. Sample plot of $\{Z[i,j], Y[i,j]\}$, where $Y[i,j]$ is the noisy wavelet coefficient, and $Z[i,j]$ is its context. A collection of $Y[i,j]$ with small values of $Z[i,j]$ have a smaller spread than those with large values of $Z[i,j]$, suggesting that context modeling provides a good variability estimate of $Y[i,j]$.

are close in value to $Z[i,j]$. To develop an intuition for this, it is helpful to examine Fig. 5, which plots the pairs $\{Z[i,j], Y[i,j]\}$, $i,j = 1, \dots, M$. The points are clustered within a cone shape centered at origin. Taking an interval of small valued $Z[i,j]$, the associated coefficients $\{Y[i,j]\}$ have a small spread; on the other hand, an interval of large valued $Z[i,j]$ has corresponding $\{Y[i,j]\}$ with a larger spread (the intervals are of different widths to capture the same number of points). This suggests that the context provides a good indication of local variability. Thus, for a given coefficient $Y[i_0, j_0]$, an interval is placed around $Z[i_0, j_0]$, and the variance of $Y[i_0, j_0]$ is estimated from the points $Y[i,j]$ whose context $Z[i,j]$ falls within this window. In particular, we take L closest points (in value) above $Z[i_0, j_0]$ and L closest points below, resulting in a total of $2L + 1$ points. We choose $L = \max(50, 0.02 \cdot M^2)$ to ensure that enough points are used to estimate the variance, but not too many points to destroy the locality of the window. Different choices of L around this value yields similar results, though too small (e.g., 10 or less) or too large (e.g., close to $M^2/2$) values of L worsen the performance significantly. Note that this is a moving window rather than the fixed classes in [16], and thus allows a continuous range of estimate values. Let \mathcal{B}_{i_0, j_0} denote the set of points $\{Y[i,j]\}$ whose context falls in the moving window. The estimate of the variance $\sigma_X^2[i_0, j_0]$ is then

$$\hat{\sigma}_X^2[i_0, j_0] = \max \left(\frac{1}{2L+1} \sum_{[k,\ell] \in \mathcal{B}_{i_0, j_0}} Y[k,\ell]^2 - \sigma_n^2, 0 \right). \quad (10)$$

The term σ_n^2 needs to be subtracted because $\{Y[i,j]\}$ are the noisy observations, and the noise is independent of the signal, with variance σ_n^2 . The threshold at location $[i_0, j_0]$ is then

$$T_B[i_0, j_0] = \frac{\sigma_n^2}{\hat{\sigma}_X^2[i_0, j_0]}. \quad (11)$$

Calculating the threshold $T_B[i,j]$ for every location $[i,j]$ yields a spatially adaptive threshold. In the implementation, the context $\{Z[i,j]\}$ are first sorted, and a moving window is placed over them, so the set \mathcal{B}_{ij} and the variance estimate $\hat{\sigma}_X^2[i,j]$ can

be updated efficiently. To update the summation in (10) only requires one addition and subtraction for the $Y[k, \ell]^2$ term, and one multiplication for the constant $1/(2L + 1)$. Thus, the arithmetic complexity incurred by computing (10) is of order M^2 , where M^2 is the number of points in the subband.

The above threshold estimation method is repeated for each subband separately, because the subbands exhibit significantly different characteristics. Up to now we have not discussed how to estimate the noise variance σ_n^2 . In some practical cases, it is possible to measure σ_n^2 based on information other than the corrupted observation. If this is not the case, as is here, we estimate it by using the robust median estimator in the highest subband of the wavelet transform

$$\hat{\sigma}_n = \text{Median}(|Y[i, j]|)/0.6745, \quad Y[i, j] \in \text{subband } HH_1 \quad (12)$$

also used in [3], [4], [6].

C. Thresholding in Overcomplete Expansion

Thresholding in the orthogonal wavelet domain has been observed to produce significantly noticeable artifacts such as Gibbs-like ringing around edges and specks in smooth regions. To ameliorate this unpleasant phenomenon, Coifman and Donoho [4] proposed the *translation-invariant (TI) denoising*. The discussion in [4] is one-dimensional (1-D), but we explain it in 2-D here. Let $Shift_{k, \ell}[g]$ denote the operation of circularly shifting the input image g by k indices in the vertical direction and ℓ indices in the horizontal, and let $Unshift_{k, \ell}[g]$ be a similar operation but in the opposite direction. Also, let $Denoise[g, T]$ denote the operation of taking the DWT of the input image g , threshold it with a chosen uniform threshold T , then transform it back to the space domain. Then TI denoising yields an output which is the average of the thresholded copies over all possible shifts: $\hat{f} = (1/N^2) \sum_{k, \ell=0}^{N-1} Unshift_{k, \ell}[Denoise[Shift_{k, \ell}[g], T]]$. The rationale is that since the orthogonal wavelet transform is a time-varying transform and thresholding the coefficients produces ringing-like phenomena, thresholding a shifted input would produce ringing at different locations, and averaging over all different shifts would yield an output with more attenuated artifacts than a single copy alone. TI denoising can be shown to be equivalent to thresholding in the shift-invariant, overcomplete representation implemented by the *nonsampled filter bank* as will be described below, up to some scaling in the thresholds. It has been shown to reduce the ringing artifacts and the specks. Thus, we proceed to extend our spatial adaptive algorithm to the nonsampled expansion.

The adaptive algorithm in the orthogonal basis described above can easily be extended to the overcomplete basis. Now consider the same orthogonal filters but used in a filter bank without down-samplers (see Fig. 6 for the 1-D case, and the 2-D case is achieved by simply extending the 1D filtering to separable filtering; refer to [14] for more detail on nonsampled filter banks). The filters are renormalized by $1/\sqrt{2}$ so that coefficient energy stays the same. This decomposition is a redundant representation, and there are correlations between the decomposition coefficients. For example, at the first level of decomposition, the odd and even

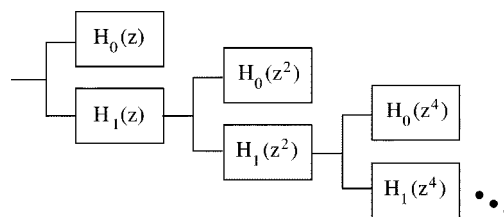


Fig. 6. One-dimensional nonsampled filter bank. One-dimensional filtering is extended to 2-D by separable filtering.

coefficients (in each direction) are correlated. Thus, we can separate the coefficients into four sets of uncorrelated coefficients, namely, $\{Y[2i, 2j]\}$, $\{Y[2i, 2j + 1]\}$, $\{Y[2i + 1, 2j]\}$ and $\{Y[2i + 1, 2j + 1]\}$. For the s th level decomposition, the coefficients can be separated into 2^{2s} sets, each containing uncorrelated coefficients, and they are $\{Y[2^s i + k_1, 2^s j + k_2]\}_{i, j}$, $k_1, k_2 = 0, 1, \dots, 2^s - 1$. Since each set contains uncorrelated coefficients, the noise are also *iid* within each set as well, and thus the adaptive algorithm can be used for each set of coefficients. This approach lets us still use the independent noise assumption and circumvent the issue of denoising correlated signal coefficients with correlated noise, which is not an easy task. That is, if the coefficients are correlated, then one can conceivably do better than thresholding each coefficient independently; one could look at the numerous correlated coefficients and do a *joint thresholding*. This is still an open problem and will be investigated. For correlated noise, [7] proved some minimax properties of using a modified *universal threshold*, $\sigma_n^{(s)} \sqrt{2 \log N}$, where $\sigma_n^{(s)}$ is the standard deviation of the noise at decomposition level s . The framework is for a deterministic signal, however, and not the Bayesian framework used here. Thus, for simplicity, we separate the coefficients into groups of uncorrelated coefficients before using the adaptive thresholding algorithm.

There are two other minor details in the implementation. First, one needs to alter the noise power σ^2 at each decomposition scale to $\sigma^2/4^s$ due to the renormalization of the filters. Second, the definition of the parent coefficient used in the neighborhood of the context is slightly changed: the parent of a coefficient in scale s is simply the coefficient at the same spatial location in scale $s + 1$.

D. Alternative Methods

There are several other possible alternative approaches which will be discussed below.

1) *Different Estimation of Variance*: An alternative common approach for estimating the local variance is to use the points in a local neighborhood around $Y[i, j]$ (as in [8]). This is a simpler way than the indirect way of first grouping coefficients with similar context, and then estimating the variance. As demonstrated by the good performance of the image coder in [8], the variance estimate from a local neighborhood yields an estimate good enough for adapting the coder. However, our experience with noisy images shows that such an estimate yields considerably more unreliable variance estimates, $\sigma_X^2[i, j]$, and also blotchy denoised image. This is because the estimate is highly sensitive to the window size we choose: a small window contains few points and thus yields

unreliable estimates; a large window adapts slowly to changing characteristics. The context-based grouping allows one to congregate those coefficients with similar context though not necessarily spatially adjacent. It also allows a large number of coefficients to be used in the variance estimation, thus yielding a more reliable estimate. Simulations show that the performance is not sensitive to the neighborhood choice \mathcal{B}_{ij} and the weight \mathbf{w} used in the context calculation, as a simple equally weighted average of the eight nearest neighbors in the same subband yields approximately the same result.

2) *One-Pass Algorithm*: The method we have proposed is a *two-pass* process: the first pass calculates the weighted average $\{Z[i, j]\}$ of the absolute values of the neighboring *noisy* coefficients, and then $\{Z[i, j]\}$ are sorted; the second pass collects the noisy coefficients with similar values of $Z[i, j]$, estimates the signal variance, $\sigma_X^2[i, j]$, from the noisy coefficients and the thresholds for the thresholding function. It is worthwhile to investigate the algorithm performance when the context modeling and the parameter estimation are performed on the *denoised* coefficients instead, since, intuitively, if the coefficients are really denoised, they should yield more reliable information. This simple intuition is, however, not as straightforward to implement as it seems. To do this in a two-pass algorithm is difficult, since $Z[i, j]$ is a weighted average of neighboring denoised coefficients, but the threshold used to denoise these coefficients are estimated from other denoised coefficients with similar context. A simple-minded alternative solution is to use a *one-pass* modification of our algorithm, where the conditioning and estimation are based on the *causal, denoised* coefficients, much along the same philosophy as one-pass compression methods conditioning on causal quantized data [8], [16]. Assume a scanning order of row by row, and initialize the first coefficient as already denoised, that is, $\hat{X}[1, 1] = Y[1, 1]$. For every new coefficient at location $[i, j]$, the context is conditioned on $Z[i, j] = \mathbf{w}^t \mathbf{u}_{ij}$ where \mathbf{u}_{ij} is now the vector containing the absolute value of denoised coefficients $\hat{X}[i, j]$ in a causal neighborhood, and the elements of \mathbf{w} are simply the equal weights. These choices are made for simplicity since the denoising performance is not too sensitive to the neighborhood selection and weight vector \mathbf{w} . The GGD parameter $\sigma_X[i, j]$ is estimated from past denoised coefficients whose contexts are similar, and $2L + 1$ coefficients are used (or all of the available coefficients so far if less than $2L + 1$ coefficients have been denoised.) Since the coefficients are already denoised, the estimation of $\sigma_X[i, j]$ is $\hat{\sigma}_X^2[i, j] = (1/2L + 1) \sum_{[k, \ell] \in \mathcal{B}_{ij}} \hat{X}[k, \ell]^2$ instead of (10). Simulations show this approach to run into problems especially when the noise power σ_n^2 is large, causing many coefficients to be denoised to zero. Having too many consecutive zero coefficients is likely to cause $\hat{\sigma}_X[i, j]$ to be zero, which then translates to an infinite threshold (i.e., $Y[i, j]$ is thresholded to zero). This in turn may cause all the subsequent coefficients to be thresholded to zero. This phenomenon is frequently encountered in backward adaptive compression methods which adapts based on causal quantized coefficients: a run of zero coefficients may cause all subsequent coefficients to be quantized to zero as well. Some work ameliorate this problem by looking ahead to identify *unpredictable sets*, coefficients whose neighbors are zero, but who should not be quantized to zero [8], [16].

This logic can be applied to the denoising framework as well. When the algorithm computes $Z[i, j]$, it identifies the locations $[k, \ell]$ in the causal neighborhood \mathcal{B}_{ij} for which $\hat{X}[k, \ell] = 0$ but $|Y[k, \ell]| \geq \sigma_n$, and these $Y[k, \ell]$ are substituted for the zero $\hat{X}[k, \ell]$ to be used in the computation of $Z[i, j]$ and $\hat{\sigma}_X[i, j]$. Simulations show the resulting images to yield worse MSE's than the previously proposed method, and they are visually considerably more noisy.

Another variation is to use the denoised coefficient for context modeling, but the observed noisy coefficients for estimating $\sigma_X^2[i, j]$ as in (10). Again, without taking some caution about the runs of zero coefficients, the variance estimate may be inadequate for several rows (recall the scanning is row by row) before having enough nonzero causal neighbors for collecting valid information. The denoised images are also similar to the ones described above, having worse MSE's than the proposed two-pass algorithm, and are visually more noisy.

3) *Heteroscedasticity Model*: A central part of our spatially adaptive algorithm is based on modeling the variance $\sigma_X^2[i, j]$ to be nonconstant and varying throughout the image. This is reminiscent of the *heteroscedasticity*, or nonconstant variance, problem in statistics. Let $\{Y[i, j]\}$ be the observed noisy wavelet coefficients, and each $Y[i, j]$ a random variable whose variance $\gamma^2[i, j]$ is nonconstant. A common approach to the heteroscedasticity problem is to model $\gamma^2[i, j]$ as a function of some design vector, \mathbf{u}_{ij} . Traditionally there are two approaches in estimating this function: parametric and nonparametric. Since we have an assumed distribution on the wavelet coefficients (i.e., GGD), the parametric approach will be used here. The readers are referred to [1], [11] and related literatures for more details on heteroscedasticity models. Using a parametric function to describe the variance $\gamma^2[i, j]$ has the advantage that it allows a compact representation of the nonconstant variance, useful for image analysis and understanding. In contrast, although the nonparametric approach described in Section II-B works well, it does not lend itself to any tractable analysis.

In the previous section, we have described the noisy coefficient $Y[i, j]$ as a sum of two random variables, $X[i, j] \sim \text{GGD}$ and $V[i, j] \sim \text{Gaussian}$. Unless $X[i, j]$ is a Gaussian distributed random variable, there is no closed form expression for the distribution of $Y[i, j]$. However, often one observes the wavelet coefficients for images to be sharply peaked at zero, better described by the Laplacian density function. Furthermore, the noisy coefficients also form a histogram which is sharply peaked at zero. Thus, for simplicity and for the sake of tractable analysis, we assume the noisy coefficient $Y[i, j]$ to be Laplacian distributed, or, alternatively, $|Y[i, j]|$ be exponentially distributed. Similar to the context modeling framework in Section II-B, let the design vector \mathbf{u}_{ij} at location $[i, j]$ be the vector containing the absolute value of the eight closest neighboring (noisy) coefficients, \mathbf{w} be the unknown regression parameter (i.e., the weights for the weighted average of the neighboring coefficients contained in \mathbf{u}_{ij}), and the variance for Y_{ij} be a function of $\mathbf{w}^t \mathbf{u}_{ij}$. Formally, our heteroscedasticity model is

$$|Y[i, j]| \sim \frac{1}{\gamma[i, j]} e^{-y/\gamma[i, j]}, \quad y \geq 0 \quad (13)$$

where the standard deviation is

$$\gamma[i, j] = K_{\theta}(\mathbf{w}^t \mathbf{u}_{ij}) \quad (14)$$

and $K_{\theta}(\cdot)$ is a smooth function such as a polynomial of order r , with unknown parameter $(r+1) \times 1$ vector θ . Modeling $\gamma[i, j]$ as a function of $\mathbf{w}^t \mathbf{u}_{ij}$ can be justified by observing that the plot of $\{(\mathbf{w}^t \mathbf{u}_{ij}, Y[i, j])\}_{i,j}$ often resides within a cone shape (see Fig. 5), implying that the variability of $Y[i, j]$ depends highly on $\mathbf{w}^t \mathbf{u}_{ij}$.

To estimate the parameters θ and \mathbf{w} , the likelihood approach is used. The negative log-likelihood of $|Y[i, j]|$ is

$$\log K_{\theta}(\mathbf{w}^t \mathbf{u}_{ij}) + \frac{|Y[i, j]|}{K_{\theta}(\mathbf{w}^t \mathbf{u}_{ij})}. \quad (15)$$

For $\{|Y[i, j]|\}_{i,j=1,\dots,N}$, the negative log-likelihood, or the likelihood function, is

$$L(\theta, \mathbf{w}) = \sum_{i,j=1}^N \left(\log K_{\theta}(\mathbf{w}^t \mathbf{u}_{ij}) + \frac{|Y[i, j]|}{K_{\theta}(\mathbf{w}^t \mathbf{u}_{ij})} \right). \quad (16)$$

The likelihood function (16) is minimized over both parameters θ and \mathbf{w} to find their optimal values. One way to do this is to start with an initial \mathbf{w} being the linear least squares estimate, \mathbf{w}_{LS} as in (8). Then θ is estimated as

$$\hat{\theta} = \arg \min_{\theta} L(\theta, \hat{\mathbf{w}}_{LS}). \quad (17)$$

The regression parameter \mathbf{w} is refined one step further as

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} L(\hat{\theta}, \mathbf{w}). \quad (18)$$

After obtaining \mathbf{w} and $\hat{\theta}$, the standard deviation of $Y[i, j]$ is estimated by $\hat{\gamma}[i, j] = K_{\hat{\theta}}(\mathbf{w}^t \mathbf{u}_{ij})$, and the variance estimate of the clean coefficient $X[i, j]$ is $\hat{\sigma}_X^2[i, j] = \max(0, \hat{\gamma}^2[i, j] - \hat{\sigma}_n^2)$. The threshold is then calculated as before to be $T_B[i, j] = \hat{\sigma}_n^2 / \hat{\sigma}_X[i, j]$.

Polynomials of order $r = 1, 2$ were experimented with, and a different set of polynomial parameters is found for each subband. Simulations show this parametric estimation of $\gamma^2[i, j]$ to differentiate well between regions of high energy (e.g., edges and textures) and smooth areas. That is, the variance estimate is larger in the edge and texture region, and smaller in the smooth regions. However, these values are not appropriate since the subsequently calculated variance estimate of $X[i, j]$, $\hat{\sigma}_X^2[i, j]$, results in zero in many subbands, which in turn translates to killing all the coefficients in the thresholding. This phenomenon may be due to the disparity between this parametric modeling of the nonconstant variance and the noisy observation modeling: in the parametric approach, the observed noisy coefficients are modeled as Laplacian distributed, whereas in the original framework, the observations are *sums* of a Laplacian and Gaussian random variable. Nevertheless, the likelihood approach to the heteroscedasticity problem may be valuable to other applications.

III. EXPERIMENTAL RESULTS

The images *Barbara* and *Lena*, of size 512×512 , are used as test images. *iid* Gaussian noise at different levels of σ_n^2

TABLE I
MSE RESULTS OF DIFFERENT DENOISING METHODS FOR VARIOUS TEST IMAGES AND σ_n VALUES

MSE/ σ_n	12.5	15	17.5	20	22.5	25
Method	<i>lena</i>					
<i>AdaptShr</i>	31.2	37.8	44.3	51.2	57.5	64.5
<i>OracleShr</i>	36.4	44.3	52.1	59.6	67.4	74.3
<i>SI-AdaptShr</i>	24.9	29.9	35.2	40.2	45.2	50.8
<i>SI-OracleShr</i>	29.8	35.9	42.3	48.7	55.7	61.2
Mihçak [10]	-	37.6	-	52.1	-	67.5
	<i>barbara</i>					
<i>AdaptShr</i>	52.4	66.2	79.7	95.6	111.0	124.0
<i>OracleShr</i>	63.1	82.6	99.6	119.3	138.0	154.4
<i>SI-AdaptShr</i>	39.5	50.4	60.7	73.2	85.3	96.2
<i>SI-OracleShr</i>	51.2	66.3	81.0	96.7	112.0	128.2
Mihçak [10]	-	61.2	-	88.7	-	117.5

are generated using *randn* in MATLAB. For the orthogonal wavelet transform, four levels of decomposition are used, and the wavelet employed is Daubechies' symmet with eight vanishing moments [5]. There are five methods that are compared, and the MSE results are shown in Table I, with the best one highlighted in bold font. The *AdaptShrink* method refers to the proposed adaptive thresholding method using the orthogonal transform DWT (the soft-thresholding function is sometimes referred to as shrinkage, and hence the name). *SI-AdaptShrink* is the adaptive thresholding using the shift-invariant (SI), nonsubsampling wavelet transform. These two are compared against the best uniform thresholding techniques (in the MSE sense) when the original uncorrupted image is assumed to be known. For uniform thresholding with DWT, in each subband, we find the *oracle* threshold T_o as

$$T_o = \arg \min_T \sum_{i,j} (\eta_T(Y[i, j]) - X[i, j])^2 \quad (19)$$

where $Y[i, j]$ and $X[i, j]$ are the wavelet coefficients of the noisy observation g and original image f , respectively. This method is referred to as *OracleShrink*. Similarly, this is extended to the nonsubsampling wavelet transform, where a different threshold is found for each set of uncorrelated coefficients within each subband (thus 2^{2s} thresholds for a subband at scale s). This method is coined *SI-OracleShrink*. Results from the denoising algorithm recently proposed in [10] is also listed in Table I for comparison. Reference [10] uses a locally estimated GGD and applies this model to the denoising problem using a Bayesian estimate. Fig. 7 shows the comparison of the different methods on a magnified region in the *Barbara* image for $\sigma = 25$ and the *lena* image for $\sigma = 22.5$. The *SI-AdaptShrink* method outperforms all the other methods in both visual quality and MSE performance. It yields significantly less ringing artifacts and blotchiness than the methods using DWT. *SI-OracleShrink* still shows significant noise in the smooth background. Thus, it



Fig. 7. Comparing the results of various denoising methods, for *lena* corrupted by noise $\sigma_n = 22.5$ and *barbara* by noise $\sigma_n = 25$. (a) Original, (b) noisy observation, (c) adaptive thresholding in DWT basis (*AdaptShrink*), (d) oracle uniform thresholding in DWT basis (*OracleShrink*), (e) spatially adaptive thresholding in overcomplete expansion (*SI-AdaptShrink*), and (f) oracle uniform thresholding in overcomplete expansion (*SI-OracleShrink*). This figure can also be found at <http://www-wavelet.eecs.berkeley.edu/~grchang/SpatialDenoise.html>.

is both the spatial adaptive thresholds and the overcomplete representation that contribute to the superior quality of *SI-AdaptShrink*. The adaptive methods denoise better especially in the flat regions, where the uniform methods yields images with much noise and specks.

IV. CONCLUSION

We have proposed a simple and effective spatially and scale-wise adaptive method for denoising via wavelet thresholding in an overcomplete expansion. The adaptivity is based on context-modeling which enables a pixel-wise estimation of the signal variance and thus of the best threshold. The issue of

spatially adapting the threshold values has not been addressed in the literature. As we have shown in this paper, adapting the threshold values to local signal energy allows us to keep much of the edge and texture details, while eliminating most of the noise in smooth regions, something that may be hard to achieve with a uniform threshold. The results show substantial improvement over the optimal uniform thresholding both in visual quality and mean squared error.

ACKNOWLEDGMENT

The authors would like to thank A. Ortega for his insightful comments.

REFERENCES

- [1] R. J. Carroll, "Adapting for heteroscedasticity in linear models," *Ann. Statist.*, vol. 10, no. 4, pp. 1224–1233, 1982.
- [2] S. G. Chang and M. Vetterli, "Spatial adaptive wavelet thresholding for image denoising," in *Proc. IEEE Int. Conf. Image Processing*, vol. 2, Santa Barbara, CA, Nov. 1997, pp. 374–377.
- [3] S. G. Chang, B. Yu, and M. Vetterli, "Image denoising via lossy compression and wavelet thresholding," *IEEE Trans. Image Processing*, vol. 9, pp. 1532–1546, Sept. 2000.
- [4] R. R. Coifman and D. L. Donoho, "Translation-invariant de-noising," in *Wavelets and Statistics*, A. Antoniadis and G. Oppenheim, Eds. Berlin, Germany: Springer-Verlag, 1995.
- [5] I. Daubechies, *Ten Lectures on Wavelets*, Vol. 61 of *Proc. CBMS-NSF Regional Conference Series in Applied Mathematics*. Philadelphia, PA: SIAM, 1992.
- [6] D. L. Donoho and I. M. Johnstone, "Ideal spatial adaptation via wavelet shrinkage," *Biometrika*, vol. 81, pp. 425–455, 1994.
- [7] I. M. Johnstone and B. W. Silverman, "Wavelet threshold estimators for data with correlated noise," *J. R. Statist. Soc.*, ser. B, vol. 59, 1997.
- [8] S. LoPresto, K. Ramchandran, and M. Orchard, "Image coding based on mixture modeling of wavelet coefficients and a fast estimation-quantization framework," in *Proc. Data Compression Conf.*, Snowbird, UT, Mar. 1997, pp. 221–230.
- [9] S. Mallat, "A theory for multiresolution signal decomposition: The wavelet representation," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 11, pp. 674–693, July 1989.
- [10] M. K. Mihçak, I. Kozintsev, and K. Ramchandran, "Spatially adaptive statistical modeling of wavelet image coefficients and its application to denoising," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Processing*, vol. 6, Mar. 1999, pp. 3253–3256.
- [11] H.-G. Müller and U. Stadtmüller, "Estimation of heteroscedasticity in regression analysis," *Ann. Statist.*, vol. 15, no. 2, pp. 610–625, 1987.
- [12] J. Shapiro, "Embedded image coding using zerotrees of wavelet coefficients," *IEEE Trans. Signal Processing*, vol. 41, pp. 3445–3462, Dec. 1993.
- [13] E. Simoncelli and E. Adelson, "Noise removal via Bayesian wavelet coring," in *Proc. IEEE Int. Conf. Image Processing*, vol. 1, Lausanne, Switzerland, Sept. 1996, pp. 379–382.
- [14] M. Vetterli and J. Kovačević, *Wavelets and Subband Coding*. Englewood Cliffs, NJ: Prentice-Hall, 1995.
- [15] P. H. Westerink, J. Biemond, and D. E. Boekee, "An optimal bit allocation algorithm for sub-band coding," in *Proc. IEEE Int. Conf. Acoustics, Speech, Signal Processing*, Dallas, TX, Apr. 1987, pp. 1378–1381.
- [16] Y. Yoo, A. Ortega, and B. Yu, "Image subband coding using context-based classification and adaptive quantization," *IEEE Trans. Image Processing*, vol. 8, pp. 1702–1715, Dec. 1999.



S. Grace Chang (S'95) received the B.S. degree from the Massachusetts Institute of Technology, Cambridge, in 1993, and the M.S. and Ph.D. degrees from the University of California, Berkeley, in 1995 and 1998, respectively, all in electrical engineering.

She was with Hewlett-Packard Laboratories, Palo Alto, CA, and is now with Hewlett-Packard Co., Grenoble, France. Her research interests include image enhancement and compression, Internet applications and content delivery, and telecommunication systems.

Dr. Chang was a recipient of the National Science Foundation Graduate Fellowship and University of California Dissertation Fellowship.



Bin Yu (A'92–SM'97) received the B.S. degree in mathematics from Peking University, China, in 1984, and the M.S. and Ph.D. degrees in statistics from the University of California, Berkeley, in 1987 and 1990, respectively.

She is an Associate Professor of statistics with the University of California, Berkeley. Her research interests include statistical inference, information theory, signal compression and denoising, bioinformatics, and remote sensing. She has published over 30 technical papers in journals such as *IEEE*

TRANSACTIONS ON INFORMATION THEORY, *IEEE TRANSACTIONS ON IMAGE PROCESSING*, *IEEE TRANSACTIONS ON GEOSCIENCE AND REMOTE SENSING*, *The Annals of Statistics*, *Annals of Probability*, *Journal of American Statistical Association*, and *Genomics*. She has held faculty positions at the University of Wisconsin, Madison, and Yale University, New Haven, CT, and was a Postdoctoral Fellow with the Mathematical Science Research Institute (MSRI), University of California, Berkeley.

Dr. Yu was in the S. S. Chern Mathematics Exchange Program between China and the U.S. in 1985. She is a Fellow of the Institute of Mathematical Statistics (IMS), and a member of The American Statistical Association (ASA). She is serving on the Board of Governors of the IEEE Information Theory Society, and as an Associate Editor for *The Annals of Statistics* and for *Statistica Sinica*.



Martin Vetterli (S'86–M'86–SM'90–F'95) received the Dipl. El.-Ing. degree from ETH Zürich (ETHZ), Zürich, Switzerland, in 1981, the M.S. degree from Stanford University, Stanford, CA, in 1982, and the Doctorat ès Science degree from EPFL Lausanne (EPFL), Lausanne, Switzerland, in 1986.

He was a Research Assistant at Stanford University and EPFL, and was with Siemens and AT&T Bell Laboratories. In 1986, he joined Columbia University, New York, where he was an Associate Professor of electrical engineering and Co-Director

of the Image and Advanced Television Laboratory. In 1993, he joined the University of California, Berkeley, where he was a Professor in the Department of Electrical Engineering and Computer Sciences until 1997, and holds now Adjunct Professor position. Since 1995, he has been a Professor of communication systems at EPFL, where he chaired the Communications Systems Division (1996–1997), and heads the Audio-Visual Communications Laboratory. He held visiting positions at ETHZ in 1990 and Stanford University in 1998. He is on the editorial boards of *Annals of Telecommunications*, *Applied and Computational Harmonic Analysis*, and the *Journal of Fourier Analysis and Applications*. He is the co-author, with J. Kovačević, of the book *Wavelets and Subband Coding* (Englewood Cliffs, NJ: Prentice-Hall, 1995). He has published about 75 journal papers on a variety of topics in signal and image processing and holds five patents. His research interests include wavelets, multirate signal processing, computational complexity, signal processing for telecommunications, digital video processing, and compression and wireless video communications.

Dr. Vetterli is a member of SIAM and was the Area Editor for Speech, Image, Video, and Signal Processing for the *IEEE TRANSACTIONS ON COMMUNICATIONS*. He received the Best Paper Award of EURASIP in 1984 for his paper on multidimensional subband coding, the Research Prize of the Brown Boverly Corporation, Switzerland, in 1986 for his doctoral thesis, the IEEE Signal Processing Society's Senior Award in 1991 and 1996 (for papers with D. LeGall and K. Ramchandran, respectively). He was a IEEE Signal Processing Distinguished Lecturer in 1999. He received the Swiss National Latsis Prize in 1996 and the SPIE Presidential award in 1999. He has been a Plenary Speaker at various conferences including the 1992 IEEE International Conference on Acoustics, Speech, and Signal Processing.