

# Adaptive Wavelet Thresholding for Image Denoising and Compression

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

**Abstract**—The first part of this paper proposes an adaptive, data-driven threshold for image denoising via wavelet soft-thresholding. The threshold is derived in a Bayesian framework, and the prior used on the wavelet coefficients is the generalized Gaussian distribution (GGD) widely used in image processing applications. The proposed threshold is simple and closed-form, and it is adaptive to each subband because it depends on data-driven estimates of the parameters. Experimental results show that the proposed method, called *BayesShrink*, is typically within 5% of the MSE of the best soft-thresholding benchmark with the image assumed known. It also outperforms Donoho and Johnstone's *SureShrink* most of the time.

The second part of the paper attempts to further validate recent claims that lossy compression can be used for denoising. The *BayesShrink* threshold can aid in the parameter selection of a coder designed with the intention of denoising, and thus achieving simultaneous denoising and compression. Specifically, the zero-zone in the quantization step of compression is analogous to the threshold value in the thresholding function. The remaining coder design parameters are chosen based on a criterion derived from Rissanen's minimum description length (MDL) principle. Experiments show that this compression method does indeed remove noise significantly, especially for large noise power. However, it introduces quantization noise and should be used only if bitrate were an additional concern to denoising.

**Index Terms**—Adaptive method, image compression, image denoising, image restoration, wavelet thresholding.

## I. INTRODUCTION

**A**N IMAGE is often corrupted by noise in its acquisition or transmission. The goal of denoising is to remove the noise while retaining as much as possible the important signal features. Traditionally, this is achieved by linear processing such as Wiener filtering. A vast literature has emerged recently on

signal denoising using nonlinear techniques, in the setting of additive white Gaussian noise. The seminal work on signal denoising via *wavelet thresholding* or *shrinkage* of Donoho and Johnstone ([13]–[16]) have shown that various wavelet thresholding schemes for denoising have near-optimal properties in the minimax sense and perform well in simulation studies of one-dimensional curve estimation. It has been shown to have better rates of convergence than linear methods for approximating functions in Besov spaces ([13], [14]). Thresholding is a nonlinear technique, yet it is very simple because it operates on one wavelet coefficient at a time. Alternative approaches to nonlinear wavelet-based denoising can be found in, for example, [1], [4], [8]–[10], [12], [18], [19], [24], [27]–[29], [32], [33], [35], and references therein.

On a seemingly unrelated front, lossy compression has been proposed for denoising in several works [6], [5], [21], [25], [28]. Concerns regarding the compression rate were explicitly addressed. This is important because any practical coder must assume a limited resource (such as bits) at its disposal for representing the data. Other works [4], [12]–[16] also addressed the connection between compression and denoising, especially with nonlinear algorithms such as wavelet thresholding in a mathematical framework. However, these latter works were not concerned with quantization and bitrates: compression results from a reduced number of nonzero wavelet coefficients, and not from an explicit design of a coder.

The intuition behind using lossy compression for denoising may be explained as follows. A signal typically has structural correlations that a good coder can exploit to yield a concise representation. White noise, however, does not have structural redundancies and thus is not easily compressible. Hence, a good compression method can provide a suitable model for distinguishing between signal and noise. The discussion will be restricted to wavelet-based coders, though these insights can be extended to other transform-domain coders as well. A concrete connection between lossy compression and denoising can easily be seen when one examines the similarity between thresholding and quantization, the latter of which is a necessary step in a practical lossy coder. That is, the quantization of wavelet coefficients *with a zero-zone* is an approximation to the thresholding function (see Fig. 1). Thus, provided that the quantization outside of the zero-zone does not introduce significant distortion, it follows that wavelet-based lossy compression achieves denoising. With this connection in mind, this paper is about wavelet thresholding for image denoising and also for lossy compression. The threshold choice aids the lossy coder to choose its zero-zone, and the resulting coder achieves simultaneous denoising and compression if such property is desired.

Manuscript received January 22, 1998; revised April 7, 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, Santa Barbara, CA, October 1997. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Patrick L. Combettes.

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)06914-1.

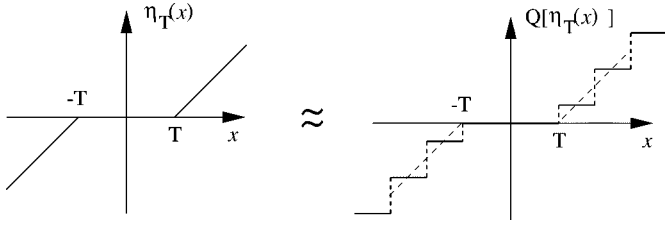


Fig. 1. Thresholding function can be approximated by quantization with a zero-zone.

The theoretical formalization of filtering additive *iid* Gaussian noise (of zero-mean and standard deviation  $\sigma$ ) via thresholding wavelet coefficients was pioneered by Donoho and Johnstone [14]. A wavelet coefficient is compared to a given threshold and is set to zero if its magnitude is less than the threshold; otherwise, it is kept or modified (depending on the thresholding rule). The threshold acts as an oracle which distinguishes between the insignificant coefficients likely due to noise, and the significant coefficients consisting of important signal structures. Thresholding rules are especially effective for signals with sparse or near-sparse representations where only a small subset of the coefficients represents all or most of the signal energy. Thresholding essentially creates a region around zero where the coefficients are considered negligible. Outside of this region, the thresholded coefficients are kept to full precision (that is, without quantization). Their most well-known thresholding methods include *VisuShrink* [14] and *SureShrink* [15]. These threshold choices enjoy asymptotic minimax optimalities over function spaces such as Besov spaces. For image denoising, however, *VisuShrink* is known to yield overly smoothed images. This is because its threshold choice,  $\sigma\sqrt{2\log M}$  (called the *universal threshold* and  $\sigma^2$  is the noise variance), can be unwarrantedly large due to its dependence on the number of samples,  $M$ , which is more than  $10^5$  for a typical test image of size  $512 \times 512$ . *SureShrink* uses a hybrid of the universal threshold and the SURE threshold, derived from minimizing Stein's unbiased risk estimator [30], and has been shown to perform well. *SureShrink* will be the main comparison to the method proposed here, and, as will be seen later in this paper, our proposed threshold often yields better result.

Since the works of Donoho and Johnstone, there has been much research on finding thresholds for nonparametric estimation in statistics. However, few are specifically tailored for images. In this paper, we propose a framework and a near-optimal threshold in this framework more suitable for image denoising. This approach can be formally described as Bayesian, but this only describes our mathematical formulation, not our philosophy. The formulation is grounded on the empirical observation that the wavelet coefficients in a subband of a natural image can be summarized adequately by a *generalized Gaussian distribution* (GGD). This observation is well-accepted in the image processing community (for example, see [20], [22], [23], [29], [34], [36]) and is used for state-of-the-art image coders in [20], [22], [36]. It follows from this observation that the average MSE (in a subband) can be approximated by the corresponding Bayesian squared error risk with the GGD as the prior applied to each in an *iid* fashion. That is, a sum is approximated by an integral. We emphasize that this is an analytical approximation and our framework is broader than assuming wavelet

coefficients are *iid* draws from a GGD. The goal is to find the soft-threshold that minimizes this Bayesian risk, and we call our method *BayesShrink*.

The proposed Bayesian risk minimization is subband-dependent. Given the signal being generalized Gaussian distributed and the noise being Gaussian, via numerical calculation a nearly optimal threshold for soft-thresholding is found to be  $T_B = \sigma^2/\sigma_X$  (where  $\sigma^2$  is the noise variance and  $\sigma_X^2$  the signal variance). This threshold gives a risk within 5% of the minimal risk over a broad range of parameters in the GGD family. To make this threshold data-driven, the parameters  $\sigma_X$  and  $\sigma$  are estimated from the observed data, one set for each subband.

To achieve simultaneous denoising and compression, the nonzero thresholded wavelet coefficients need to be quantized. Uniform quantizer and centroid reconstruction is used on the GGD. The design parameters of the coder, such as the number of quantization levels and binwidths, are decided based on a criterion derived from Rissanen's *minimum description length* (MDL) principle [26]. This criterion balances the tradeoff between the compression rate and distortion, and yields a nice interpretation of operating at a fixed slope on the rate-distortion curve.

The paper is organized as follows. In Section II, the wavelet thresholding idea is introduced. Section II-A explains the derivation of the *BayesShrink* threshold by minimizing a Bayesian risk with squared error. The lossy compression based on the MDL criterion is explained in Section III. Experimental results on several test images are shown in Section IV and compared with *SureShrink*. To benchmark against the best possible performance of a threshold estimate, the comparisons also include *OracleShrink*, the best soft-thresholding estimate obtainable assuming the original image known, and *OracleThresh*, the best hard-thresholding counterpart. The *BayesShrink* method often comes to within 5% of the MSEs of *OracleShrink*, and is better than *SureShrink* up to 8% most of the time, or is within 1% if it is worse. Furthermore, the *BayesShrink* threshold is very easy to compute. *BayesShrink* with the additional MDL-based compression, as expected, introduces quantization noise to the image. This distortion may negate the denoising achieved by thresholding, especially when  $\sigma$  is small. However, for larger values of  $\sigma$ , the MSE due to the lossy compression is still significantly lower than that of the noisy image, while fewer bits are used to code the image, thus achieving both denoising and compression.

## II. WAVELET THRESHOLDING AND THRESHOLD SELECTION

Let the signal be  $\{f_{ij}, i, j = 1, \dots, N\}$ , where  $N$  is some integer power of 2. It has been corrupted by additive noise and one observes

$$g_{ij} = f_{ij} + \varepsilon_{ij}, \quad i, j = 1, \dots, N \quad (1)$$

where  $\{\varepsilon_{ij}\}$  are independent and identically distributed (*iid*) as normal  $N(0, \sigma^2)$  and independent of  $\{f_{ij}\}$ . The goal is to remove the noise, or “denoise”  $\{g_{ij}\}$ , and to obtain an estimate  $\{\hat{f}_{ij}\}$  of  $\{f_{ij}\}$  which minimizes the mean squared error (MSE),

$$\text{MSE}(\hat{\mathbf{f}}) = \frac{1}{N^2} \sum_{i,j=1}^N (\hat{f}_{ij} - f_{ij})^2. \quad (2)$$

LL <sub>3</sub>	HL <sub>3</sub>	HL <sub>2</sub>	HL <sub>1</sub>
LH <sub>3</sub>	HH <sub>3</sub>		
LH <sub>2</sub>		HH <sub>2</sub>	
LH <sub>1</sub>		HH <sub>1</sub>	

Fig. 2. Subbands of the 2-D orthogonal wavelet transform.

Let  $\mathbf{g} = \{g_{ij}\}_{i,j}$ ,  $\mathbf{f} = \{f_{ij}\}_{i,j}$ , and  $\boldsymbol{\varepsilon} = \{\varepsilon_{ij}\}_{i,j}$ ; 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 two-dimensional dyadic orthogonal wavelet transform operator, and similarly  $\mathbf{X} = \mathcal{W}\mathbf{f}$  and  $\mathbf{V} = \mathcal{W}\boldsymbol{\varepsilon}$ . The readers are referred to references such as [23], [31] for details of the two-dimensional orthogonal wavelet transform. It is convenient to label the subbands of the transform as in Fig. 2. 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  $N/2^J > 1$ . Note that since the transform is orthogonal,  $\{V_{ij}\}$  are also iid  $N(0, \sigma^2)$ .

The wavelet-thresholding denoising method 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.

There are two thresholding methods frequently used. The *soft-threshold* function (also called the shrinkage function)

$$\eta_T(x) = \text{sgn}(x) \cdot \max(|x| - T, 0) \quad (3)$$

takes the argument and shrinks it toward zero by the *threshold*  $T$ . The other popular alternative is the *hard-threshold* function

$$\psi_T(x) = x \cdot \mathbf{1}\{|x| > T\} \quad (4)$$

which keeps the input if it is larger than the threshold  $T$ ; otherwise, it is set to zero. The wavelet thresholding procedure removes noise by thresholding *only* the wavelet coefficients of the detail subbands, while keeping the low resolution coefficients unaltered.

The soft-thresholding rule is chosen over hard-thresholding for several reasons. First, soft-thresholding has been shown to achieve near-optimal minimax rate over a large range of Besov spaces [12], [14]. Second, for the generalized Gaussian prior assumed in this work, the optimal soft-thresholding estimator yields a smaller risk than the optimal hard-thresholding estimator (to be shown later in this section). Lastly, in practice, the soft-thresholding method yields more visually pleasant images over hard-thresholding because the latter is discontinuous and yields abrupt artifacts in the recovered images, especially when

the noise energy is significant. In what follows, soft-thresholding will be the primary focus.

While the idea of thresholding is simple and effective, finding a good threshold is not an easy task. For one-dimensional (1-D) deterministic signal of length  $M$ , Donoho and Johnstone [14] proposed for *VisuShrink* the universal threshold,  $T_U = \sigma\sqrt{2\log M}$ , which results in an estimate asymptotically optimal in the minimax sense (minimizing the maximum error over all possible  $M$ -sample signals). One other notable threshold is the SURE threshold [15], derived from minimizing Stein's unbiased risk estimate [30] when soft-thresholding is used. The *SureShrink* method is a hybrid of the universal and the SURE threshold, with the choice being dependent on the energy of the particular subband [15]. The SURE threshold is data-driven, does not depend on  $M$  explicitly, and *SureShrink* estimates it in a subband-adaptive manner. Moreover, *SureShrink* has yielded good image denoising performance and comes close to the true minimum MSE of the optimal soft-threshold estimator (cf. [4], [12]), and thus will be the main comparison to our proposed method.

In the statistical Bayesian literature, many works have concentrated on deriving the best threshold (or shrinkage factor) based on priors such as the Laplacian and a mixture of Gaussians (cf. [1], [8], [9], [18], [24], [27], [29], [32], [35]). With an integral approximation to the pixel-wise MSE distortion measure as discussed earlier, the formulation here is also Bayesian for finding the best soft-thresholding rule under the generalized Gaussian prior. A related work is [27] where the hard-thresholding rule is investigated for signals with Laplacian and Gaussian distributions.

The GGD has been used in many subband or wavelet-based image processing applications [2], [20], [22], [23], [29], [34], [36]. In [29], it was observed that a GGD with the shape parameter  $\beta$  ranging from 0.5 to 1 [see (1)] can adequately describe the wavelet coefficients of a large set of natural images. Our experience with images supports the same conclusion. Fig. 3 shows the histogram of the wavelet coefficients of the images shown in Fig. 9, against the generalized Gaussian curve, with the parameters labeled (the estimation of the parameters will be explained later in the text.) A heuristic can be set forward to explain why there are a large number of "small" coefficients but relatively few "large" coefficients as the GGD suggests: the small ones correspond to smooth regions in a natural image and the large ones to edges or textures.

#### A. Adaptive Threshold for BayesShrink

The GGD, following [20], is

$$GG_{\sigma_X, \beta}(x) = C(\sigma_X, \beta) \exp\{-[\alpha(\sigma_X, \beta)|x|]^\beta\} \quad (5)$$

$-\infty < x < \infty, \sigma_X > 0, \beta > 0$ , where

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

and

$$C(\sigma_X, \beta) = \frac{\beta \cdot \alpha(\sigma_X, \beta)}{2\Gamma\left(\frac{1}{\beta}\right)}$$

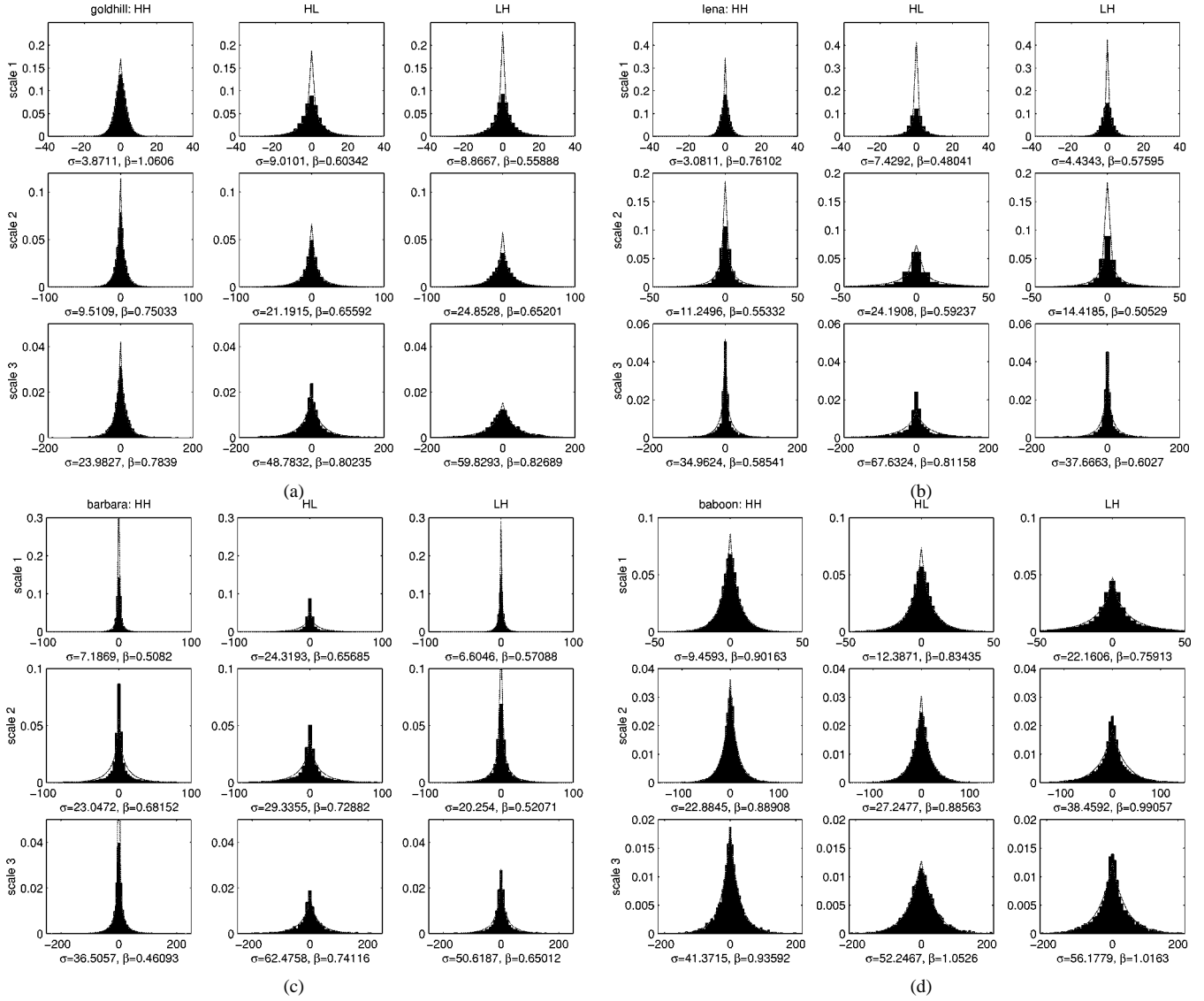


Fig. 3. Histogram of the wavelet coefficients of four test images. For each image, from top to bottom it is fine to coarse scales: from left to right, they are the *HH*, *HL*, and *LH* subbands, respectively.

and  $\Gamma(t) = \int_0^\infty e^{-u} u^{t-1} du$  is the gamma function. The parameter  $\sigma_X$  is the standard deviation and  $\beta$  is the *shape* parameter. For a given set of parameters, the objective is to find a soft-threshold  $T$  which minimizes the *Bayes risk*,

$$r(T) = E(\hat{X} - X)^2 = E_X E_{Y|X}(\hat{X} - X)^2 \quad (6)$$

where  $\hat{X} = \eta_T(Y)$ ,  $Y|X \sim N(x, \sigma^2)$  and  $X \sim GG_{\sigma_X, \beta}$ . Denote the optimal threshold by  $T^*$ ,

$$T^*(\sigma_X, \beta) = \arg \min_T r(T) \quad (7)$$

which is a function of the parameters  $\sigma_X$  and  $\beta$ . To our knowledge, there is no closed form solution for  $T^*$  for this chosen prior, thus numerical calculation is used to find its value.<sup>1</sup>

Before examining the general case, it is insightful to consider two special cases of the GGD: the Gaussian ( $\beta = 2$ ) and the

<sup>1</sup>It was observed that for the numerical calculation, it is more robust to obtain the value of  $T^*$  from locating the zero-crossing of the derivative,  $r'(T)$ , than from minimizing  $r(T)$  directly. In the Laplacian case, a recent work [17] derives an analytical equation for  $T^*$  to satisfy and calculates  $T^*$  by solving such an equation.

Laplacian ( $\beta = 1$ ) distributions. The Laplacian case is particularly interesting, because it is analytically more tractable and is often used in image processing applications.

*Case 1:* (Gaussian)  $X \sim N(0, \sigma_X^2)$  with  $\beta = 2$ . It is straightforward to verify that

$$E_X E_{Y|X}(\hat{X} - X)^2 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (\eta_T(y) - x)^2 p(y|x) p(x) dy dx \quad (8)$$

$$= \sigma_X^2 w\left(\frac{\sigma_X^2}{\sigma^2}, \frac{T}{\sigma}\right) \quad (9)$$

where

$$w(\sigma_X^2, T) = \sigma_X^2 + 2(T^2 + 1 - \sigma_X^2) \Phi\left(\frac{T}{\sqrt{1 + \sigma_X^2}}\right) - 2T(1 + \sigma_X^2) \phi(T, 1 + \sigma_X^2) \quad (10)$$

with the standard normal density function  $\phi(x, \sigma^2) = (1/\sqrt{2\pi\sigma^2}) \exp(-(x^2/2\sigma^2))$  and the survival function of the standard normal  $\Phi(x) = \int_x^\infty \phi(t, 1) dt$ .

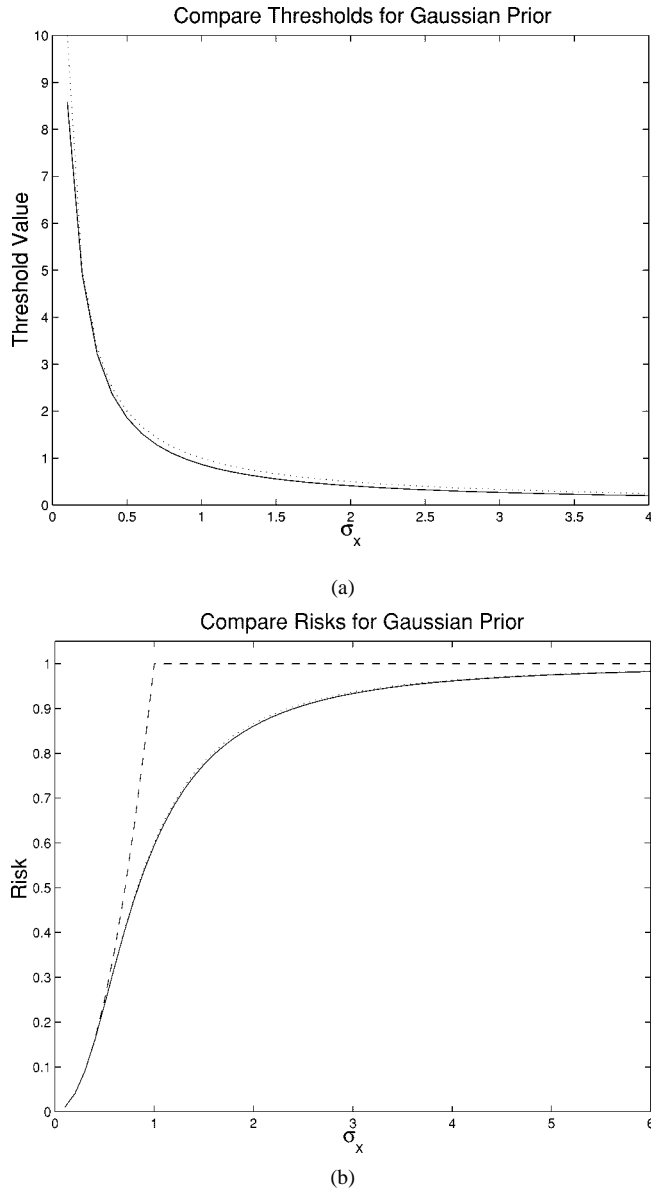


Fig. 4. Thresholding for the Gaussian prior, with  $\sigma = 1$ . (a) Compare the optimal threshold  $T^*(\sigma_X, 2)$  (solid —) and the threshold  $T_B(\sigma_X)$  (dotted  $\cdots$ ) against the standard deviation  $\sigma_X$  on the horizontal axis. (b) Compare the risk of using optimal soft-thresholding (—),  $T_B$  for soft-thresholding ( $\cdots$ ), and optimal hard-thresholding (---).

Assuming  $\sigma = 1$  for the time being, the value of  $T^*(\sigma_X, 2)$  is found numerically for different values of  $\sigma_X$  and is plotted against  $\sigma_X$  in Fig. 4(a), with

$$T_B(\sigma_X) = \frac{1}{\sigma_X} \quad (11)$$

superimposed on top. It is clear that this simple and closed-form expression,  $T_B(\sigma_X)$ , is very close to the numerically found  $T^*(\sigma_X, 2)$ . The expected risks of  $T^*(\sigma_X, 2)$  and  $T_B(\sigma_X)$  are shown in Fig. 4(b) for  $\sigma = 1$ , where the maximum deviation of  $r(T_B(\sigma_X))$  is less than 1% of the optimal risk,  $r(T^*(\sigma_X, 2))$ . For general  $\sigma$ , it is an easy scaling exercise to see that (11) becomes

$$T_B(\sigma_X) = \frac{\sigma^2}{\sigma_X}. \quad (12)$$

For a further comparison, the risk for hard-thresholding is also calculated. After some algebra, it can be shown that the risk for hard-thresholding is

$$r_h(T) = \sigma^2 + (\sigma^2 - \sigma_X^2) \cdot \left( 2T\phi(T, \sigma_X^2 + \sigma^2) + 2\bar{\Phi}\left(\frac{T}{\sqrt{\sigma_X^2 + \sigma^2}}\right) - 1 \right). \quad (13)$$

By setting to zero the derivative of (13) with respect to  $T$ , the optimal threshold is found to be

$$T_h^*(\sigma_X, 2) = \begin{cases} 0, & \text{if } \sigma_X > \sigma \\ \infty, & \text{if } \sigma_X < \sigma \\ \text{anything}, & \text{if } \sigma_X = \sigma. \end{cases} \quad (14)$$

with the associated risk

$$r_h(T_h^*) = \begin{cases} \sigma^2, & \text{if } \sigma_X > \sigma \\ \sigma_X^2, & \text{if } \sigma_X \leq \sigma. \end{cases} \quad (15)$$

Fig. 4(b) shows that both the optimal and near-optimal soft-threshold estimators,  $\eta_{T^*}(\cdot)$  and  $\eta_{T_B}(\cdot)$ , achieve lower risks than the optimal hard-threshold estimator.

The threshold  $T_B = \sigma^2/\sigma_X$  is not only nearly optimal but also has an intuitive appeal. The normalized threshold,  $T_B/\sigma$ , is inversely proportional to  $\sigma_X$ , the standard deviation of  $X$ , and proportional to  $\sigma$ , the noise standard deviation. When  $\sigma/\sigma_X \ll 1$ , the signal is much stronger than the noise,  $T_B/\sigma$  is chosen to be small in order to preserve most of the signal and remove some of the noise; vice versa, when  $\sigma/\sigma_X \gg 1$ , the noise dominates and the normalized threshold is chosen to be large to remove the noise which has overwhelmed the signal. Thus, this threshold choice adapts to both the signal and noise characteristics as reflected in the parameters  $\sigma$  and  $\sigma_X$ .

*Case 2: (Laplacian)* With  $\beta = 1$ , the GGD becomes Laplacian:  $X \sim \text{LAP}(x) = (1/\sqrt{2}\sigma_X) \exp\{-(\sqrt{2}/\sigma_X)|x|\}$ . Again for the time being let  $\sigma = 1$ . The optimal threshold  $T^*(\sigma_X, 1)$  found numerically is plotted against the standard deviation  $\sigma_X$  on the horizontal axis in Fig. 5(a). The curve of  $T^*(\sigma_X, 1)$  (in solid line —) is compared with  $T_B(\sigma_X) = 1/\sigma_X$  (in dotted line  $\cdots$ ) in Fig. 5(a). Their corresponding expected risks are shown in Fig. 5(b), and the deviation of  $r(T_B(\sigma_X))$  from the  $r(T^*)$  is less than 0.8%. This suggests that  $T_B(\sigma_X)$  also works well in the Laplacian case. For general  $\sigma$ , (12) holds again.

The threshold choice  $T_B^h(\sigma_X) = 2\sqrt{2}\sigma^2/\sigma_X$  was found independently in [27] for approximating the optimal hard-thresholding using the Laplacian prior. Fig. 5(a) compares the optimal hard-threshold,  $T_h^*(\sigma_X, 1)$ , and  $T_B^h(\sigma_X)$  to the soft-thresholds  $T^*(\sigma_X, 1)$  and  $T_B(\sigma_X)$ . The corresponding risks are plotted in Fig. 5(b), which shows the soft-thresholding rule to yield a lower risk for this chosen prior. In fact, for  $\sigma_X$  larger than approximately  $1.3\sigma$ , the risk of the approximate hard-threshold is worse than if no thresholding were performed (which yields a risk of  $\sigma^2$ ).

When  $\sigma_X$  tends to infinity or the SNR is going to infinity, an asymptotic approximation of  $T^*(\sigma_X, 1)$  is derived in [17] in this Laplacian case to be  $\sqrt{2}/\sigma_X$ . However, in the same article, this asymptotic approximation is outperformed in seven test images by the our proposed threshold,  $T_B(\sigma_X)$ .

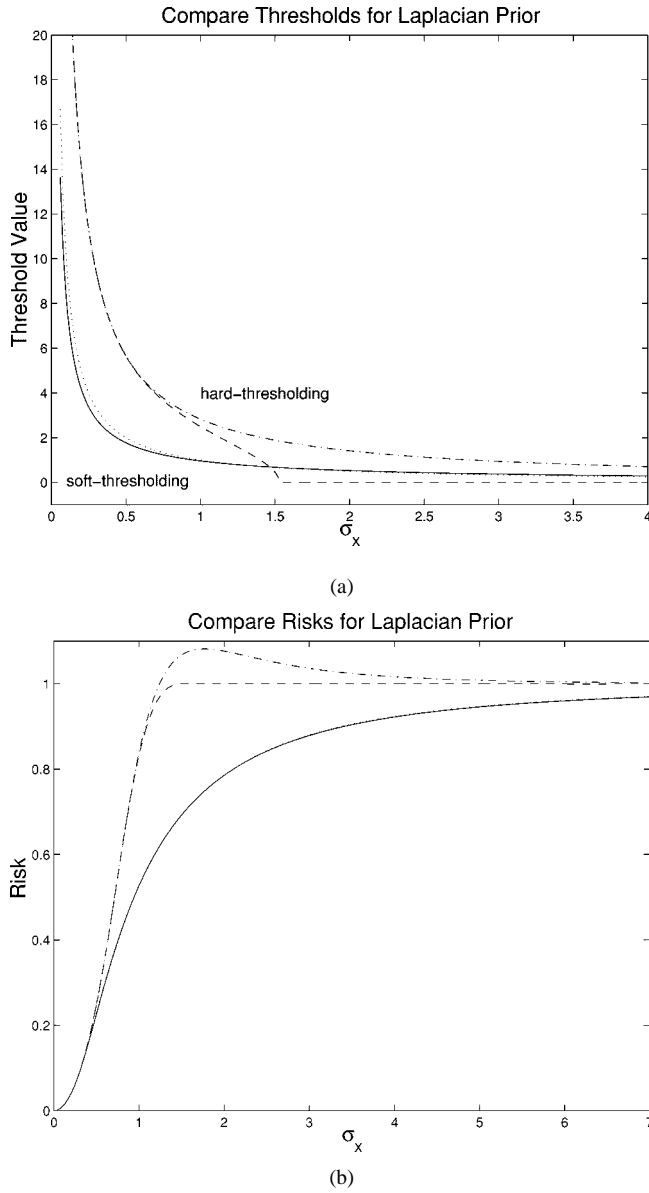


Fig. 5. Thresholding for the Laplacian prior, with  $\sigma = 1$ . (a) Compare the optimal soft-threshold  $T^*(\sigma_X, 1)$  (—), the *BayesShrink* threshold  $T_B(\sigma_X)$  ( $\cdots$ ), the optimal hard-threshold  $T_h^*(\sigma_X, 1)$  (---), and the threshold  $T_{Bh}$  (---) against the standard deviation on the horizontal axis. (b) Their corresponding risks.

With these insights from the special cases, the discussion now returns to the general case of GGD.

**Case 3: (Generalized Gaussian)** The proposed threshold  $T_B(\sigma_X)$  in (12) has been found to work well for the general case. Let  $\sigma = 1$ . In Fig. 6(a), each dotted line ( $\cdots$ ) is the optimal threshold  $T^*(\sigma_X, \beta)$  for a given fixed  $\beta$ , plotted against  $\sigma_X$  on the horizontal axis. The values  $\beta = 0.6, 1, 2, 3, 4$  are shown. The threshold,  $T_B(\sigma_X) = 1/\sigma_X$ , is plotted with the solid line (—). The curve of the optimal threshold that lies closest to  $T_B(\sigma_X)$  is for  $T^*(\sigma_X, \beta = 1)$ , the Laplacian case, while other curves deviate from  $T_B$  as  $\beta$  moves away from 1. Fig. 6(b) shows the corresponding risks. The deviation between the optimal risk  $r(T^*)$  and  $r(T_B)$  grows as  $\beta$  moves away from 1, but the error is still within 5% of the optimal  $r(T^*)$  for the curves shown in Fig. 6(b). Because the threshold  $T_B$  depends

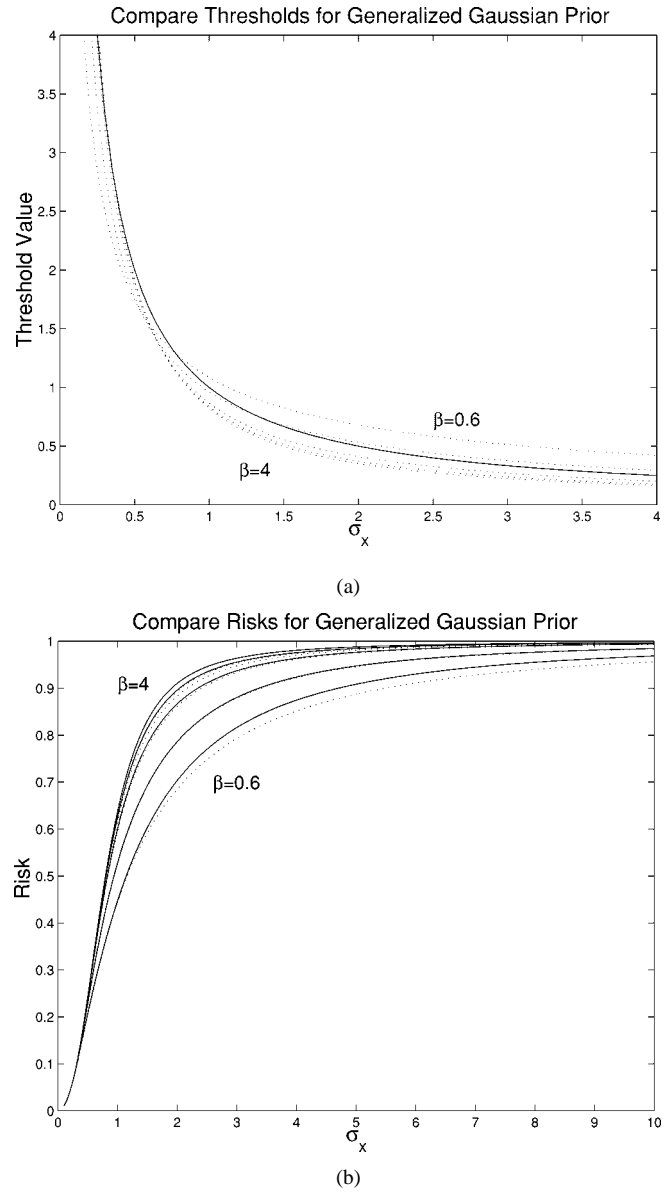


Fig. 6. Thresholding for the generalized Gaussian prior, with  $\sigma = 1$ . (a) Compare the approximation  $T_B(\sigma_X) = \sigma^2/\sigma_X$  (—) with the optimal threshold  $T^*(\sigma_X, \beta)$  for  $\beta = 0.6, 1, 2, 3, 4$  ( $\cdots$ ). The horizontal axis is the standard deviation,  $\sigma_X$ . (b) The optimal risks are in ( $\cdots$ ), and the approximation in (—).

only on the standard deviation and not on the shape parameter  $\beta$ , it may not yield a good approximation for values of  $\beta$  other than the range tested here, and the threshold may need to be modified to incorporate  $\beta$ . However, since for the wavelet coefficients typical values of  $\beta$  falls in the range  $[0.5, 1]$ , this simple form of the threshold  $T_B$  is appropriate for our purpose. For a fixed set of parameters, the curve of the risk (as a function of the threshold  $T$ ) is very flat near the optimal threshold  $T^*$ , implying that the error is not sensitive to a slight perturbation near  $T^*$ .

#### B. Parameter Estimation for Data-Driven Adaptive Threshold

This section focuses on the estimation of the GGD parameters,  $\sigma_X$  and  $\beta$ , which in turn yields a data-driven estimate of  $T_B(\sigma_X)$  that is adaptive to different subband characteristics.

The noise variance  $\sigma^2$  needs to be estimated first. In some situations, it may be possible to measure  $\sigma^2$  based on information other than the corrupted image. If such is not the case, it is estimated from the subband  $HH_1$  by the robust median estimator, also used in [14], [15],

$$\hat{\sigma} = \frac{\text{Median}(|Y_{ij}|)}{0.6745}, \quad Y_{ij} \in \text{subband } HH_1. \quad (16)$$

The parameter  $\beta$  does not explicitly enter into the expression of  $T_B(\sigma_X)$ , only the signal standard deviation,  $\sigma_X$ , does. Therefore it suffices to estimate directly  $\sigma_X$  or  $\sigma_X^2$ .

Recall the observation model is  $Y = X + V$ , with  $X$  and  $V$  independent of each other, hence

$$\sigma_Y^2 = \sigma_X^2 + \sigma^2 \quad (17)$$

where  $\sigma_Y^2$  is the variance of  $Y$ . Since  $Y$  is modeled as zero-mean,  $\sigma_Y^2$  can be found empirically by

$$\hat{\sigma}_Y^2 = \frac{1}{n^2} \sum_{i,j=1}^n Y_{ij}^2 \quad (18)$$

where  $n \times n$  is the size of the subband under consideration. Thus

$$\hat{T}_B(\hat{\sigma}_X) = \frac{\hat{\sigma}^2}{\hat{\sigma}_X} \quad (19)$$

where

$$\hat{\sigma}_X = \sqrt{\max(\hat{\sigma}_Y^2 - \hat{\sigma}^2, 0)}. \quad (20)$$

In the case that  $\hat{\sigma}^2 \geq \hat{\sigma}_Y^2$ ,  $\hat{\sigma}_X$  is taken to be 0. That is,  $\hat{T}_B(\hat{\sigma}_X)$  is  $\infty$ , or, in practice,  $\hat{T}_B(\hat{\sigma}_X) = \max(|Y_{ij}|)$ , and all coefficients are set to 0. This happens at times when  $\sigma$  is large (for example,  $\sigma > 20$  for a grayscale image).

To summarize, we refer to our method as *BayesShrink* which performs soft-thresholding, with the data-driven, subband-dependent threshold,

$$\hat{T}_B(\hat{\sigma}_X) = \hat{\sigma}^2 / \hat{\sigma}_X.$$

### III. MDL PRINCIPLE FOR COMPRESSION-BASED DENOISING: THE MDLQ CRITERION

Recall our hypothesis is that compression achieves denoising because the zero-zone in the quantization step (typical in compression methods) corresponds to thresholding in denoising. For the purpose of compression, after using the adaptive threshold  $\hat{T}_B(\hat{\sigma}_X)$  for the zero-zone, there still remains the questions of how to quantize the coefficients outside of the zero-zone and how to code them. Fig. 7 illustrates the block diagram of the compression method. It shows that the coder needs to decide on the design parameters  $m$ ,  $\Delta$  (the number of quantization

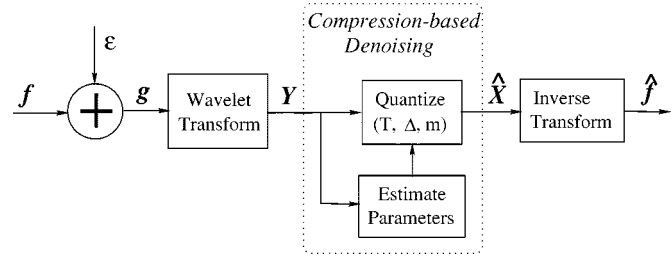


Fig. 7. Schematic for compression-based denoising. Denoising is achieved in the wavelet transform domain by lossy-compression, which involves the design of parameters  $T$ ,  $m$ , and  $\Delta$ , relating to the zero-zone width, the number of quantization levels, and the quantization binwidth, respectively.

bins and the binwidth, respectively), in addition to the zero-zone threshold  $\hat{T}_B$ . The choice of these parameters is discussed next.

When compressing a signal, two important objectives are to be kept in mind. On the one hand, the distortion between the compressed signal and the original should be kept low; on the other hand, the description of the compressed signal should use as few bits as possible to code. Typically, these two objectives are conflicting, thus a suitable criterion is needed to reach a compromise. Rissanen's MDL principle allows a tradeoff between these two objectives [26].

Let  $\mathcal{M}$  be a library or class of models from which the "best" one is chosen to represent the data. According to the MDL principle, given a sequence of observations, the "best" model is the one that yields the shortest description length for describing the data using the model, where the description length can be interpreted as the number of bits needed for encoding. This description can be accomplished by a two-part code: one part to describe the model and the other the description of the data using the model.

More precisely, given the set of observations  $\mathbf{Y}$ , we wish to find a model  $\hat{\mathbf{X}}$  to describe it. The MDL principle chooses  $\hat{\mathbf{X}}$  which minimizes the two-part code-length,

$$L(\mathbf{Y}, \hat{\mathbf{X}}) = L(\mathbf{Y}|\hat{\mathbf{X}}) + L(\hat{\mathbf{X}}) \quad (21)$$

where  $L(\mathbf{Y}|\hat{\mathbf{X}})$  is the code-length for  $\mathbf{Y}$  based on  $\hat{\mathbf{X}}$ , and  $L(\hat{\mathbf{X}})$  is the code-length for  $\hat{\mathbf{X}}$ .

In Saito's simultaneous compression and denoising method [28] for a length- $M$  one-dimensional signal, the hard-threshold function was used to generate the models  $\hat{\mathbf{X}} = \psi_T(\mathbf{Y})$ , where the number  $K$  of nonzero coefficients to retain is determined by minimizing the MDL criterion. The first term  $L(\mathbf{Y}|\hat{\mathbf{X}})$  is the idealized code-length with the normal distribution [see (23)], and the second term  $L(\hat{\mathbf{X}})$  is taken to be  $(3/2)K \log_2 M$ , of which  $K \log_2 M$  are the bits needed to indicate the location of each nonzero coefficient (assuming an uniform indexing) and  $(1/2)\log_2 M$  for the value of each of the  $K$  coefficients [see [26] for justification on using  $(1/2)\log_2 M$  bits to store the coefficient value]. Although compression has been achieved in the sense that a fewer number of nonzero coefficients are kept, [28] does not address the quantization step necessary in a practical compression setting.

In the following section, an MDL-based quantization criterion will be developed by minimizing  $L(\mathbf{Y}, \hat{\mathbf{X}})$  with the restriction that  $\hat{\mathbf{X}}$  belongs to the set of quantized signals.

### A. Derivation of the MDLQ Criterion

Consider a particular subband of size  $n \times n$ . Since the noisy wavelet transform coefficients are  $\mathbf{Y} = \mathbf{X} + \mathbf{V}$ , where  $V_{ij}$  are iid  $N(0, \sigma^2)$ , then  $Y_{ij}|X_{ij} \sim N(X_{ij}, \sigma^2)$ . Thus,

$$L(\mathbf{Y}|\mathbf{X}) = - \sum_{i,j=1}^n \log p(Y_{ij}|X_{ij}) \quad (22)$$

$$= \frac{1}{2\sigma^2 \log 2} \sum_{i,j=1}^n (Y_{ij} - X_{ij})^2 + \frac{1}{2} \log(2\pi\sigma^2 n^2). \quad (23)$$

The second term in (23) is a constant, and thus is ignored in the minimization. Expression (23) appears also in [21], [28]. The approach described here differs from theirs in the estimate  $\hat{\mathbf{X}}$ .

Let  $\mathcal{M}$  be the set of quantized coefficients  $\hat{\mathbf{X}}^Q$ , and  $\hat{\mathbf{X}}$  be constrained in  $\mathcal{M}$ . Plugging in  $\hat{\mathbf{X}}^Q$  as  $\mathbf{X}$  in (23) (with constant terms removed) gives

$$L(\mathbf{Y}|\hat{\mathbf{X}}^Q) = - \sum_{i,j=1}^n \log p(Y_{ij}|\hat{X}_{ij}^Q) \quad (24)$$

$$= \frac{1}{2\sigma^2 \log 2} \sum_{i,j=1}^n (Y_{ij} - \hat{X}_{ij}^Q)^2. \quad (25)$$

There are many possible ways to quantize and encode  $\hat{\mathbf{X}}$ . One way is the *uniform threshold quantizer* (UTQ) with centroid reconstruction based on the generalized Gaussian distribution. The parameters of the GGD can be estimated from the observed noisy coefficients as described below, which is a variant of that described in [29].

For noiseless observations,  $\sigma_X^2$  is estimated as

$$\hat{\sigma}_X^2 = \frac{1}{n^2} \sum_{i,j=1}^n X_{ij}^2 \quad (26)$$

and  $\beta$  is solved from

$$\kappa_X = \frac{\Gamma\left(\frac{1}{\beta}\right) \Gamma\left(\frac{5}{\beta}\right)}{\Gamma^2\left(\frac{3}{\beta}\right)} \quad (27)$$

where  $\kappa_X$  is the kurtosis of the GGD and is estimated as

$$\hat{\kappa}_X = \frac{1}{\hat{\sigma}_X^4} \frac{1}{n^2} \sum_{i,j=1}^n X_{ij}^4.$$

The parameter values listed in Fig. 3 are estimated this way.

When the image is corrupted by additive Gaussian noise, the second and fourth moments have the following relations:

$$\begin{aligned} \sigma_Y^2 &= \sigma_X^2 + \sigma^2, \\ \kappa_Y &= \frac{1}{\sigma_Y^4} \left( 6\sigma^2\sigma_Y^2 - 3\sigma^4 + (\sigma_Y^2 - \sigma^2)^2 \frac{\Gamma\left(\frac{1}{\beta}\right) \Gamma\left(\frac{5}{\beta}\right)}{\Gamma^2\left(\frac{3}{\beta}\right)} \right). \end{aligned} \quad (28)$$

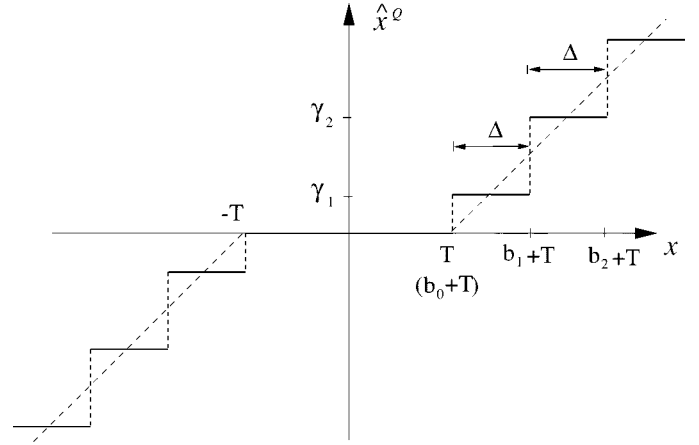


Fig. 8. Illustrating the quantizer.

The noise variance,  $\sigma^2$ , is estimated via (16). The second moment,  $\sigma_Y^2$ , and the kurtosis,  $\kappa_Y$ , can be measured from the observations  $\{Y_{ij}\}$ . The parameter  $\sigma_X$  is estimated as in (20) and  $\beta$  is then solved from (28).

Once the GGD parameters have been estimated, the quantizer has sufficient information to perform the quantization. The quantizer, shown in Fig. 8, consists of  $m$  levels of bins of equal size  $\Delta$  on each side, resulting in a total of  $2m + 1$  quantization bins (one zero-zone plus  $m$  symmetric levels on each positive and negative side). These bins are indexed as  $\ell = -m, \dots, -1, 0, 1, \dots, m$ . Consider the positive side and let  $b_0, b_1, \dots, b_m$  denote the boundaries of the quantization bins, with centroid reconstruction values  $\gamma_1, \gamma_2, \dots, \gamma_m$ . The value of  $\gamma_\ell$  with boundaries  $b_{\ell-1}$  and  $b_\ell$  is

$$\gamma_\ell = \frac{\int_{b_{\ell-1}}^{b_\ell} x G G_{\sigma_X, \beta}(x) dx}{\int_{b_{\ell-1}}^{b_\ell} G G_{\sigma_X, \beta}(x) dx}. \quad (29)$$

Equation (29) is calculated using numerical integration (e.g. the trapezoidal rule) since it does not have a closed-form solution. Note that  $b_0 = 0$ , and during quantization,  $b_m$  is taken to be  $\infty$ .

The negative side is quantized in a symmetric way. The quantized coefficients are denoted by  $\{\hat{X}_{ij}^Q\}$ . Note that the zero coefficients resulting from thresholding are kept as zeros, and that the subsequent quantization of the nonzero coefficients does not set any additional coefficients to zero. On average, the smallest number of bits needed to code  $\hat{\mathbf{X}}^Q$  is the Shannon code. Thus the code-length for coding the bin indices is

$$L(\hat{\mathbf{X}}^Q|m, \Delta) = - \sum_{\ell=-m}^m K_\ell \log \frac{K_\ell}{n^2} \quad (30)$$

where  $K_\ell$  is the number of coefficients in bin  $\ell$ . The additional parameters  $m$  and  $\Delta$  need to be coded also, but it is supposed that any positive values are equally likely, thus a fixed number of bits are allocated for  $L(m, \Delta)$  for each subband (8 bytes were used in the experiment).



Now we state the model selection criterion, *MDLQ*:

$$\begin{aligned} MDLQ(\mathbf{X}^Q, m, \Delta) &= \frac{1}{2\sigma^2 \log 2} \sum_{i,j=1}^n (Y_{ij} - \hat{X}_{ij}^Q)^2 + L(\hat{\mathbf{X}}^Q | m, \Delta), \\ \text{with } L(\hat{\mathbf{X}}^Q | m, \Delta) &= - \sum_{\ell=-m}^m K_\ell \log \frac{K_\ell}{n^2}. \end{aligned} \quad (31)$$

To find the best model, (32) is minimized over  $m$  and  $\Delta$  to find the corresponding set of quantized coefficients. In the implementation, which is by no means optimized, this is done by searching over  $m = 1, 2, \dots$ , and for each  $m$ , minimizing (32) over a reasonable range of  $\Delta$  [such as  $\Delta \in [\max(|Y_{ij}|)/(m+5), \max(|Y_{ij}|)/(m-1.5)]]$ . Typically, once a minimum (in  $m$ ) has been reached, the *MDLQ* cost increase monotonically, thus the search can terminate soon after a minimum has been detected.

This *MDLQ* compression with *BayesShrink* zero-zone selection is applied to each subband independently. The steps discussed in this section are summarized as follows.

- Estimate the noise variance  $\sigma^2$ , and the GGD standard deviation  $\sigma_X$ .
- Calculate the threshold  $\hat{T}_B$ , and soft-threshold the wavelet coefficients.
- To quantize the nonzero coefficients, minimize (32) over  $m$  and  $\Delta$  to find the corresponding quantized coefficients  $\hat{\mathbf{X}}^Q$ , which is the compressed, denoised estimate of  $\mathbf{X}$ .

The coarsest subband  $LL_J$  is quantized differently in that it is not thresholded, and the quantization with (32) assumes the uniform distribution. The  $LL_J$  coefficients are essentially local averages of the image, and are not characterized by distributions with a peak at zero, thus the uniform distribution is used for generality. With the mean subtracted, the uniform distribution is assumed to be symmetric about zero. Every quantization bin (including the zero-zone) is of width  $\Delta$ , and the reconstruction values are the midpoints of the intervals.

The *MDLQ* criterion in (32) has the additional interpretation of operating at a specified point on the rate-distortion (R-D) curve, as also pointed out by Liu and Moulin [21]. For a given coder, one can obtain a set of operational rate-distortion points  $(R, D)$ . When there is a rate or distortion constraint, the constraint problem can be formulated into a minimization problem with a Lagrange multiplier,  $\lambda D + R$ . In this case, (32) can be interpreted as operating at  $\lambda = (1/2\sigma^2 \log 2)$ . Natarajan [25] and Liu and Moulin [21] both proposed to use compression for denoising. The former operates at a constrained distortion,  $D \leq \sigma^2$ , and the latter operates at  $\lambda = (1/2\sigma^2 \log 2)$  on the R-D curve. Both works recommend the use of “any reasonable coder” while our coder is designed specifically with the purpose of denoising.

#### IV. EXPERIMENTAL RESULTS AND DISCUSSION

The  $512 \times 512$  grayscale images “goldhill,” “lena,” “barbara” and “baboon” are used as test images with different noise levels  $\sigma = 10, 20, 30, 35$ . The original images are shown in



Fig. 9. Original images. From top left, clockwise: goldhill, lena, barbara and baboon.

Fig. 9. The wavelet transform employs Daubechies’ least asymmetric compactly-supported wavelet with eight vanishing moments [11] with four scales of orthogonal decomposition.

To assess the performance of *BayesShrink*, it is compared with *SureShrink* [15]. The 1-D implementation of *SureShrink* can be obtained from the WaveLab toolkit [3], and the 2-D extension is straightforward. To gauge the best possible performance of a soft-threshold estimator, these methods are also benchmarked against what we call *OracleShrink*, which is the truly optimal soft-thresholding estimator assuming the original image is known. The threshold of *OracleShrink* in each subband is

$$T_{OS} = \arg \min_T \sum_{i,j=1}^n (\eta_T(Y_{ij}) - X_{ij})^2 \quad (32)$$

with  $X_{ij}$  known. To further justify the choice of soft-thresholding over hard-thresholding, another benchmark, *OracleThresh*, is also computed. *OracleThresh* is the best possible performance of a hard-threshold estimator, with subband-adaptive thresholds, each of which is defined as

$$T_{OT} = \arg \min_T \sum_{i,j=1}^n (\psi_T(Y_{ij}) - X_{ij})^2 \quad (33)$$

with  $X_{ij}$  known. The MSEs from the various methods are compared in Table I, and the data are collected from an average of five runs. The columns refer to, respectively, *OracleShrink*, *SureShrink*, *BayesShrink*, *BayesShrink* with MDLQ-based compression, *OracleThresh*, Wiener filtering, and the bitrate (in bpp, or bits-per-pixel) of the MDLQ-compressed image. Since the main benchmark is against *SureShrink*, the better one of the *SureShrink* and *BayesShrink* is highlighted in bold font for each test set. The MSEs resulting from *BayesShrink* comes to within

TABLE I  
FOR VARIOUS TEST IMAGES AND  $\sigma$  VALUES, LISTS MSE OF (1) *OracleShrink*, (2) *SureShrink*, (3) *BayesShrink*, (4) *BayesShrink* WITH MDLQ-COMPRESSION, (5) *OracleThresh*, AND (6) WIENER FILTERING THE LAST COLUMN SHOWS THE BITRATE (BITS PER PIXEL) OF THE COMPRESSED IMAGE OF (4). AVERAGED OVER FIVE RUNS

	<i>OracleShrink</i>	<i>SureShrink</i>	<i>BayesShrink</i>	<i>BayesShrink</i> + Compress	<i>OracleThresh</i>	Wiener	bitrate (bpp)
<i>goldhill</i>							
$\sigma=10$	41.28	42.26	<b>41.98</b>	59.21	54.11	42.87	1.055
$\sigma=20$	86.35	93.21	<b>88.59</b>	112.43	108.32	97.06	0.453
$\sigma=30$	124.80	128.98	<b>126.40</b>	155.60	152.22	189.75	0.270
$\sigma=35$	140.55	151.19	<b>141.97</b>	174.06	172.09	250.22	0.225
<i>lena</i>							
$\sigma=10$	28.31	<b>29.21</b>	29.65	40.71	34.93	28.52	0.747
$\sigma=20$	59.56	63.95	<b>61.73</b>	81.66	72.57	82.11	0.373
$\sigma=30$	89.74	94.13	<b>92.06</b>	119.37	110.16	175.36	0.234
$\sigma=35$	104.27	107.19	<b>106.45</b>	138.62	127.68	236.63	0.201
<i>barbara</i>							
$\sigma=10$	45.96	56.21	<b>51.27</b>	81.43	58.04	67.84	1.205
$\sigma=20$	118.11	<b>121.19</b>	121.52	170.20	150.94	136.12	0.852
$\sigma=30$	190.63	201.09	<b>192.60</b>	250.86	253.10	241.74	0.623
$\sigma=35$	226.30	246.07	<b>229.66</b>	289.77	300.11	308.51	0.516
<i>baboon</i>							
$\sigma=10$	66.37	85.76	<b>80.04</b>	129.71	90.27	115.15	1.343
$\sigma=20$	170.28	185.16	<b>180.15</b>	242.62	232.60	182.29	0.853
$\sigma=30$	263.09	<b>269.88</b>	270.18	332.66	349.62	284.14	0.583
$\sigma=35$	302.71	314.08	<b>309.92</b>	367.72	393.19	346.79	0.477

5% of *OracleShrink* for the smoother images *goldhill* and *lena*, and are most of the time within 6% for highly detailed images such as *barbara* and *baboon* (though it may suffer up to 20% for small  $\sigma$ ). *BayesShrink* outperforms *SureShrink* most of the time, up to approximately 8%. We observed in the experiments that using solely the SURE threshold yields excellent performance (sometimes yielding even lower MSE than *BayesShrink* by up to 1–2%). However, the hybrid method of *SureShrink* results at times in the choice of the universal threshold which can be too large. As illustrated in Table I, all three soft-thresholding methods outperforms significantly the best hard-thresholding rule, *OracleThresh*.

It is not surprising that the SURE threshold and the *BayesShrink* threshold yield similar performances. The SURE threshold can also be viewed as an approximate optimal soft-threshold in terms of MSE. For a particular subband of size  $n \times n$ , following [15],

$$\text{Sure}(T, Y) = n^2 - 2 \sum_{i,j=1}^n I_{\{|Y_{ij}| \leq T\}} + \sum_{i,j=1}^n (|Y_{ij}| \wedge T)^2 \quad (34)$$

where  $a \wedge b$  denotes  $\min(a, b)$ , and the SURE threshold is defined to be the value of  $T$  minimizing  $\text{Sure}(T, Y)$ .

Recall

$$Y_{ij} = X_{ij} + V_{ij}, \quad j = 1, \dots, n$$

and  $V_{ij}$ s are *iid*  $N(0, \sigma^2)$ . Conditioning on  $\mathbf{X} = \mathbf{x}$ , by Stein's result,

$$E_Y[\|\eta_T(Y) - \mathbf{x}\|^2 | \mathbf{X} = \mathbf{x}] = E_Y[\text{Sure}(T, Y) | \mathbf{X} = \mathbf{x}]. \quad (35)$$

Moreover, as we have done before, if the distribution of  $X$ s is approximated by a GGD, then the distribution of  $Y$ s is approximated by the mixture distribution of GGD and  $N(0, \sigma^2)$ ; or  $Y = X + V$  while  $X$  follows a GGD and is independent of  $V$ .

By the Law of Large Numbers,

$$\frac{1}{n^2} \text{Sure}(T, Y) \approx 1 - 2E_Y I_{\{|Y| \leq T\}} + E_Y (|Y| \wedge T)^2. \quad (36)$$

Taking expectation with respect to the GGD on both sides of (35), the risk can be written as

$$\begin{aligned} r(T) &= \frac{1}{n^2} E_Y \|\eta_T(Y) - \mathbf{x}\|^2 \\ &= \frac{1}{n^2} E_Y \text{Sure}(T, Y) \\ &= 1 - 2E_Y I_{\{|Y| \leq T\}} + E_Y (|Y| \wedge T)^2. \end{aligned} \quad (37)$$

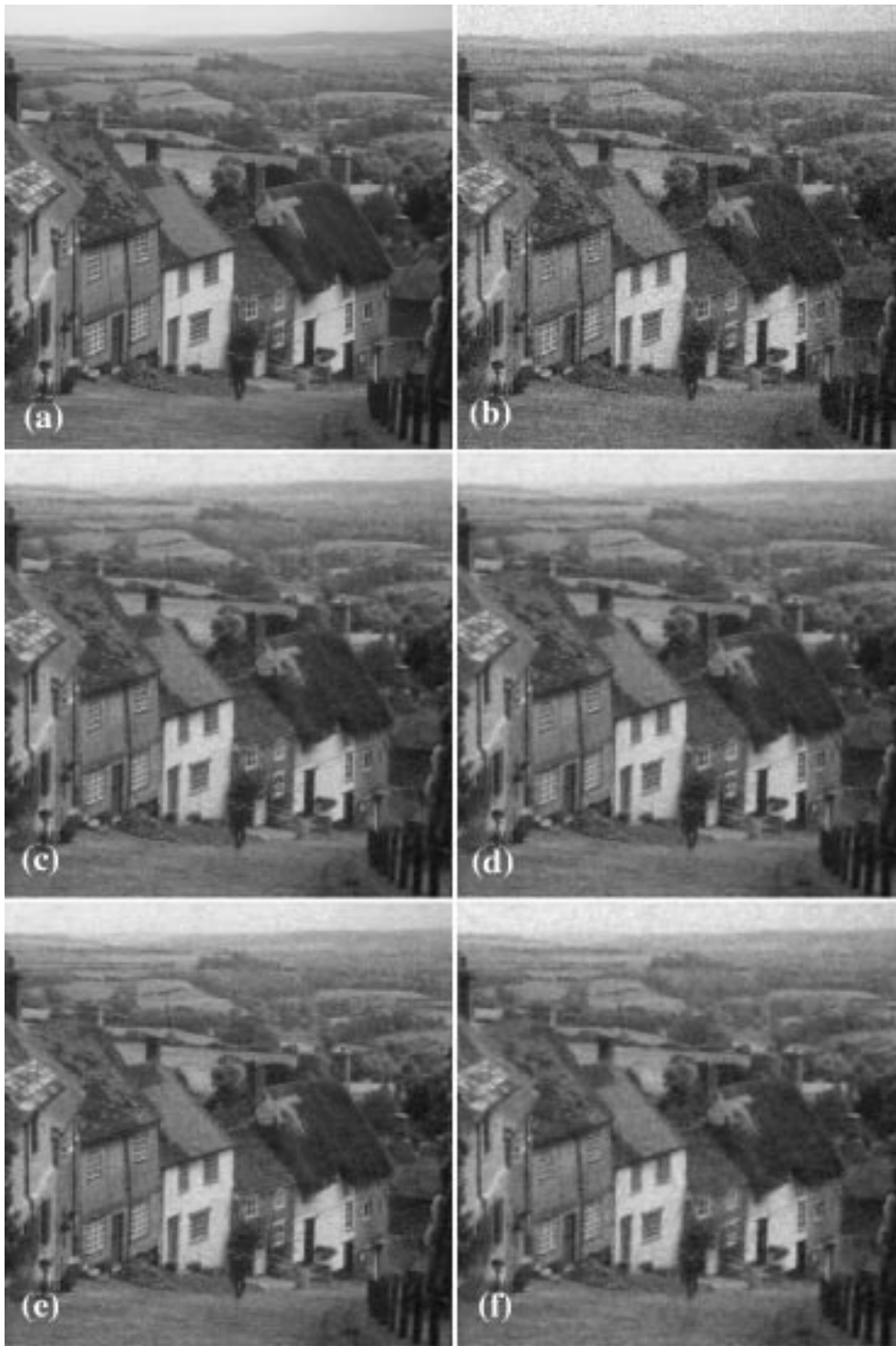


Fig. 10. Comparing the performance of the various methods on *goldhill* with  $\sigma = 20$ . (a) Original. (b) Noisy image,  $\sigma = 20$ . (c) *OracleShrink*. (d) *SureShrink*. (e) *BayesShrink*. (f) *BayesShrink* followed by MDLQ compression.

Comparing (37) with (36), one can conclude that  $(1/n^2)\text{Sure}(T, \mathbf{Y})$  is a data-based approximation to  $r(T)$ , and the SURE threshold, which minimizes  $\text{Sure}(T, Y)$ , is an alternative to *BayesShrink* for minimizing the Bayesian risk.

We have also made comparisons with the Wiener filter, the best linear filtering possible. The version used is the adaptive filter, *wiener2*, in the MATLAB image processing toolbox, using the default settings ( $3 \times 3$  local window size, and the

unknown noise power is estimated). The MSE results are shown in Table I, and they are considerably worse than the nonlinear thresholding methods, especially when  $\sigma$  is large. The image quality is also not as good as those of the thresholding methods.

The MDLQ-based compression step introduces quantization noise which is quite visible. As shown in the last column of Table I, the coder achieves a lower bitrate, but at the expense

of increasing the MSE. The MSE can be even worse than the noisy observation for small values of  $\sigma$ , especially for the highly detailed images. This is because the quantization noise is significant compared to the additive Gaussian noise. For larger  $\sigma$ , the compressed images can achieve noise reduction up to approximately 75% in terms of MSE. Furthermore, the bitrates are significantly less than the original 8 bpp for grayscale images. Thus, compression does achieve denoising and the proposed MDLQ-based compression can be used if simultaneous denoising and compression is a desired feature. If only the best denoising performance were the goal, obviously using solely *BayesShrink* is preferred.

Note that the first-order entropy coding,  $L(\hat{\mathbf{X}}^Q|m, \Delta)$ , for the bitrate of the quantized coefficients is a rather loose estimate. With more sophisticated coding methods (e.g. predictive coding, pixel classification), the same bitrate could yield a higher number of quantization level  $m$ , thus resulting in a lower MSE and enhancing the performance of the MDLQ-based compression-denoise.

A fair assessment of the MDLQ scheme for quantization after thresholding is the R-D curve used in Hansen and Yu [17] (see <http://cm.bell-labs.com/stat/binyu/publications.html>). This R-D curve is calculated using noiseless coefficients, and yields the best possible in terms of R-D tradeoff when the quantization is restricted to equal-binwidth. It thus gives an idea on how effective MDLQ is in choosing the tradeoff with respect to the optimal. The closeness of the MDLQ point to this R-D lower-bound curve indicates that MDLQ chooses a good R-D tradeoff without the knowledge of the noiseless coefficients required in deriving this R-D curve.

Fig. 10 shows the resulting images of each denoising method for *goldhill* and  $\sigma = 20$  (a zoomed-in section of the image is displayed in order to show the details). Table II compares the threshold values for each subband chosen by *OracleShrink*, *SureShrink* and *BayesShrink*, averaged over five runs. It is clear that the *BayesShrink* threshold selection is comparable to the SURE threshold and to the true optimal threshold  $T_{OS}$ . Some of the unexpectedly large threshold values in *SureShrink* comes from the universal threshold, not the SURE threshold, and these are placed in parentheses in the table. Table II(c) lists the thresholds of *BayesShrink*, and the thresholds in parentheses correspond to the case when  $\hat{T}_B = \max(|Y_{ij}|)$ , and all coefficients have been set to zero. Table III tabulates the values of  $m$  chosen by MDLQ for each subband of the *goldhill* image,  $\sigma = 20$ , averaged over five runs. The MDLQ criterion allocates more levels in the coarser, more important levels, as would be the case in a practical subband coding situation. A value of  $m = 0$  indicates that the coefficients have already been thresholded to zero, and there is nothing to code.

The results for *lena* and  $\sigma = 10$  are also shown. Fig. 11 shows the same sequences for a zoomed-in portion of *lena* with noise  $\sigma = 10$ . The corresponding results of threshold selections and MDLQ parameters for *lena* with noise  $\sigma = 10$  are listed in Tables IV and V. Interested readers can obtain a better view of the images at the website, <http://www-wavelet.eecs.berkeley.edu/~grchang/compressDenoise/>.

TABLE II  
THE THRESHOLD VALUES OF *OracleShrink*, *SureShrink*, AND *BayesShrink*, RESPECTIVELY, (AVERAGED OVER FIVE RUNS) FOR THE DIFFERENT SUBBANDS OF GOLDHILL, WITH NOISE STRENGTH  $\sigma = 20$

Scales	Orientations		
	HH	HL	LH
1 (finest)	84.90	37.36	38.04
2	37.12	21.16	18.50
3	17.31	11.51	8.45
4 (coarsest)	7.84	6.24	3.40

(a) *OracleShrink*:  $T_{OS}$

Scales	Orientations		
	HH	HL	LH
1 (finest)	(95.66)	(95.66)	(95.66)
2	89.48	21.78	18.86
3	18.52	11.46	9.14
4 (coarsest)	8.44	5.58	3.56

(b) *SureShrink*: SURE or (universal)

Scales	Orientations		
	HH	HL	LH
1 (finest)	(95.85)	49.69	51.34
2	46.37	19.78	16.84
3	17.43	8.45	6.91
4 (coarsest)	7.03	2.86	3.13

(c) *BayesShrink*:  $T_B$

TABLE III  
THE VALUE OF  $m$  (AVERAGED OVER FIVE RUNS) FOR THE DIFFERENT SUBBANDS OF GOLDHILL, WITH NOISE STRENGTH  $\sigma = 20$

Scales	Orientations			
	HH	HL	LH	LL
1 (finest)	0.0	4.0	4.2	
2	3.0	4.2	3.0	
3	3.4	5.2	5.0	
4 (coarsest)	5.0	15.4	9.8	27.4

## V. CONCLUSION

Two main issues regarding image denoising were addressed in this paper. Firstly, an adaptive threshold for wavelet thresholding images was proposed, based on the GGD modeling of subband coefficients, and test results showed excellent performance. Secondly, a coder was designed specifically for simultaneous compression and denoising. The proposed *BayesShrink* threshold specifies the zero-zone of the quantization step of this coder, and this zero-zone is the main agent in the coder which removes the noise. Although the setting in this paper was in the



Fig. 11. Comparing the performance of the various methods on *lena* with  $\sigma = 10$ . (a) Original. (b) Noisy image,  $\sigma = 10$ . (c) *OracleShrink*. (d) *SureShrink*. (e) *BayesShrink*. (f) *BayesShrink* followed by MDLQ compression.

wavelet domain, the idea can be extended to other transform domains such as DCT, which also relies on the energy compaction and sparse representation properties to achieve good compression.

There are several interesting directions worth pursuing. The current compression selects the threshold (i.e. zero-zone size)  $T_B$  and the quantization bin size  $\Delta$  in a two-stage process. In

typical image coders, however, the zero-zone is chosen to be either the same size or twice the size as other bins. Thus it would be interesting to jointly select these two values and analyze their dependencies on each other. Furthermore, a more sophisticated coder is likely to produce better compressed images than the current scheme, which uses the first order entropy to code the bin indices. With an improved coder, an increase in the number

TABLE IV  
THE THRESHOLD VALUES OF *OracleShrink*, *SureShrink*, AND *BayesShrink*,  
RESPECTIVELY, (AVERAGED OVER FIVE RUNS) FOR THE DIFFERENT  
SUBBANDS OF LENA, WITH NOISE STRENGTH  $\sigma = 10$

Scales	Orientations		
	HH	HL	LH
1 (finest)	25.73	14.13	19.16
2	11.21	7.66	10.18
3	6.15	3.58	5.89
4 (coarsest)	3.13	1.37	2.48

(a) *OracleShrink*:  $T_{OS}$

Scales	Orientations		
	HH	HL	LH
1 (finest)	(48.99)	15.93	(48.99)
2	12.54	8.57	11.23
3	6.65	4.03	6.62
4 (coarsest)	3.30	1.06	2.79

(b) *SureShrink*: SURE or (universal)

Scales	Orientations		
	HH	HL	LH
1 (finest)	(48.83)	15.72	32.20
2	9.93	4.50	7.67
3	3.10	1.60	2.88
4 (coarsest)	1.21	0.58	1.25

(c) *BayesShrink*:  $T_B$

TABLE V  
THE VALUE OF  $m$  (AVERAGED OVER FIVE RUNS) FOR THE DIFFERENT  
SUBBANDS OF LENA, WITH NOISE STRENGTH  $\sigma = 10$

Scales	Orientations			
	HH	HL	LH	LL
1 (finest)	0.0	5.2	5.6	
2	4.4	8.6	5.4	
3	9.0	11.6	8.0	
4 (coarsest)	18.6	31.8	14.2	45.0

of quantization bins would not increase the bitrate penalty by much, and thus the coefficients would be quantized at a finer resolution than the current method. Lastly, the model family  $\mathcal{M}$  could be expanded. For example, one could use a collection of wavelet bases for the wavelet decomposition, rather than using just one chosen wavelet, to allow possibly better representations of the signals.

In our other work [7], it was demonstrated that *spatially* adaptive thresholds greatly improves the denoising performance over uniform thresholds. That is, the threshold value changes for *each* coefficient. The threshold selection uses the context-modeling

idea prevalent in coding methods, thus it would be interesting to extend this spatially adaptive threshold to the compression framework, without incurring too much overhead. This would likely improve the denoising performance.

## REFERENCES

- [1] F. Abramovich, T. Sapatinas, and B. W. Silverman, "Wavelet thresholding via a Bayesian approach," *J. R. Statist. Soc., ser. B*, vol. 60, pp. 725–749, 1998.
- [2] M. Antonini, M. Barlaud, P. Mathieu, and I. Daubechies, "Image coding using wavelet transform," *IEEE Trans. Image Processing*, vol. 1, no. 2, pp. 205–220, 1992.
- [3] J. Buckheit, S. Chen, D. Donoho, I. Johnstone, and J. Scargle, "WaveLab Toolkit," <http://www-stat.stanford.edu:80/~wavelab/>.
- [4] A. Chambolle, R. A. DeVore, N. Lee, and B. J. Lucier, "Nonlinear wavelet image processing: Variational problems, compression, and noise removal through wavelet shrinkage," *IEEE Trans. Image Processing*, vol. 7, pp. 319–335, 1998.
- [5] S. G. Chang, B. Yu, and M. Vetterli, "Bridging compression to wavelet thresholding as a denoising method," in *Proc. Conf. Information Sciences Systems*, Baltimore, MD, Mar. 1997, pp. 568–573.
- [6] —, "Image denoising via lossy compression and wavelet thresholding," in *Proc. IEEE Int. Conf. Image Processing*, vol. 1, Santa Barbara, CA, Nov. 1997, pp. 604–607.
- [7] —, "Spatially adaptive wavelet thresholding with context modeling for image denoising," *IEEE Trans. Image Processing*, vol. 9, pp. 1522–1531, Sept. 2000.
- [8] H. Chipman, E. Kolaczyk, and R. McCulloch, "Adaptive bayesian wavelet shrinkage," *J. Amer. Statist. Assoc.*, vol. 92, no. 440, pp. 1413–1421, 1997.
- [9] M. Clyde, G. Parmigiani, and B. Vidakovic, "Multiple shrinkage and subset selection in wavelets," *Biometrika*, vol. 85, pp. 391–402, 1998.
- [10] M. S. Crouse, R. D. Nowak, and R. G. Baraniuk, "Wavelet-based statistical signal processing using hidden Markov models," *IEEE Trans. Signal Processing*, vol. 46, pp. 886–902, Apr. 1998.
- [11] I. Daubechies, *Ten Lectures on Wavelets, Vol. 61 of Proc. CBMS-NSF Regional Conference Series in Applied Mathematics*. Philadelphia, PA: SIAM, 1992.
- [12] R. A. DeVore and B. J. Lucier, "Fast wavelet techniques for near-optimal image processing," in *IEEE Military Communications Conf. Rec.* San Diego, Oct. 11–14, 1992, vol. 3, pp. 1129–1135.
- [13] D. L. Donoho, "De-noising by soft-thresholding," *IEEE Trans. Inform. Theory*, vol. 41, pp. 613–627, May 1995.
- [14] D. L. Donoho and I. M. Johnstone, "Ideal spatial adaptation via wavelet shrinkage," *Biometrika*, vol. 81, pp. 425–455, 1994.
- [15] —, "Adapting to unknown smoothness via wavelet shrinkage," *Journal of the American Statistical Assoc.*, vol. 90, no. 432, pp. 1200–1224, December 1995.
- [16] —, "Wavelet shrinkage: Asymptopia?," *J. R. Stat. Soc. B*, ser. B, vol. 57, no. 2, pp. 301–369, 1995.
- [17] M. Hansen and B. Yu, "Wavelet thresholding via MDL: Simultaneous denoising and compression," 1999, submitted for publication.
- [18] M. Jansen, M. Malfait, and A. Bultheel, "Generalized cross validation for wavelet thresholding," *Signal Process.*, vol. 56, pp. 33–44, Jan. 1997.
- [19] I. M. Johnstone and B. W. Silverman, "Wavelet threshold estimators for data with correlated noise," *J. R. Statist. Soc.*, vol. 59, 1997.
- [20] R. L. Joshi, V. J. Crump, and T. R. Fisher, "Image subband coding using arithmetic and trellis coded quantization," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 5, pp. 515–523, Dec. 1995.
- [21] J. Liu and P. Moulin, "Complexity-regularized image denoising," *Proc. IEEE Int. Conf. Image Processing*, vol. 2, pp. 370–373, Oct. 1997.
- [22] S. M. LoPresto, K. Ramchandran, and M. T. 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.
- [23] S. Mallat, "A theory for multiresolution signal decomposition: The wavelet representation," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 11, pp. 674–693, July 1989.
- [24] G. Nason, "Choice of the threshold parameter in wavelet function estimation," in *Wavelets in Statistics*, A. Antoniadis and G. Oppenheim, Eds. Berlin, Germany: Springer-Verlag, 1995.
- [25] B. K. Natarajan, "Filtering random noise from deterministic signals via data compression," *IEEE Trans. Signal Processing*, vol. 43, pp. 2595–2605, Nov. 1995.

- [26] J. Rissanen, *Stochastic Complexity in Statistical Inquiry*. Singapore: World Scientific, 1989.
- [27] F. Ruggeri and B. Vidakovic, "A Bayesian decision theoretic approach to wavelet thresholding," *Statist. Sinica*, vol. 9, no. 1, pp. 183–197, 1999.
- [28] N. Saito, "Simultaneous noise suppression and signal compression using a library of orthonormal bases and the minimum description length criterion," in *Wavelets in Geophysics*, E. Foufoula-Georgiou and P. Kumar, Eds. New York: Academic, 1994, pp. 299–324.
- [29] E. Simoncelli and E. Adelson, "Noise removal via Bayesian wavelet coring," *Proc. IEEE Int. Conf. Image Processing*, vol. 1, pp. 379–382, Sept. 1996.
- [30] C. M. Stein, "Estimation of the mean of a multivariate normal distribution," *Ann. Statist.*, vol. 9, no. 6, pp. 1135–1151, 1981.
- [31] M. Vetterli and J. Kovačević, *Wavelets and Subband Coding*. Englewood Cliffs, NJ: Prentice-Hall, 1995.
- [32] B. Vidakovic, "Nonlinear wavelet shrinkage with Bayes rules and Bayes factors," *J. Amer. Statist. Assoc.*, vol. 93, no. 441, pp. 173–179, 1998.
- [33] Y. Wang, "Function estimation via wavelet shrinkage for long-memory data," *Ann. Statist.*, vol. 24, pp. 466–484, 1996.
- [34] 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.
- [35] N. Weyrich and G. T. Warhola, "De-noising using wavelets and cross-validation," Dept. of Mathematics and Statistics, Air Force Inst. of Tech., AFIT/ENC, OH, Tech. Rep. AFIT/EN/TR/94-01, 1994.
- [36] 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 a 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 Dr.Sci. degree from EPF 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 Boveri 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.