

# A New SURE Approach to Image Denoising: Inter-scale Orthonormal Wavelet Thresholding

Florian Luisier, Thierry Blu and Michael Unser

## Abstract

This paper introduces a new approach to orthonormal wavelet image denoising. Instead of postulating a statistical model for the wavelet coefficients, we directly parametrize the denoising process as a sum of elementary nonlinear processes with unknown weights. We then minimize an estimate of the mean square error between the clean image and the denoised one.

The key point is that we have at our disposal a very accurate, statistically unbiased, MSE estimate—Stein’s Unbiased Risk Estimate (SURE)—that depends on the noisy image alone, not on the clean one. Like the MSE, this estimate is quadratic in the unknown weights and its minimization amounts to solving a linear system of equations. The existence of this *a priori* estimate makes it unnecessary to devise a specific statistical model for the wavelet coefficients. Instead, and contrary to the custom in the literature, these coefficients are not considered random anymore.

We describe an interscale orthonormal wavelet thresholding algorithm based on this new approach and show its near optimal performance—both regarding quality and CPU requirement—by comparing with the results of three state-of-the-art nonredundant denoising algorithms on a large set of test images. An interesting fallout of this study is the development of a new, group-delay based, parent-child prediction in a wavelet dyadic tree.

## Index Terms

Image denoising, orthonormal wavelet transform, SURE minimization, inter-scale dependencies.

The authors are with the Biomedical Imaging Group (BIG), Swiss Federal Institute of Technology (EPFL), CH-1015 Lausanne, Switzerland (e-mail: florian.luisier@epfl.ch; thierry.blu@epfl.ch; michael.unser@epfl.ch).

## I. INTRODUCTION

During acquisition and transmission, images are often corrupted by additive noise that can be modeled as Gaussian most of the time. The main aim of an image denoising algorithm is then to reduce the noise level, while preserving the image features. The multi-resolution analysis performed by the wavelet transform has been shown to be a powerful tool to achieve these goals. Indeed, in the wavelet domain, the noise is uniformly spread throughout the coefficients, while most of the image information is concentrated in the few largest ones (sparsity of the wavelet representation).

The most straightforward way of distinguishing information from noise in the wavelet domain consists in thresholding the wavelet coefficients. Of the various thresholding strategies, *soft-thresholding* is the most popular and has been theoretically justified by Donoho and Johnstone [1]. These authors have shown that the shrinkage rule is near-optimal in the minimax sense and provided the expression of the optimal threshold value  $T$ —called the “universal threshold”—as a function of the noise power  $\sigma^2$  when the number of samples  $N$  is large:  $T = \sqrt{2\sigma^2 \log N}$ . The use of the universal threshold to denoise images in the wavelet domain is known as *VisuShrink* [2].

Yet, despite its theoretical appeal, minimax is different from mean squared error (MSE) as a measure of error. A lot of work has been done to propose alternative thresholding strategies that behave better in terms of MSE than *VisuShrink* [3]–[6]. Donoho and Johnstone themselves acknowledged this flaw and suggested to choose the optimal threshold value  $T$  by minimizing Stein’s unbiased risk estimator (SURE) [7] when the data fail to be sparse enough for the minimax theory to be valid. This hybrid approach has been coined *SureShrink* by their authors [1]. Without challenging the soft-thresholding strategy, alternative threshold value selections have been proposed as well. One of the most popular was proposed by Chang *et al.*, who derived their threshold in a Bayesian framework, assuming a generalized Gaussian distribution for the wavelet coefficients. This solution to the wavelet denoising problem is known as *BayesShrink* [8] and has a better MSE performance than *SureShrink*.

Beyond the pointwise approach, more recent investigations have shown that substantially larger denoising gains can be obtained by considering the intra- and inter-scale correlations of the wavelet coefficients. In addition, increasing the redundancy of the wavelet transform is strongly beneficial to the denoising performances, a point to which we will come back later. We have selected three such techniques reflecting the state-of-the-art in wavelet denoising, against which we will compare our results:

- *Portilla et al.* [9]<sup>1</sup>: Their main idea is to model the neighborhoods of coefficients at adjacent positions and scales as a Gaussian scale mixture (GSM); the wavelet estimator is then a Bayes least squares (BLS). Their denoising method, consequently called *BLS-GSM*, is the most efficient up-to-date approach.
- *Pižurica et al.* [10]<sup>2</sup>: Assuming a generalized Laplacian prior for the noise-free data, their approach called *ProbShrink* is driven by the estimation of the probability that a given coefficient contains significant information—notation of “signal of interest”.
- *Sendur et al.* [11], [12]<sup>3</sup>: Their method, called *BiShrink*, is based on new non-Gaussian bivariate distributions to model inter-scale dependencies. A non-linear bivariate shrinkage function using the maximum a posteriori (MAP) estimator is then derived. In a second paper, these authors have extended their approach by taking into account the intra-scale variability of wavelet coefficients.

These techniques have been devised for both redundant and non-redundant transforms.

Despite reports on the superior denoising performances of redundant transforms [13], [14], we will only consider critically sampled wavelet transforms in this paper. The rationale behind our choice is that, since there is no *added* information—only *repeated* information—in redundant transforms, we believe that, eventually, a nonredundant transform may match the performance of redundant ones. This would potentially be very promising since the major drawback of redundant transforms are their memory and CPU time requirements which limits their routine use for very large images and, above all, usual volumes of data.

More than a specific denoising algorithm, this paper is about a powerful new method for optimizing *beforehand*—unaware of the clean image—the performance of a denoising method. Here, we want in particular to promote Stein’s Unbiased Risk Estimate (SURE) which is nothing less than an *a priori* estimation of the MSE resulting from an arbitrary processing of noisy data. This estimate turns out to be more accurate as more data are available, which is the case of images. Wavelet denoising methods routinely involve a statistical description of the coefficient distribution [15], an estimation of the—always nonlinear—statistical parameters and then, a search for the best denoising algorithm for this type of

<sup>1</sup>available at <http://decsai.ugr.es/~javier/denoise/software/index.htm>, with a  $3 \times 3$  neighborhood as suggested by the authors.

<sup>2</sup>available at <http://telin.ugent.be/~sanja/>, with a  $3 \times 3$  neighborhood and a threshold value  $T = \sigma$  as suggested by the authors.

<sup>3</sup>available at <http://taco.poly.edu/WaveletSoftware/denoise2.html>, with a  $7 \times 7$  neighborhood as suggested by the authors.

statistics. In contrast, by taking advantage of Stein’s MSE estimate, our method goes directly to the last step, without caring for the statistical description: in short, we do not make any explicit hypotheses on the clean image. In fact, we do not consider it as a random process at all; the randomness in our formulation follows from the Gaussian white noise alone.

Our approach consists thus in parametrizing the denoising method and choosing the parameters that minimize this MSE estimate. Previous techniques using the SURE required the minimization of complicated expressions for few nonlinear parameters [16], [17] or the use of parallel block iterative convex programming [18]. What makes our approach more tractable and efficient, is precisely the parametrizing method: a *linear combination* of nonlinear denoising functions—thresholding functions. Because of this “linear” choice, the minimization of the MSE estimate merely amounts to solving a linear system of equations, whose size is the number of weights in the linear combination. Obviously, the number of parameters, or degrees of freedom, is not a challenge and highly complicated thresholding behaviors can be obtained this way. In the context of image denoising, a univariate linear parametrization combined with an implicit SURE minimization was already evoked in [19] (*sigmoidal filtering*).

Because of the particular simplicity of Stein’s estimate for pointwise denoising functions, we will not exploit the full potential of the theory in this paper and will only consider interscale pointwise thresholding in the orthonormal wavelet transform. This excludes any intra-scale considerations. Yet, we will show that our denoising method performs better than the nonredundant versions of the state-of-the-art methods [9], [10], [12] on almost all tested images, to the noteworthy exception of *Barbara*, which may require intra-scale processing. Without any optimization attempts in our implementation, the comparison of computation times already show how economical our method is.

The paper is organized as follows: in Section II, we expose the SURE theory for functions of one or several statistically independent variables, and sketch the principles of our parametrization strategy; in Section III, we show how these principles can be exploited to build an efficient pointwise thresholding function that outperforms all known pointwise techniques; in Section IV, we extend the approach to a thresholding function that involves coarser scale parents as well. On this occasion, we develop a new formula to build a parent coefficient out of parent subbands; and finally, we compare our denoising method to the best available nonredundant techniques (section V). Both the competitiveness and robustness of our method validate our new approach as an attractive solution for image denoising.

## II. THEORETICAL ELEMENTS

### A. Problem setting

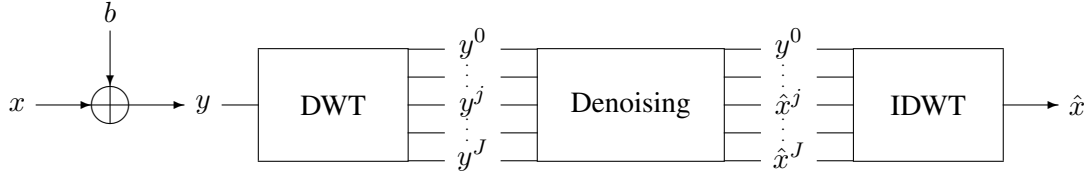


Fig. 1. Principle of wavelet denoising.

Wavelet denoising consists in three main stages (see Figure 1):

- i Perform a discrete wavelet transform (DWT) to the noisy data  $y = (y_n)_{n \in [1, N]}$  which is the sum of the noise-free data  $x = (x_n)_{n \in [1, N]}$  and the noise  $b = (b_n)_{n \in [1, N]}$ ;
- ii Denoise  $J$  noisy wavelet subimages  $y^j = x^j + b^j = (y_n^j)_{n \in [1, N_j]}$ ,  $j \in [1, J]$  by computing  $J$  estimates  $\hat{x}^j$  of the noise-free highpass subbands  $x^j$ ;
- iii Reconstruct the denoised image by applying the inverse discrete wavelet transform (IDWT) on the processed highpass wavelet subimages  $\hat{x}^j$  to obtain an estimate  $\hat{x}$  of the noise-free data  $x$ .

One can make two important remarks that set the context in which we will develop our denoising method:

- We will only consider additive Gaussian white noise following a normal law defined by a zero mean and a known<sup>4</sup>  $\sigma^2$  variance; i.e.  $b \sim \mathcal{N}(0, \sigma^2)$ .
- We will only consider *orthonormal* wavelet transform; the consequences are:
  - the mean-square error (MSE) in the space domain is a weighted sum of the MSE of each individual subband:

$$\underbrace{\langle |\hat{x} - x|^2 \rangle}_{\text{MSE}} = \sum_{j=0}^J \frac{N_j}{N} \underbrace{\langle |\hat{x}^j - x^j|^2 \rangle}_{\text{MSE}^j} \quad (1)$$

where we have introduced the notation:

$$\langle u \rangle = \frac{1}{N} \sum_{n=1}^N u_n \quad (2)$$

for the statistical mean estimate.

- the noise remains white and Gaussian with same statistics in the orthonormal wavelet domain, i.e.  $b^j \sim \mathcal{N}(0, \sigma^2)$ .

<sup>4</sup>In practice, the noise standard deviation can be accurately estimated using a robust median estimator [1].

This allows us to apply a new denoising function independently in every highpass subband, which means that our solution is subband-adaptive like most of the successful wavelet denoising approaches.

### B. Stein's unbiased MSE estimate (SURE)

In denoising applications, the performance is often measured in terms of peak signal-to-noise ratio (PSNR), which can be defined as follows<sup>5</sup>:

$$\text{PSNR} = 10 \log_{10} \left( \frac{\max(x^2)}{\langle |\hat{x} - x|^2 \rangle} \right) \quad (3)$$

Since the noise is a random process, we introduce an expectation operator  $\mathcal{E}\{\cdot\}$  to guess the potential results obtained after processing the noisy data  $y$ . Note that the noise-free data  $x$  is not modeled as a random process; thus:  $\mathcal{E}\{x\} = x$ .

The aim of image denoising is naturally to maximize the PSNR and thus to minimize the MSE defined in (1). In this paper, we choose to estimate each  $x^j$  by a pointwise function of  $y^j$ :

$$(\hat{x}_n^j)_{n \in [1, N_j]} = \left( \theta^j(y_n^j) \right)_{n \in [1, N_j]}$$

From now on, we will drop the subband index  $j$  since a new denoising function is independently applied in each individual subband. Our goal is to find a function  $\theta$  that minimizes

$$\text{MSE} = \langle |\theta(y) - x|^2 \rangle = \langle \theta(y)^2 \rangle - 2 \langle x\theta(y) \rangle + \langle x^2 \rangle \quad (4)$$

In practice, we only have access to the noisy signal  $y = x + b$ , and not to the original signal  $x$ . In (4), we thus need to remove the explicit dependence on  $x$ . Note that, since  $\langle x^2 \rangle$  has no influence in the minimization process, we do not need to estimate it. The remaining problematic term is only  $\langle x\theta(y) \rangle$ . However, the following theorem, a version of which was proposed by Stein in [7], allows us to overcome this difficulty:

*Theorem 1:* Let  $\theta : \mathbb{R} \rightarrow \mathbb{R}$  be a (weakly) differentiable function that does not explode at infinity<sup>6</sup>. Then, the following random variable:

$$\begin{aligned} \epsilon &= \langle \theta(y)^2 - 2y\theta(y) + 2\sigma^2\theta'(y) \rangle + \langle x^2 \rangle \\ &= \underbrace{\frac{1}{N} \sum_{n=1}^N \left( \theta^2(y_n) - 2y_n\theta(y_n) + 2\sigma^2\theta'(y_n) \right)}_{\tilde{\epsilon}} + \langle x^2 \rangle \end{aligned} \quad (5)$$

<sup>5</sup>For 8-bit images, usually  $\max(x^2) = 255^2$ .

<sup>6</sup>Typically, such that  $|\theta(z)| \leq \text{Const} \cdot \exp(az^2)$  for  $a < \frac{1}{2\sigma^2}$ .

is an unbiased estimator of the MSE, i.e.:

$$\mathcal{E}\{\epsilon\} = \mathcal{E}\{<|\theta(y) - x|^2>\}$$

**Proof:**

We can develop the square error between  $x_n$  and its estimate  $\theta(y_n)$  as:

$$\begin{aligned} \mathcal{E}\{|\theta(y_n) - x_n|^2\} &= \mathcal{E}\{\theta^2(y_n)\} - 2\mathcal{E}\{x_n\theta(y_n)\} + \mathcal{E}\{x_n^2\} \\ &= \mathcal{E}\{\theta^2(y_n)\} - 2\mathcal{E}\{y_n\theta(y_n)\} + 2\mathcal{E}\{b_n\theta(y_n)\} + x_n^2 \end{aligned}$$

where each term is well-defined thanks to the hypothesis on  $\theta$ .

We then use the fact that the Gaussian probability density  $q(b_n)$  satisfies  $b_n q(b_n) = -\sigma^2 q'(b_n)$  to evaluate  $\mathcal{E}\{b_n\theta(y_n)\}$ :

$$\begin{aligned} \mathcal{E}\{b_n\theta(y_n)\} &= \int \theta(x_n + b_n) b_n q(b_n) db_n \\ &= -\sigma^2 \int \theta(x_n + b_n) q'(b_n) db_n \\ &= \sigma^2 \int \theta'(x_n + b_n) q(b_n) db_n \quad (\text{by parts}) \\ &= \sigma^2 \mathcal{E}\{\theta'(y_n)\} \end{aligned} \tag{6}$$

Note that the integrated part  $\left[\sigma^2 \theta(x_n + b_n) q(b_n)\right]_{-\infty}^{+\infty}$  vanishes by hypothesis. This is known as Stein's Lemma [7] and leads to

$$\mathcal{E}\{|\theta(y_n) - x_n|^2\} = \mathcal{E}\{\theta^2(y_n)\} - 2\mathcal{E}\{y_n\theta(y_n)\} + 2\sigma^2 \mathcal{E}\{\theta'(y_n)\} + x_n^2.$$

Since the expectation of a sum is equal to the sum of the expectations, we can deduce that:

$$\mathcal{E}\{<|\theta(y) - x|^2>\} = \mathcal{E}\{<\theta^2(y)>\} - 2\mathcal{E}\{<y\theta(y)>\} + 2\sigma^2 \mathcal{E}\{<\theta'(y)>\} + <x^2>$$

■

As said before, there is no need to estimate  $<x^2>$ , since this term will disappear in the minimization. So, in practice, we will consider  $\tilde{\epsilon}$  which is the only part of the MSE estimate that depends on the choice of the denoising function  $\theta$ .

Note that Theorem 1 is still valid if  $\theta(y)$  is replaced by a two-variable denoising function  $\theta(y, z)$  where  $z$  is random, but independent<sup>7</sup> of  $y$ . In particular, in an orthonormal wavelet transform—which

<sup>7</sup>We recall that the randomness of  $y = x + b$  only results from the Gaussian white noise  $b$ , because no statistical model is assumed on the noise-free data  $x$ .

transforms Gaussian white noise into Gaussian white noise— $z$  can be any wavelet coefficients other than  $y$  itself.

The result given by Theorem 1 becomes particularly interesting in image denoising applications, where the number of samples is large. Indeed, by the law of large numbers, the standard deviation of  $\epsilon$  is small; i.e., the estimate  $\epsilon$  is close to its expectation which is the MSE of the denoising procedure. As a result, we can use  $\epsilon$  as if it were the true MSE. The next section shows how to use Theorem 1 efficiently.

### C. A SURE approach to image denoising

Our denoising approach amounts to minimizing  $\epsilon$  over a range of reasonable denoising functions  $\theta$ . We claim that this will result in the minimization of the MSE over the same range of functions, up to a small random error inversely proportional to the square root of the number of samples. Before defining more precisely which denoising functions we consider reasonable, we can illustrate the search for the optimal value  $T$  by applying Theorem 1 when  $\theta$  is the well-known *soft-thresholding* function defined by:

$$\theta(y) = \text{sign}(y)(|y| - T)_+ \quad (7)$$

where:  $(x)_+ = \max(x, 0)$ .

By Theorem 1, the following expression has to be minimized over  $T$ :

$$\tilde{\epsilon}(T) = \langle (2\sigma^2 + T^2 - y^2)(|y| - T)_+^0 \rangle \quad (8)$$

The last expression has its minimum exactly for the same  $T$  as the following formula:

$$\text{SURE}(T; y) = \sigma^2 - \frac{1}{N} \left( 2\sigma^2 \cdot \#\{n : |y_n| \leq T\} - \sum_{n=1}^N \min(|y_n|, T)^2 \right) \quad (9)$$

which appears in [1].

The estimated optimal threshold value is then:  $\tilde{T}_{\text{opt}} = \text{argmin}_T (\text{SURE}(T; y)) = \text{argmin}_T (\tilde{\epsilon}(T))$ .

We must notice here that the so-called *SureShrink* procedure developed by Donoho and Johnstone in [1] uses in fact an hybrid scheme between the SURE theory and the universal threshold (asymptotically optimal when the data exhibit a high level of sparsity). Their minimization of  $\text{SURE}(T; y)$  is thus restricted to  $T \in [0; T_{\text{univ}}]$ , where  $T_{\text{univ}} = \sqrt{2\sigma^2 \log N}$  is the universal threshold. Our opinion, however, is that this restriction is unnecessary—and often suboptimal—in image denoising applications where quality is measured by a mean-square criterium. This is because even though natural images have small wavelet coefficients, these are not vanishing as required by the strict sparsity results. It may even be argued that these small coefficients convey important texture information, and should thus not be set to zero.

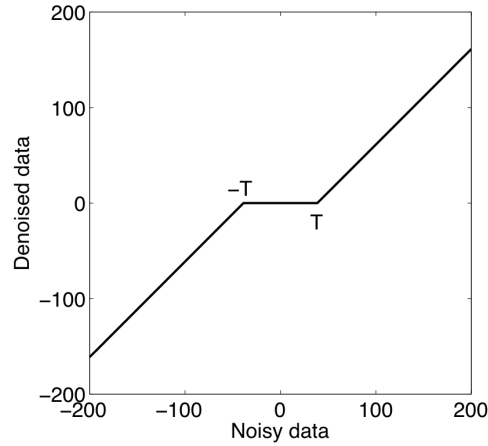


Fig. 2. The *soft-thresholding* function.

As we can verify on Figure 3, the estimate of Theorem 1 is statistically very reliable and robust, making it completely suitable for an accurate estimation of the optimal threshold.

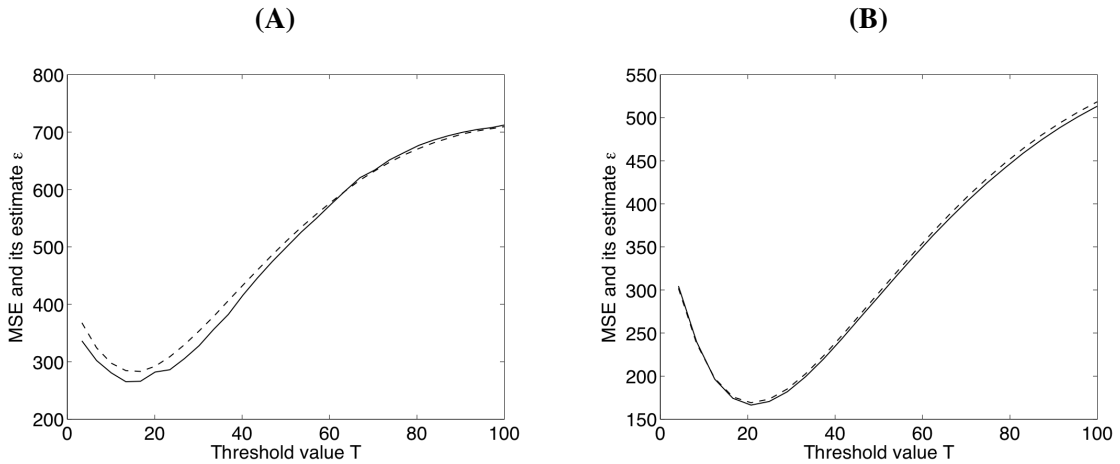


Fig. 3. Statistical accuracy of Theorem 1 illustrated with the *soft-threshold*: the true MSE is in dashed lines, while its estimate  $\epsilon$  is in solid line. (A)  $N = 32 \times 32$  samples and  $\sigma = 20$ . (B)  $N = 256 \times 256$  samples and  $\sigma = 20$ . The variance of the estimator decreases when the number of samples  $N$  increases, making Theorem 1 statistically reliable for image denoising applications.

The *soft-thresholding* function (see Figure 2) exhibits two main drawbacks: first, it only depends on a single parameter  $T$  and thus, its shape is not very flexible; second, this dependency is not linear. The consequence of these two remarks is that the sensitivity of the *soft-thresholding* function with respect to the value of  $T$  is high, and that finding the optimal threshold requires a nonlinear search algorithm.

In order to mitigate this issue, we choose to build a denoising function that depends *linearly* on a set of parameters—degrees of freedom—which we will determine exactly by minimizing  $\epsilon$ . The exact minimization is especially simple (linear) because the MSE estimate  $\epsilon$  has a quadratic form, much like the true MSE. The key idea is thus to build a linearly parameterized denoising function of the form:

$$\theta(y) = \sum_{k=1}^K a_k \varphi_k(y) \quad (10)$$

where  $K$  is the number of parameters.

If we introduce (10) into the estimate of the MSE given in Theorem 1 and perform differentiations over the  $a_k$ , we obtain for all  $k \in [1; K]$ :

$$\begin{aligned} 0 &= \frac{1}{2} \frac{\partial \epsilon}{\partial a_k} = \langle \theta(y) \varphi_k(y) - y \varphi_k(y) + \sigma^2 \varphi'_k(y) \rangle \\ &\quad \Updownarrow \\ \sum_{l=1}^K \underbrace{\langle \varphi_k(y) \varphi_l(y) \rangle}_{M_{k,l}} a_l - \underbrace{\langle y \varphi_k(y) - \sigma^2 \varphi'_k(y) \rangle}_{c_k} &= 0 \end{aligned}$$

These equations can be summarized in matrix form as  $\mathbf{M}\mathbf{a} = \mathbf{c}$ , where  $\mathbf{a} = [a_1 \dots a_K]^T$  and  $\mathbf{c} = [c_1 \dots c_K]^T$  are vectors of size  $K \times 1$ , and  $\mathbf{M} = [M_{k,l}]_{1 \leq k, l \leq K}$  is a matrix of size  $K \times K$ . This linear system is solved for  $\mathbf{a}$  by

$$\mathbf{a} = \mathbf{M}^{-1} \mathbf{c}, \quad (11)$$

which makes our approach very simple to implement. Note that, since we are only interested in the minimum of  $\epsilon$ , we are ensured that there will always be a solution. When several solutions are admissible (e.g., when  $\text{rank}(\mathbf{M}) < K$ ) any one of them will be acceptable—in particular, the one provided by the pseudoinverse of  $\mathbf{M}$ . When this degeneracy occurs, we will conclude that the parameters  $a_k$  belong to some linear subspace and thus, that some of them are useless (the function is “over-parameterized”). Of course, it is desirable to keep the number of degrees of freedom  $K$  as low as possible in order for the estimate  $\epsilon$  to keep a small variance.

### III. EFFICIENT SURE-BASED POINTWISE THRESHOLDING

In the previous section, we have proposed a general form of denoising functions (10). The difficulty is now to choose suitable basis functions  $\varphi_k$  that will determine the shape of our denoising function. Therefore, we want the denoising function  $\theta$  to satisfy the following properties:

- differentiability: required to apply Theorem 1;
- anti-symmetry: the wavelet coefficients are not expected to exhibit a sign preference;

- linear behavior for large coefficients: because  $\theta(y)$  should asymptotically tend to  $y$ .

After trying several types of  $\varphi_k$ , we have found that all of them give quite similar results, when the above conditions are satisfied. We have thus decided to retain the following pointwise denoising function:

$$\theta(y) = \sum_{k=1}^K a_k y e^{-(k-1) \frac{y^2}{2T^2}} \quad (12)$$

We choose derivatives of Gaussians (DOG) because they decay quite fast, which ensures a linear behavior close to the identity for large coefficients (see Figure 4).

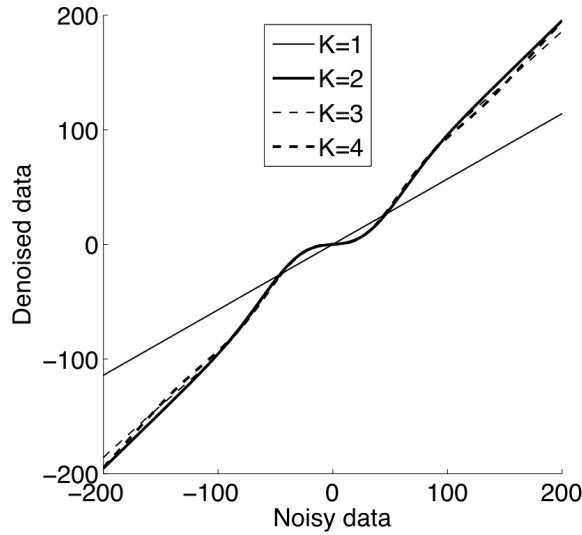


Fig. 4. The shape of our denoising function (12) in a particular subband, for various  $K$  and optimized  $a_k$ 's and  $T$ .

In addition to the linear coefficients, our denoising function contains two nonlinear dependencies: the number of terms  $K$  and the parameter  $T$ . We will see later that they can be fixed independently of the image.

If we consider only one parameter ( $K = 1$ ), our denoising function simply becomes  $\theta(y) = a_1 y$ , which is the simplest linear pointwise denoising function. The direct minimization of the estimate  $\epsilon$  provides

$$a_1 = 1 - \frac{\sigma^2}{\langle y^2 \rangle} \quad (13)$$

which is known as the James-Stein estimator [20].

Practical tests (with optimization over the parameter  $T$ , independently in each subband) on various images and with various noise levels have shown that, as soon as  $K \geq 2$ , the results become quite similar.

It thus appears that it is sufficient to keep as few as  $K = 2$  terms in (12). This is confirmed in Figure 4, which shows that the shape of our denoising function is nearly insensitive to the variation of  $K \geq 2$ .

Moreover, the optimal value of the parameter  $T$  is closely linked to the standard deviation  $\sigma$  of the noise and in a lesser way to the number of parameters  $K$ . Its interpretation is quite similar as in the case of the *soft-threshold*: it manages the transition between low SNR to high SNR coefficients. In our case though, the variations of the minimal  $\epsilon$  (over  $a_k$ ) when  $T$  changes are quite small (see Figure 5), because our denoising function is much more flexible than the *soft-threshold*. This sensitivity becomes even smaller as the number of parameters  $K$  increases. In fact, this indicates that some parameters are in that case useless.

To summarize, we have shown that both the number of terms  $K$  and the parameter  $T$  have only a minor influence on the quality of the denoising process. This indicates that these two parameters do not have to be optimized; instead, they can be fixed once for all, independently of the type of image. From a practical point of view, we suggest to use  $K = 2$  terms and  $T = \sqrt{6} \sigma$  (see Figure 5), leading to the following pointwise thresholding function:

$$\theta_0(y; \mathbf{a}) = \left( a_1 + a_2 e^{-\frac{y^2}{12\sigma^2}} \right) y \quad (14)$$

Now, it is interesting to evaluate the efficiency of our denoising function (14) and the accuracy of our minimization process based on an estimate  $\epsilon$  of the MSE. We propose to compare our results with the best results that can be reached by the popular *soft-threshold* with an optimal threshold choice (*OracleShrink*). Two main observations naturally come out of Table I :

- i SURE is a reliable estimate of the MSE, since the resulting average loss in PSNR is within 0.02 dB for all images.
- ii Our sum of DOG (14) gives better PSNRs than the optimal *soft-threshold*.

#### IV. EFFICIENT SURE-BASED INTER-SCALE THRESHOLDING

The integration of inter-scale information has been shown to improve the denoising quality, both visually and in terms of PSNR [9], [11], [21]. However, the gain brought is often modest, especially considering the additional complications involved by this processing [9]. In this section, we reformulate the problem by first building a loose prediction  $y_p$  of wavelet coefficients  $y$  out of a suitably filtered version of the lowpass subband at the same scale, and then by including this predictor in an explicit pointwise denoising function. Apart from the specific denoising problem addressed in this paper, we believe more generally that other applications (e.g. compression, detection, segmentation) could benefit as well from the theory that leads to this predictor.

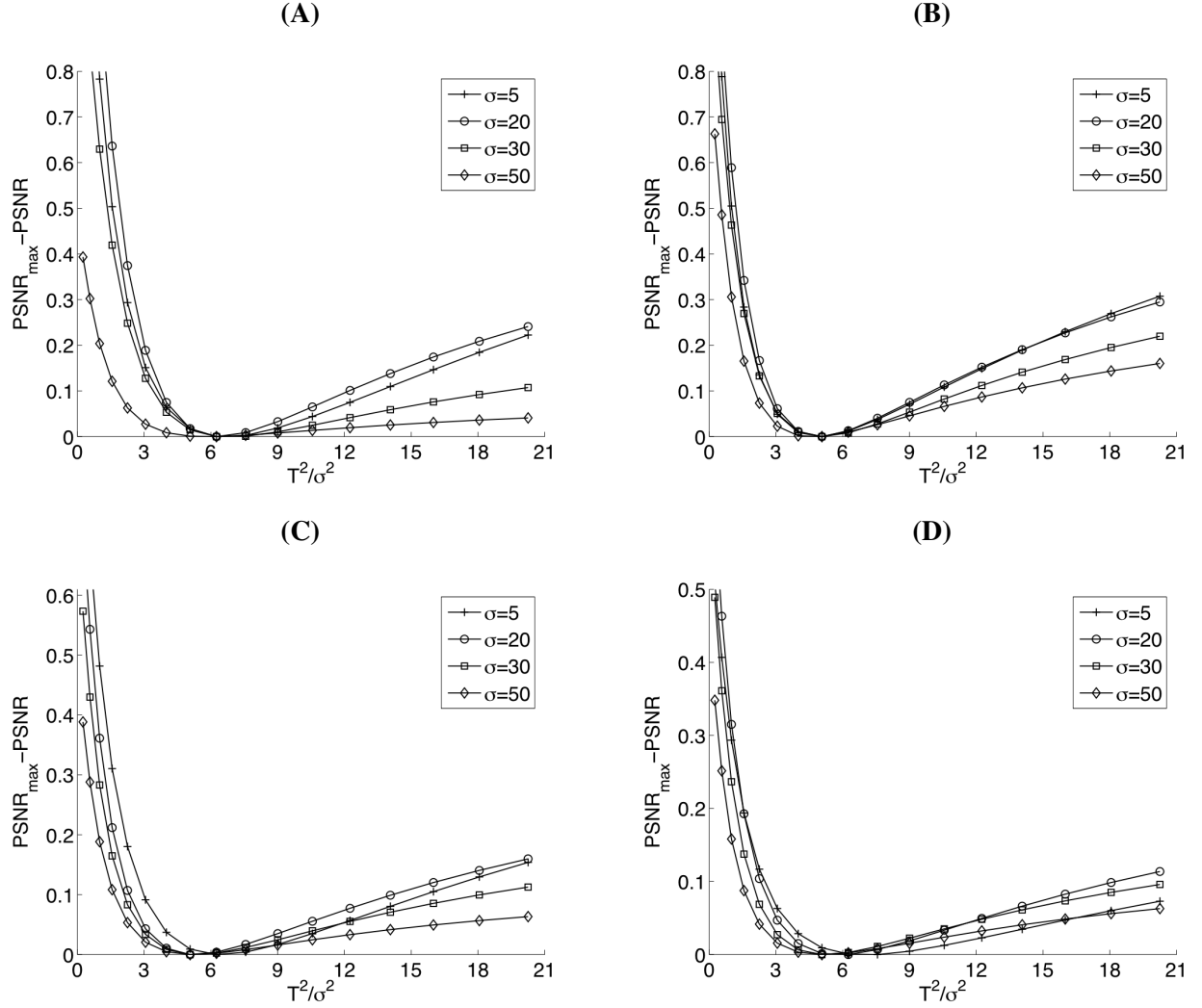


Fig. 5. Sensitivity of our denoising function (14) with respect to variations of  $T$ . (A) *Peppers*  $256 \times 256$ . (B) *MIT*  $256 \times 256$ . (C) *Lena*  $512 \times 512$ . (D) *Boat*  $512 \times 512$ . We can notice that for all images and for the whole range of input PSNR the maximum of the PSNR is reached for  $\frac{T^2}{\sigma^2} \simeq 6$ .

#### A. Building the inter-scale prediction

The wavelet coefficients that lie on the same dyadic tree (see Figure 6) are well-known to be large together in the neighborhood of image discontinuities. What can thus be predicted with reasonably good accuracy are the position of large wavelet coefficients out of parents at lower resolutions. However, getting the actual values of the finer resolution scale coefficients seem somewhat out of reach. This suggests that the best we can get out of between-scale correlations is a segmentation between regions of large and small coefficients. This comes back to the idea of signal of interest proposed by Pižurica *et al.* in [10].

TABLE I

COMPARISON OF OUR SUM OF DOG (14) WITH THE ORACLE *soft-threshold* (NON-REDUNDANT *sym8*, 4 ITERATIONS).

$\sigma$	5	10	20	30	50	5	10	20	30	50
Method	Boat $512 \times 512$					Goldhill $512 \times 512$				
<i>OracleShrink</i>	36.09	32.11	28.64	26.81	24.79	35.99	31.97	28.75	27.18	25.45
<b>Sum of DOG (Oracle)</b>	<b>36.35</b>	<b>32.37</b>	<b>28.85</b>	<b>27.03</b>	<b>25.01</b>	<b>36.21</b>	<b>32.25</b>	<b>28.99</b>	<b>27.42</b>	<b>25.67</b>
<b>Sum of DOG (SURE)</b>	<b>36.35</b>	<b>32.37</b>	<b>28.85</b>	27.02	25.00	<b>36.21</b>	<b>32.25</b>	<b>28.99</b>	27.41	25.66
Method	Peppers $256 \times 256$					Bridge $256 \times 256$				
<i>OracleShrink</i>	36.38	32.06	28.03	25.84	23.34	34.83	29.81	25.77	23.93	22.06
<b>Sum of DOG (Oracle)</b>	<b>36.67</b>	<b>32.36</b>	<b>28.28</b>	<b>25.97</b>	<b>23.47</b>	<b>34.89</b>	<b>30.00</b>	<b>26.10</b>	<b>24.29</b>	<b>22.40</b>
<b>Sum of DOG (SURE)</b>	<b>36.67</b>	32.35	28.27	25.95	23.45	<b>34.89</b>	<b>30.00</b>	26.09	24.28	22.39

Note: output PSNRs have been averaged over ten noise realizations.

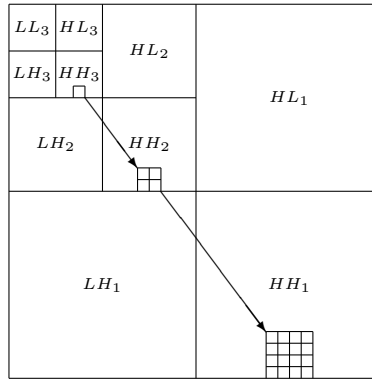


Fig. 6. Three stages of a fully decimated orthogonal wavelet transform and the so-called parent-child relationship.

In a critically sampled orthonormal wavelet decomposition, the parent subband is half-size the child subband. The usual way of putting the two subbands in correspondence is simply to expand the parent by a factor two. Unfortunately, this approach does not take into account the potential—noninteger—shift caused by the filters of the DWT. We thus propose a more sophisticated solution, which addresses this issue and ensures the alignment of image features between the child and its parent.

Our idea comes from the following observation: let  $LH_j$  and  $LL_j$  be, respectively, bandpass and lowpass outputs at iteration  $j$  of the filterbank. Then, if the *group delay*<sup>8</sup> between the bandpass and the lowpass filters are *equal*, no shift between the features of  $LH_j$  and  $LL_j$  will occur. Of course, depending on the amplitude response of the filters, some features may be attenuated, blurred, or enhanced, but their

<sup>8</sup>i.e., the frequency gradient of the phase response, with a minus sign.

location will remain unchanged. When the group delays differ, which is the general case, we thus propose to filter the lowpass subband  $LL_j$  in order to *compensate for the group delay difference* with  $LH_j$ . This operation is depicted in Figure 7 (A):  $LL_j$  is filtered in the three bandpass “directions” by adequately designed filters  $W_{HL}$ ,  $W_{HH}$  and  $W_{LH}$ , providing aligned—i.e., group delay compensated—subbands with  $HL_j$ ,  $HH_j$  and  $LH_j$ .

Because the filters considered in this paper are separable, we only have to consider 1D group delay compensation (GDC).

*Definition 1:* We say that two filters  $H(z)$  and  $G(z)$  are group delay compensated if and only if the group delay of the quotient filter  $H(z)/G(z)$  is zero identically; i.e., if and only if there exists a (anti-)symmetric filter  $R(z) = \pm R(z^{-1})$  such that  $H(z) = G(z)R(z)$ .

The following result shows how to choose a GDC filter in a standard orthonormal filterbank.

*Theorem 2:* For the output of the dyadic orthonormal filterbank of Figure 7 (B) to be group delay compensated, it is necessary and sufficient that:

$$W(z^2) = G(z^{-1})G(-z^{-1})(1 + \epsilon z^{-2})R(z^2) \quad (15)$$

where  $\epsilon = \pm 1$  and  $R(z) = R(z^{-1})$  is arbitrary.

**Proof:**

Group delay compensation between the two filterbank branches is equivalent to (see Figure 7 (B))

$$H(z^{-1})W(z^2) = G(z^{-1})R_1(z) \quad (16)$$

where  $R_1(z) = \epsilon R_1(z^{-1})$  is an arbitrary symmetric ( $\epsilon = 1$ ) or anti-symmetric ( $\epsilon = -1$ ) filter.

Because the filters  $H$  and  $G$  are orthonormal, we have that  $H(z^{-1}) = zG(-z)$ , and thus (16) can be rearranged as:

$$W(z^2) = \frac{G(z^{-1})R_1(z)}{zG(-z)} = G(z^{-1})G(-z^{-1}) \underbrace{\frac{z^{-1}R_1(z)}{G(-z)G(-z^{-1})}}_{R_2(z)} \quad (17)$$

We observe that  $R_2(z)$  is an even polynomial because both  $G(z^{-1})G(-z^{-1})$  and  $W(z^2)$  are. If we denote  $R_2(z) = (1 + \epsilon z^{-2})R(z^2)$ , then the symmetry of  $R_1(z)$  implies that

$$\begin{aligned} R(z^{-2}) &= \frac{zR_1(z^{-1})}{(1 + \epsilon z^2)G(-z)G(-z^{-1})} \\ &= \frac{\epsilon zR_1(z)}{(1 + \epsilon z^2)G(-z)G(-z^{-1})} \\ &= \frac{z^{-1}R_1(z)}{(1 + \epsilon z^{-2})G(-z)G(-z^{-1})} \\ &= R(z^2) \end{aligned}$$

i.e.,  $R(z)$  is an arbitrary zero-phase filter.

After substitution in (17), this finally leads us to the formulation (15), as an equivalent characterization of the group delay compensation in the filterbank of Figure 7 (B). ■

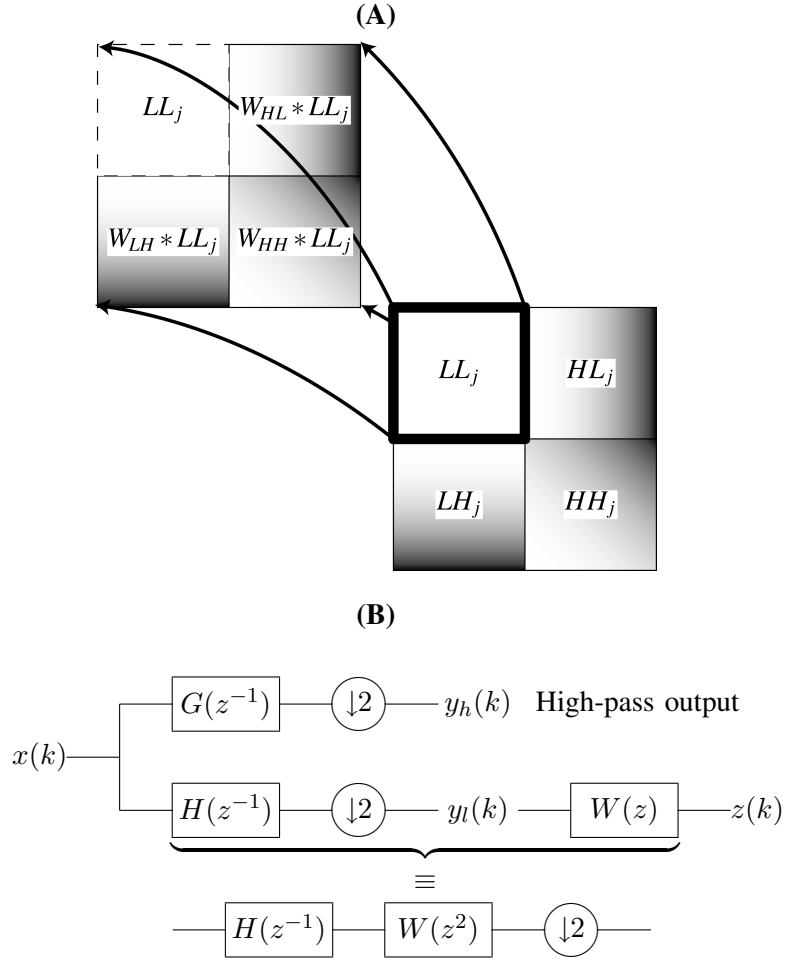


Fig. 7. One way of obtaining the whole parent information out of the lowpass subband: (A) 2D illustration. (B) 1D filterbank illustration.

In addition to (15), the GDC filter  $W(z)$  has to satisfy a few constraints:

- *Energy preservation*, i.e.,  $\sum_{n \in \mathbb{Z}} w_n^2 = 1$ , in order for the amplitude of the two outputs to be comparable;
- *Highpass behavior*, in order for the filtered lowpass image to “look like” the bandpass target;
- *Shortest possible response*, in order to minimize the enlargement of image features.

We can give a simple GDC filter in the case of symmetric filters. The shortest highpass  $W(z)$  satisfying the GDC condition is in fact the simple gradient filter:  $W(z) = z - 1$ . If the symmetry is not centered

at the origin but at a position  $n_0$ , then  $W(z) = z^{-n_0}(z - 1)$ . This type of solution is still adequate for near-symmetric filters such as the Daubechies *symlets* [22]. When the lowpass filter is not symmetric, we can simply take  $R(z^2) = 1$  in (15).

Finally, in order to increase the homogeneity inside regions of similar magnitude coefficients, we apply a 2D-smoothing filter—a normalized Gaussian kernel  $G(x) = \frac{1}{\sqrt{2\pi}}e^{-\frac{x^2}{2}}$ —onto the absolute value of the GDC output. In the rest of the paper, we will refer to the so-built inter-scale predictor by  $y_p$ .

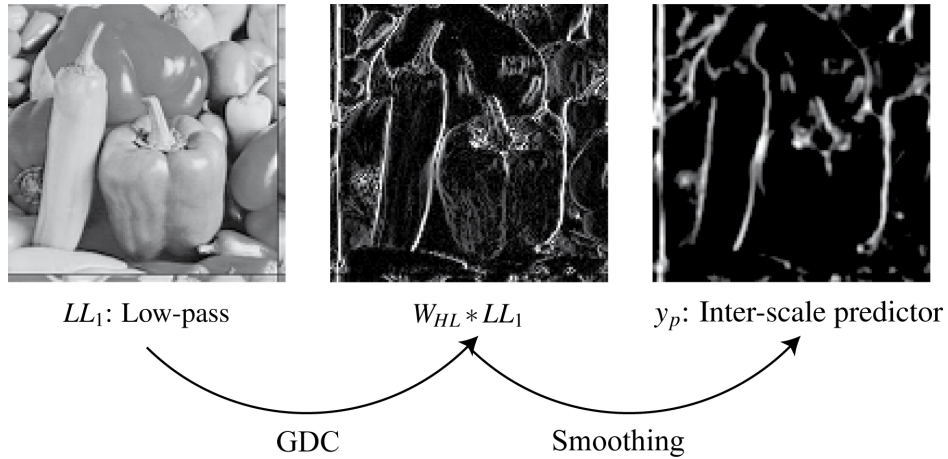


Fig. 8. Building an efficient interscale predictor, illustrated with a particular subband ( $HL_1$ ) of the noise-free *Peppers* image.

### B. Integrating the inter-scale predictor

Now that we have built the inter-scale predictor  $y_p$ , we have to suitably integrate it into our pointwise denoising function. As mentioned before, this inter-scale predictor does not tell us much about the actual value of its corresponding child wavelet coefficients. It only gives an indication on its expected magnitude. Here, we thus propose to use the parent  $y_p$  as a discriminator between high SNR wavelet coefficients and low SNR wavelet coefficients, leading to the following general pointwise denoising function:

$$\theta(y, y_p) = f(y_p) \sum_{k=1}^K a_k \varphi_k(y) + (1 - f(y_p)) \sum_{k=1}^K b_k \varphi_k(y) \quad (18)$$

The linear parameters  $a_k$  and  $b_k$  are then solved for by minimizing the MSE estimate  $\epsilon$  defined in Theorem 1, for the linear parameters  $a_k$  and  $b_k$ . The optimal coefficients are obtained in the same way as in section II-C and involve a solution similar to (11).

A first thought choice for the function  $f$  in (18) is simply the Heaviside function

$$H(y_p) = \begin{cases} 1, & \text{if } |y_p| \geq T \\ 0, & \text{if } |y_p| < T \end{cases} \quad (19)$$

where  $T$  can be interpreted as a decision factor. However, since the classification will not be perfect (i.e. some small parent coefficients may correspond to high magnitude child coefficients, and vice-versa), it is more appropriate to use a smoother decision function. We thus propose to use instead:

$$f(y_p) = e^{-\frac{y_p^2}{2T^2}} \quad (20)$$

As in the univariate case (section III), we suggest to use a sum of DOG with  $K = 2$  terms for each class of wavelet coefficients and<sup>9</sup>  $T = \sqrt{6} \sigma$ , leading to the following bivariate denoising function:

$$\begin{aligned} \theta(y, y_p; \mathbf{a}, \mathbf{b}) &= e^{-\frac{y_p^2}{12\sigma^2}} \theta_0(y; \mathbf{a}) + \left(1 - e^{-\frac{y_p^2}{12\sigma^2}}\right) \theta_0(y; \mathbf{b}) \\ &= e^{-\frac{y_p^2}{12\sigma^2}} \left(a_1 + a_2 e^{-\frac{y^2}{12\sigma^2}}\right) y + \left(1 - e^{-\frac{y_p^2}{12\sigma^2}}\right) \left(b_1 + b_2 e^{-\frac{y^2}{12\sigma^2}}\right) y \end{aligned} \quad (21)$$

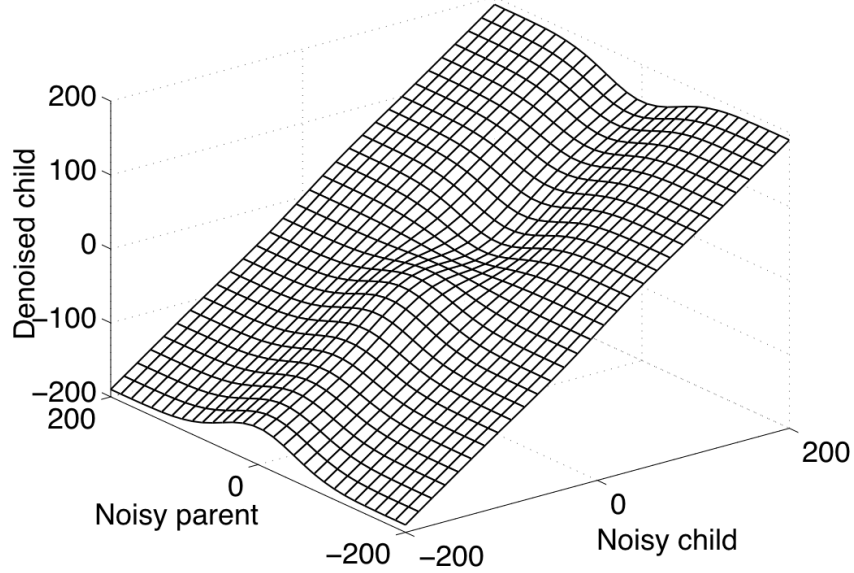


Fig. 9. 3D surface plot of a possible realization of our inter-scale thresholding function (21).

<sup>9</sup>Side investigations have shown that the  $T$  needed in (20) and the one optimized in section III can be chosen identical for optimal performances and equal to  $\sqrt{6} \sigma$ .

Table II quantifies the improvement introduced by this new way of integrating the inter-scale information, as compared to the usual expansion of the parent subband.

TABLE II  
DENOISING PERFORMANCE IMPROVEMENT BROUGHT BY OUR INTERSCALE STRATEGY (NON-REDUNDANT *sym8*, 4 ITERATIONS).

$\sigma$	5	10	20	30	50	100	5	10	20	30	50	100
Method	Peppers $256 \times 256$						House $256 \times 256$					
<i>Expansion by 2</i>	36.76	32.49	28.46	26.21	23.62	20.92	37.50	33.59	30.03	28.07	25.78	22.92
<i>Proposed</i>	<b>37.17</b>	<b>33.18</b>	<b>29.33</b>	<b>27.13</b>	<b>24.43</b>	<b>21.32</b>	<b>37.88</b>	<b>34.29</b>	<b>30.93</b>	<b>28.98</b>	<b>26.58</b>	<b>23.51</b>

Note: output PSNRs have been averaged over ten noise realizations.

## V. EXPERIMENTAL RESULTS

In this section, we compare our inter-scale dependent thresholding function (21) with some of the best state-of-the-art techniques: Sendur's *et al.* bivariate MAP estimator with local variance estimation, Portilla's *BLS-GSM* and Pižurica's *ProbShrink*.

In all comparisons, we use a critically sampled orthonormal wavelet basis with eight vanishing moments (*sym8*) over four decomposition stages.

### A. PSNR comparisons

We have tested the various denoising methods for a representative set of standard 8-bit grayscale images such as *Al*, *Barbara*, *Boat*, *Crowd*, *Goldhill* (size  $512 \times 512$ ) and *Peppers*, *House*, *Bridge* (size  $256 \times 256$ ), corrupted by simulated additive Gaussian white noise at eight different power levels  $\sigma \in [5, 10, 15, 20, 25, 30, 50, 100]$ , which corresponds to PSNR decibel values  $[34.15, 28.13, 24.61, 22.11, 20.17, 18.59, 14.15, 8.13]$ . The denoising process has been performed over ten different noise realizations for each standard deviation and the resulting PSNRs averaged over these ten runs. The parameters of each method have been set according to the values given by their respective authors in the corresponding referred papers. Variations in output PSNRs are thus only due to the denoising techniques themselves. This reliable comparison was only possible thanks to the kindness of the various authors who have provided their respective Matlab codes on their personal websites.

Table III summarizes the results obtained. To the noteworthy exception of *Barbara*, our results are already competitive with the best techniques available that consider non-redundant orthonormal transforms.

We stress again that our processing consists in a simple pointwise threshold, driven by interscale information; i.e., without taking intra-scale dependencies into consideration, contrary to the best performing methods (*ProbShrink*, *BiShrink* and *BLS-GSM*).

When looking closer at the results, we observe that:

- Our method outperforms the classical *BayesShrink* by more than +1 dB on average.
- Our method gives better results than Sendur's *BiShrink*  $7 \times 7$  which integrates both the inter- and the intra-scale dependencies (average gain of +0.6 dB).
- Our method gives better results than Pižurica's *ProbShrink*  $3 \times 3$  which integrates the intra-scale dependencies (average gain of +0.4 dB).
- We obtain similar or sometimes even better results than Portilla's *BLS-GSM*  $3 \times 3$  for most of the images.
- For the *Barbara* image, our method is among the worst performers together with the pointwise *BayesShrink*. Our explanation for this is that some local information (especially the texture in Barbara's trousers) is completely lost at coarser scales (see Figure 10). Inter-scale correlations may be too weak for this image, which indicates that an efficient denoising process may require intrascale information as well.
- The gap between our non-redundant SURE-based approach and the best up-to-date redundant results lies in the range of 0.5-1 dB for most images.

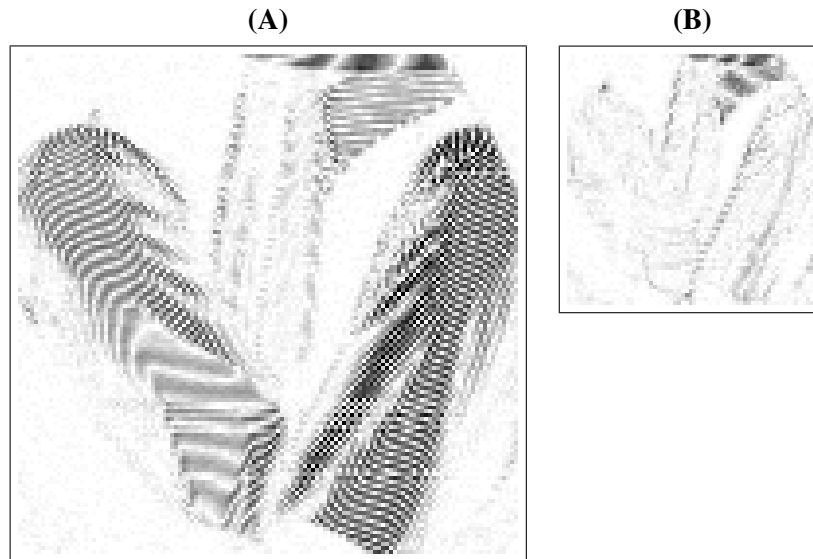


Fig. 10. (A) A zoom at Barbara's trousers at the finest scale of an orthonormal wavelet transform: the stripes are clearly visible. (B) A zoom at Barbara's trousers at the next coarser scale: the stripes are not visible anymore.

TABLE III

COMPARISON OF SOME OF THE MOST EFFICIENT DENOISING METHODS (NON-REDUNDANT  $\text{sym}8$ , 4 ITERATIONS).

$\sigma$	5	10	15	20	25	30	50	100	5	10	15	20	25	30	50	100
<b>Input PSNR</b>	34.15	28.13	24.61	22.11	20.17	18.59	14.15	8.13	34.15	28.13	24.61	22.11	20.17	18.59	14.15	8.13
<b>Method</b>	<b>Peppers <math>256 \times 256</math></b>								<b>House <math>256 \times 256</math></b>							
<i>BayesShrink</i>	35.83	31.49	29.30	27.85	26.72	25.73	23.17	20.73	36.91	32.92	30.81	29.42	28.44	27.66	25.49	22.87
<i>BiShrink <math>7 \times 7</math></i>	36.61	32.55	30.25	28.66	27.47	26.51	23.89	20.80	37.54	33.60	31.56	30.16	29.07	28.20	25.83	22.84
<i>ProbShrink <math>3 \times 3</math></i>	36.72	32.68	30.41	28.85	27.67	26.70	23.85	20.85	37.59	33.84	31.74	30.29	29.20	28.35	25.99	23.17
<i>BLS-GSM <math>3 \times 3</math></i>	36.80	32.86	30.62	29.07	27.90	26.97	24.40	20.88	<b>38.01</b>	34.26	32.23	30.79	29.65	28.72	26.15	22.97
<b>Our bivariate (21)</b>	<b>37.17</b>	<b>33.18</b>	<b>30.91</b>	<b>29.33</b>	<b>28.12</b>	<b>27.13</b>	<b>24.43</b>	<b>21.32</b>	37.88	<b>34.29</b>	<b>32.32</b>	<b>30.93</b>	<b>29.86</b>	<b>28.98</b>	<b>26.58</b>	<b>23.51</b>
<i>Best redundant</i>	37.09	33.33	31.26	29.84	28.74	27.84	25.30	21.98	38.43	35.04	33.23	31.91	30.87	30.01	27.62	24.53
<b>Method</b>	<b>AI <math>512 \times 512</math></b>								<b>Bridge <math>256 \times 256</math></b>							
<i>BayesShrink</i>	37.77	34.17	32.10	30.67	29.63	28.84	26.67	23.84	34.81	29.80	27.30	25.75	24.69	23.90	22.04	19.99
<i>BiShrink <math>7 \times 7</math></i>	38.01	34.50	32.57	31.23	30.21	29.39	27.09	24.01	34.94	29.93	27.38	25.81	24.75	23.97	22.11	19.97
<i>ProbShrink <math>3 \times 3</math></i>	38.11	34.58	32.64	31.28	30.08	29.32	27.18	24.24	34.59	29.61	27.20	25.74	24.73	23.97	22.10	20.08
<i>BLS-GSM <math>3 \times 3</math></i>	38.38	34.83	32.93	31.58	30.53	29.68	27.35	24.20	34.98	29.98	27.50	26.02	25.01	24.25	22.34	20.00
<b>Our bivariate (21)</b>	<b>38.43</b>	<b>34.90</b>	<b>32.97</b>	<b>31.64</b>	<b>30.64</b>	<b>29.84</b>	<b>27.61</b>	<b>24.56</b>	<b>35.06</b>	<b>30.22</b>	<b>27.84</b>	<b>26.36</b>	<b>25.33</b>	<b>24.56</b>	<b>22.60</b>	<b>20.35</b>
<i>Best redundant</i>	38.90	35.46	33.66	32.42	31.46	30.67	28.46	25.51	35.23	30.46	28.07	26.60	25.58	24.83	22.98	20.78
<b>Method</b>	<b>Barbara <math>512 \times 512</math></b>								<b>Boat <math>512 \times 512</math></b>							
<i>BayesShrink</i>	35.78	31.25	28.86	27.32	26.22	25.34	23.14	21.36	35.99	31.98	29.94	28.55	27.52	26.71	24.74	22.44
<i>BiShrink <math>7 \times 7</math></i>	36.76	32.52	30.14	28.51	27.29	26.33	23.91	21.47	36.18	32.46	30.47	29.08	28.03	27.20	25.05	22.52
<i>ProbShrink <math>3 \times 3</math></i>	36.75	32.48	30.04	28.40	27.20	26.27	23.86	21.58	36.20	32.53	30.50	29.11	28.05	27.22	25.12	22.69
<i>BLS-GSM <math>3 \times 3</math></i>	<b>37.05</b>	<b>32.89</b>	<b>30.54</b>	<b>28.93</b>	<b>27.72</b>	<b>26.76</b>	<b>24.25</b>	21.53	36.46	32.89	<b>30.89</b>	<b>29.49</b>	28.43	27.58	25.34	22.64
<b>Our bivariate (21)</b>	36.71	32.18	29.66	27.98	26.76	25.83	23.70	<b>21.76</b>	<b>36.70</b>	<b>32.90</b>	30.85	29.47	<b>28.44</b>	<b>27.63</b>	<b>25.50</b>	<b>22.97</b>
<i>Best redundant</i>	37.69	33.90	31.71	30.16	28.96	27.99	25.32	22.47	36.94	33.53	31.64	30.32	29.30	28.48	26.28	23.65
<b>Method</b>	<b>Crowd <math>512 \times 512</math></b>								<b>Goldhill <math>512 \times 512</math></b>							
<i>BayesShrink</i>	34.60	29.31	26.53	24.73	23.45	22.47	20.07	17.46	35.93	31.94	29.96	28.69	27.79	27.13	25.41	23.32
<i>BiShrink <math>7 \times 7</math></i>	34.71	29.48	26.70	24.88	23.57	22.57	20.13	17.40	36.17	32.27	30.32	29.07	28.15	27.44	25.57	23.26
<i>ProbShrink <math>3 \times 3</math></i>	34.42	29.29	26.59	24.83	23.56	22.58	20.15	17.43	36.07	32.30	30.35	29.07	28.13	27.43	25.62	23.47
<i>BLS-GSM <math>3 \times 3</math></i>	34.79	29.63	26.91	25.12	23.84	22.85	20.39	17.51	36.37	32.61	30.68	29.41	28.47	27.73	25.73	23.30
<b>Our bivariate (21)</b>	<b>34.86</b>	<b>29.77</b>	<b>27.11</b>	<b>25.38</b>	<b>24.13</b>	<b>23.17</b>	<b>20.75</b>	<b>17.97</b>	<b>36.53</b>	<b>32.69</b>	<b>30.76</b>	<b>29.52</b>	<b>28.60</b>	<b>27.89</b>	<b>26.06</b>	<b>23.82</b>
<i>Best redundant</i>	34.96	30.05	27.49	25.83	24.63	23.69	21.24	18.33	36.88	33.24	31.37	30.13	29.22	28.50	26.60	24.30

Note: output PSNRs have been averaged over ten noise realizations. The best redundant results are obtained using the *BLS-GSM  $3 \times 3$*  with an 8-orientations full steerable pyramid; results slightly differ from the ones published in [9], because no boundary extension has been applied here.

It is instructive to compare the results (see Figure 11) obtained with our inter-scale dependent thresholding function (21), with the ones obtained with our simple univariate denoising function (14). The improvement (often more than +1 dB) is quite significant for most standard images (see Figure 11). Yet, for images that have a substantial high frequency contents the integration of inter-scale dependencies does not lead to such an impressive gain. On the same graphs, we have also included the results obtained with the *OracleShrink*, showing a systematic underperformance with regards to even our simple univariate denoising function.

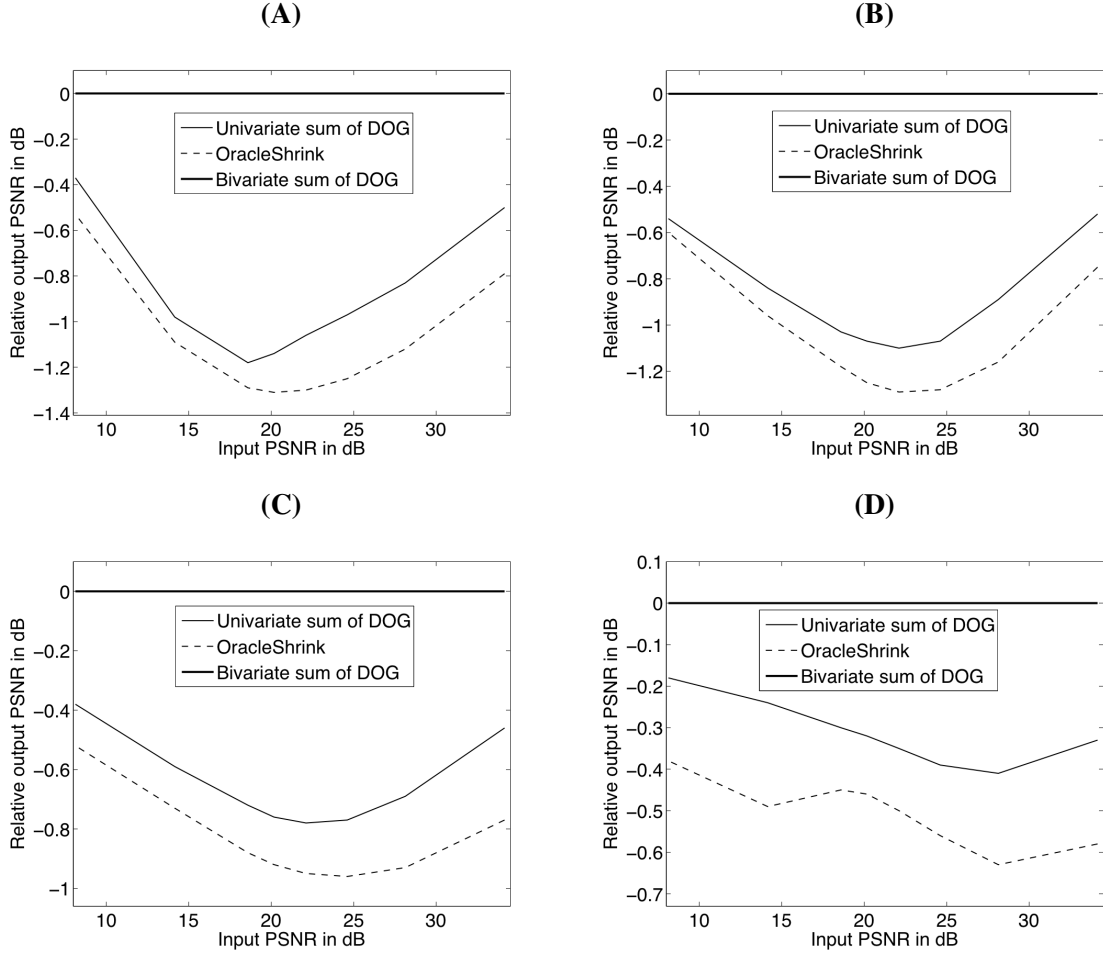


Fig. 11. Comparison of our inter-scale dependent thresholding function (21) with the best possible soft-threshold *OracleShrink* and with our simple univariate denoising function (14). (A) *Peppers* 256 × 256. (B) *House* 256 × 256. (C) *Lena* 512 × 512. (D) *Barbara* 512 × 512.

### B. Visual quality

Although there is no consensual objective way to judge the visual quality of a denoised image, two important criteria are widely used: the visibility of processing artifacts and the conservation of image edges. Processing artifacts usually result from a modification of the spatial correlation between wavelet coefficients (often caused by the zeroing of small neighboring coefficients) and are likely to be reduced by taking into account intra-scale dependencies. Instead, image edge distortions usually arise from modifications of the interscale coefficient correlations. The amplitude of these modifications is likely to be reduced by a careful consideration of interscale dependencies in the denoising function.

Since our algorithm only includes interscale considerations, we expect it to be specifically robust to

noise with regards to edge preservation. Additionally, we would like to stress that our method exhibits the fewest number of artifacts, which we attribute to the fact that we are never forcing any wavelet coefficients to zero. These observations are illustrated in Figures 12 and 13.

### C. Computation time

It is also interesting to evaluate the various denoising methods from a practical point of view: the computation time. Indeed, the results achieved by overcomplete representation are admittedly superior than the ones obtained by critically sampled wavelet transforms, but their weakness is the time they require (nearly 27 s on a Power Mac G5 workstation with 1.8 GHz PowerPC 970 CPU for  $256 \times 256$  images to obtain the redundant results reported in Table III). With our simple univariate method (14), the whole denoising process (including four iterations of an orthonormal wavelet transform) lasts approximately 0.4 s for  $256 \times 256$  images (1.6 s for  $512 \times 512$  images), using a similar workstation. With our inter-scale dependent thresholding function (21), the whole denoising task takes between 0.6 - 0.7 s for  $256 \times 256$  images and about 2.7 s for  $512 \times 512$  images. To compare with, Portilla's BLS-GSM with a  $3 \times 3$  window size lasts approximately 10 s for  $512 \times 512$  images, using the same orthonormal transform. Besides giving competitive results, our method is thus also much faster.

Table IV summarizes the relative computation time of the various methods considered in this paper. Note that the main part of the *ProbShrink* is contained in a pre-compiled file, making its execution time a bit faster than the other algorithms which are fully implemented in Matlab.

TABLE IV  
RELATIVE COMPUTATION TIME OF VARIOUS DENOISING TECHNIQUES.

Method	Unit of time [U]	
	$256 \times 256$ images	$512 \times 512$ images
<i>BayesShrink</i>	1.0	3.9
<i>BiShrink</i> $7 \times 7$	1.4	5.4
<i>ProbShrink</i> $3 \times 3$	2.8	6.6
<i>BLS-GSM</i> $3 \times 3$	7.8	30.0
Univariate sum of DOG (14)	1.2	4.5
Bivariate sum of DOG (21)	2.0	7.9
Redundant <i>BLS-GSM</i> $3 \times 3$	81.5	311.8

Note: The computation times have been averaged over twenty runs.

## VI. CONCLUSION

We have presented a new approach to orthonormal wavelet image denoising that does not need any prior statistical modelization of the wavelet coefficients. This approach is made possible thanks to the existence of an efficient estimate of the MSE between noisy and clean image—the SURE—that is based on the noisy data alone. Its minimization over a set of denoising processes automatically provides a near-optimal solution in the sense of the *a posteriori* MSE. For efficiency reasons, we have chosen this set to be a linear span of basic nonlinear mappings.

Using this approach, we have designed an image denoising algorithm that takes into account interscale dependencies, but discards intra-scale correlations. In order to compensate for features misalignment, we have developed a rigorous procedure based on the relative group delay between the scaling and wavelet filters—*group delay compensation*. The information brought by this new inter-scale predictor is used to classify smoothly between high-and low-SNR wavelet coefficients.

The comparison of the denoising results obtained with our algorithm, and with the best state-of-the-art nonredundant techniques (that integrate both inter- and intra-scale dependencies), demonstrate the efficiency of our SURE-based approach which gave the best output PSNRs for most of the images. The visual quality of our denoised images is moreover characterized by fewer artifacts than the other methods.

We are currently working on an efficient integration of the intra-scale correlations within the SURE-based approach. Our goal is to show that the consideration of inter- and intra-scale dependencies brings denoising gains that rival the quality of the best redundant techniques such as *BLS-GSM*.

## ACKNOWLEDGMENT

This work was supported in part by the Swiss National Science Foundation under grant 200020-109415.

## REFERENCES

- [1] D. L. Donoho and I. M. Johnstone, “Adapting to Unknown Smoothness via Wavelet Shrinkage,” *Journal of the American Statistical Association*, vol. 90, no. 432, pp. 1200–1224, December 1995.
- [2] —, “Ideal Spatial Adaptation via Wavelet Shrinkage,” *Biometrika*, vol. 81, pp. 425–455, 1994.
- [3] L. Breiman, “Better Subset Regression Using the Non-Negative Garrote,” *Technometrics*, vol. 37, no. 4, pp. 373–384, November 1995.
- [4] N. G. Kingsbury, “Image Processing with Complex Wavelets,” *Phil. Trans. R. Soc. A.*, September 1999.
- [5] H.-Y. Gao and A. G. Bruce, “Waveshrink with Firm Shrinkage,” *Statistica Sinica*, vol. 7, pp. 855–874, 1997.
- [6] —, “Wavelet Shrinkage Denoising Using the Non-Negative Garrote,” *J. Comput. Graph. Stat.*, vol. 7, no. 4, pp. 469–488, 1998.

- [7] C. Stein, "Estimation of the Mean of a Multivariate Normal Distribution," *The Annals of Statistics*, vol. 9, pp. 1135–1151, 1981.
- [8] S. G. Chang, B. Yu, and M. Vetterli, "Adaptive Wavelet Thresholding for Image Denoising and Compression," *IEEE Transactions on Image Processing*, vol. 9, no. 9, September 2000.
- [9] J. Portilla, V. Strela, M. J. Wainwright, and E. P. Simoncelli, "Image Denoising using Scale Mixtures of Gaussians in the Wavelet Domain," *IEEE Transactions on Image Processing*, vol. 12, no. 11, November 2003.
- [10] A. Pižurica and W. Philips, "Estimating the Probability of the Presence of a Signal of Interest in Multiresolution Single- and Multiband Image Denoising," *IEEE Transactions on Image Processing*, vol. 15, no. 3, March 2006.
- [11] L. Sendur and I. W. Selesnick, "Bivariate Shrinkage Functions for Wavelet-Based Denoising Exploiting Interscale Dependency," *IEEE Transactions on Signal Processing*, vol. 50, no. 11, November 2002.
- [12] —, "Bivariate Shrinkage With Local Variance Estimation," *IEEE Signal Processing Letters*, vol. 9, no. 12, December 2002.
- [13] N. G. Kingsbury, "Complex Wavelets for Shift Invariant Analysis and Filtering of Signals," *Journal of Applied Computational Harmonic Analysis*, vol. 10, no. 3, pp. 234–253, May 2001.
- [14] J.-L. Starck, E. J. Candès, and D. L. Donoho, "The Curvelet Transform for Image Denoising," *IEEE Transactions on Image Processing*, vol. 11, no. 6, June 2002.
- [15] M. S. Crouse, R. D. Nowak, and R. G. Baraniuk, "Wavelet-based Signal Processing Using Hidden Markov Models," *IEEE Transactions on Signal Processing*, vol. 46, no. 4, pp. 886–902, April 1998.
- [16] X.-P. Zhang and M. D. Desai, "Adaptive Denoising based on SURE Risk," *IEEE Signal Processing Letters*, vol. 5, no. 10, October 1998.
- [17] A. Benazza-Benyahia and J.-C. Pesquet, "Building Robust Wavelet Estimators for Multicomponent Images Using Stein's Principle," *IEEE Transactions on Image Processing*, vol. 14, no. 11, pp. 1814–1830, November 2005.
- [18] P. L. Combettes and J.-C. Pesquet, "Wavelet-constrained Image Restoration," *International Journal on Wavelets, Multiresolution and Information Processing*, vol. 2, no. 4, pp. 371–389, December 2004.
- [19] J.-C. Pesquet and D. Leporini, "A New Wavelet Estimator for Image Denoising," *Sixth International Conference on Image Processing and its Applications*, vol. 1, pp. 249–253, July 14–17 1997.
- [20] W. James and C. Stein, "Estimation with Quadratic Loss," *Proc. Fourth Berkeley Symp. Math. Statist. Probab.*, vol. 1, pp. 361–379, 1961.
- [21] S. G. Chang, B. Yu, and M. Vetterli, "Spatially Adaptive Wavelet Thresholding with Context Modeling for Image Denoising," *IEEE Transactions on Image Processing*, vol. 9, no. 9, September 2000.
- [22] I. Daubechies, "Ten Lectures on Wavelets," *CBMS-NSF Regional Conference series in Applied Mathematics*, vol. 61 of Proc., March 1992.
- [23] M. K. Mihçak, Kozintsev, K. Ramchandran, and P. Moulin, "Low-Complexity Image Denoising Based on Statistical Modeling of Wavelet Coefficients," *IEEE Signal Processing Letters*, vol. 6, no. 12, December 1999.
- [24] F. Abramovitch, T. Sapatinas, and B. W. Silverman, "Wavelet Thresholding via a Bayesian Approach," *Journal of the Royal Statistical Society. Series B*, vol. 60, no. 4, pp. 725–749, 1998.
- [25] J. S. Lee, "Digital Image Enhancement and Noise filtering by use of local Statistics," *IEEE Pattern Anal. Machine Intell.*, vol. PAMI-2, pp. 165–168, March 1980.
- [26] E. P. Simoncelli, *Bayesian Interference in Wavelet based Models*, ser. Lecture Notes in Statistics. New York: Springer-Verlag, March 1999, vol. 141, ch. 18, pp. 291–308.

- [27] B. Vidakovic, *Statistical Modeling by Wavelets*. Wiley-Interscience, April 1999.

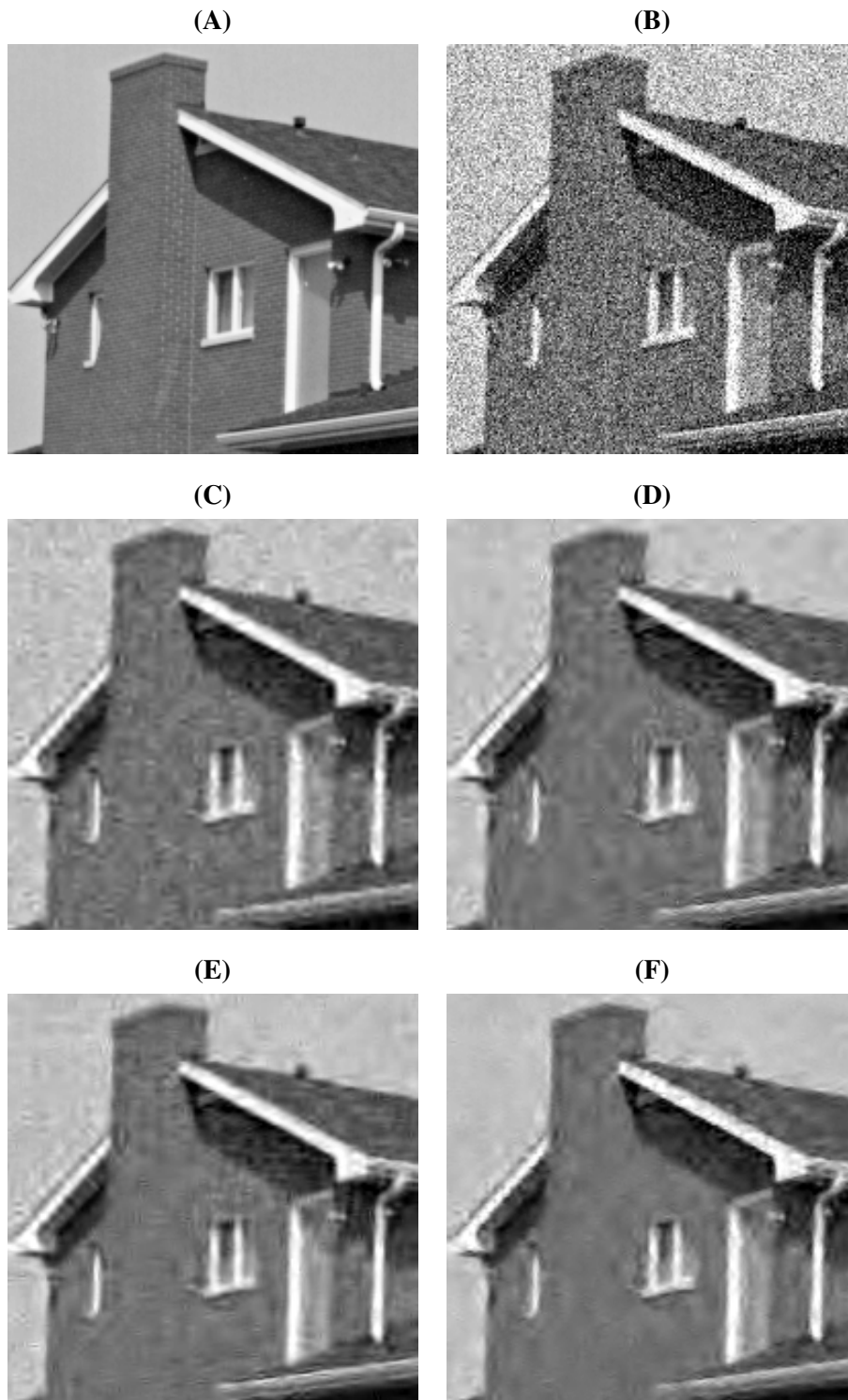


Fig. 12. (A) Part of the noise-free  $256 \times 256$  *House* image. (B) A noisy version of it: PSNR = 18.59 dB. (C) Denoised result using the *BayesShrink*: PSNR = 27.57 dB. (D) Denoised result using the *BiShrink*  $7 \times 7$ : PSNR = 28.19 dB. (E) Denoised result using the *BLS-GSM*  $3 \times 3$ : PSNR = 28.73 dB. (F) Denoised result using our inter-scale dependent thresholding function (21): PSNR = 28.96 dB.

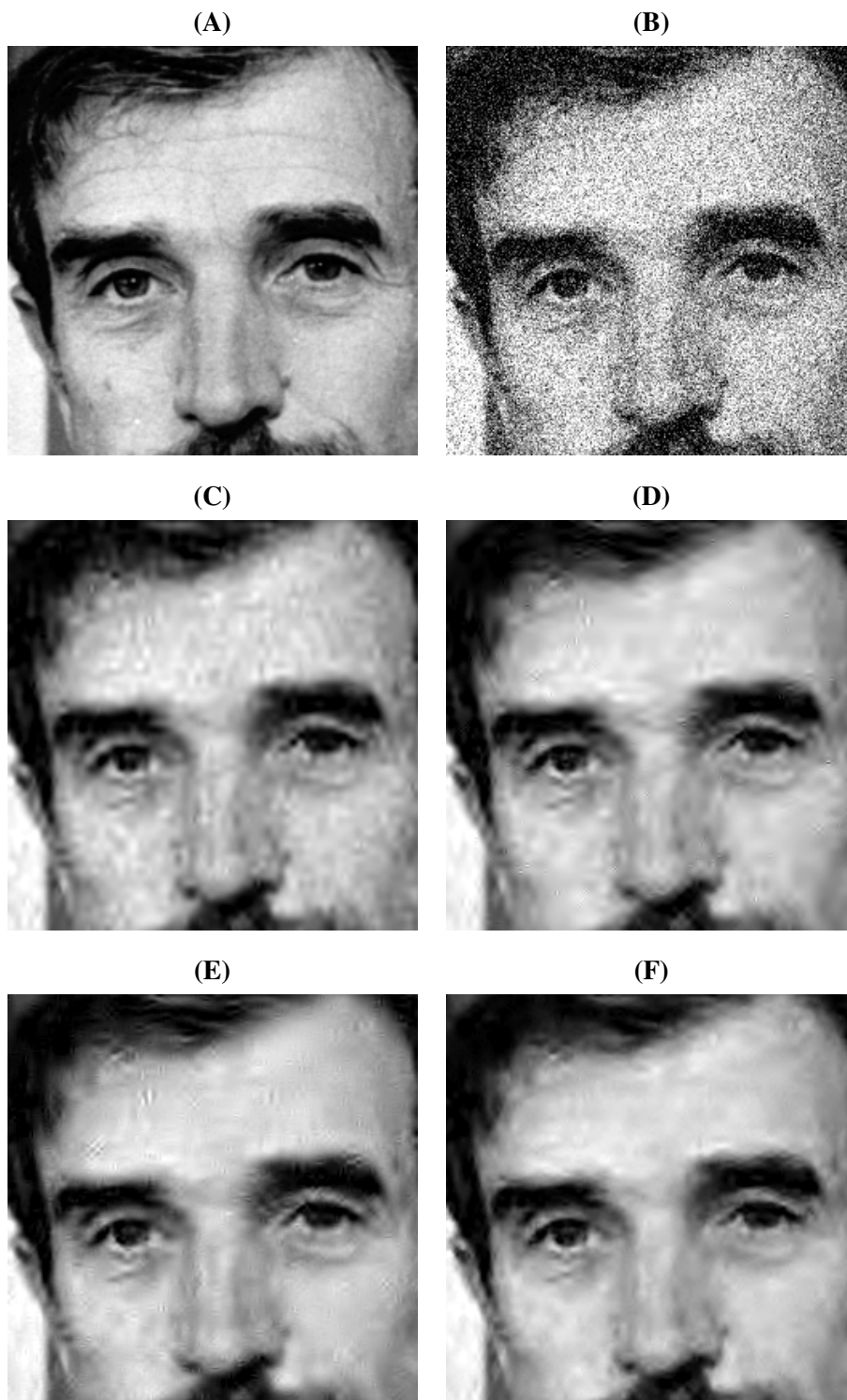


Fig. 13. (A) Part of the noise-free  $512 \times 512$  *AI* image. (B) A noisy version of it: PSNR = 14.15 dB. (C) Denoised result using the *BayesShrink*: PSNR = 26.71 dB. (D) Denoised result using the *BiShrink*  $7 \times 7$ : PSNR = 27.12 dB. (E) Denoised result using the *BLS-GSM*  $3 \times 3$ : PSNR = 27.34 dB. (F) Denoised result using our inter-scale dependent thresholding function (21): PSNR = 27.66 dB.