

Analysis of Low Bit Rate Image Transform Coding

Stéphane Mallat, *Member, IEEE*, and Frédéric Falzon

Abstract—Calculations based on high-resolution quantizations prove that the distortion rate $D(\bar{R})$ of an image transform coding is proportional to $2^{-2\bar{R}}$ when \bar{R} is large enough. In wavelet and block cosine bases, we show that if $\bar{R} < 1$ bit/pixel, then $D(\bar{R})$ varies like $\bar{R}^{1-2\gamma}$, where γ remains of the order of 1 for most natural images. The improved performance of embedded codings in wavelet bases is analyzed. At low bit rates, we show that the compression performance of an orthonormal basis depends mostly on its ability to approximate images with a few nonzero vectors.

Index Terms— Distortion-rate, image compression, JPEG, wavelet basis.

I. INTRODUCTION

IF THE SIGNALS to be encoded are realizations of a Gaussian process, under the high-resolution quantization hypothesis, we know nearly everything about the performance of a transform coding. For an average of \bar{R} bits per pixel, the mean-square error $D(\bar{R})$ varies proportionally to $2^{-2\bar{R}}$ with a constant that depends on the bit allocation and the basis. Current image transform coders operate below 1 bit per pixel. For such low bit rates, the high-resolution quantization assumption yields an incorrect estimate of the distortion rate $D(\bar{R})$. In this range, we show that $D(\bar{R})$ depends mostly on the error D_0 when approximating signals with a limited number of vectors selected from the orthogonal basis. We verify that in wavelet and block cosine bases, $D(\bar{R})$ varies like $\bar{R}^{1-2\gamma}$, where γ remains of the order of 1 for most “natural” images.

Transform coding algorithms can be improved by an embedding strategy that sends the larger amplitude coefficients first and then progressively refines their quantization. This improvement is analyzed mathematically and evaluated numerically for the wavelet zero-tree algorithm of Shapiro [13]. Embedded coders outperform classical transform coders when there is some prior information on which basis vectors produce large, average, or small decomposition coefficients for typical signals.

This paper begins with a brief review of entropy-constrained scalar quantizers and high bit rate transform coding. Section

III analyzes the distortion rate at low bit rates and gives numerical examples with a wavelet transform coder and JPEG. The performance of embedded transform coders and their applications to wavelets is studied in Section IV.

II. HIGH BIT-RATE COMPRESSION

The class of signals to be encoded is represented by a random vector Y of size N . Although these signals may be multidimensional-like images, they are indexed by an integer n to simplify notations: $Y[n]$. A transform coder decomposes these signals in an orthonormal basis $\mathcal{B} = \{g_m\}_{0 \leq m < N}$

$$Y = \sum_{m=0}^{N-1} A[m] g_m.$$

Each coefficient $A[m]$ is a random variable defined by

$$A[m] = \langle Y, g_m \rangle = \sum_{n=0}^{N-1} Y[n] g_m^*[n].$$

To construct a finite code, each coefficient $A[m]$ is approximated by a quantized variable $\hat{A}[m]$. We concentrate on scalar quantizations, which are most often used for transform coding. The next section reviews important results concerned with minimizing the quantization error.

A. Entropy-Constrained Scalar Quantization

A scalar quantizer Q approximates a real random variable X by a quantized variable $\hat{X} = Q(X)$, which takes its values in a finite set. Suppose that X takes its values in $[a, b]$, which may correspond to the whole real axis. We decompose $[a, b]$ into K intervals $(y_{k-1}, y_k]_{1 \leq k \leq K}$ of variable lengths, with $y_0 = a$ and $y_K = b$. If $x \in (y_{k-1}, y_k]$, then $Q(x) = x_k$. We denote

$$p_k = \Pr\{X \in (y_{k-1}, y_k]\} = \Pr\{\hat{X} = x_k\}.$$

The Shannon theorem [5] proves that the entropy

$$\mathcal{H}(\hat{X}) = - \sum_{k=1}^K p_k \log_2 p_k$$

is a lower bound of the average number of bits per symbol used to encode the values of \hat{X} . Arithmetic entropy coding [16] achieves an average bit rate that can be arbitrarily close to the entropy lower bound; therefore, we shall consider that this lower bound is reached. An *entropy constrained scalar quantizer* is designed to minimize $\mathcal{H}(\hat{X})$ for a fixed mean-square distortion $D = E\{(X - \hat{X})^2\}$.

Manuscript received February 15, 1997; revised November 30, 1997. This work was supported by the French Centre National d'Etudes Spatiales and AFOSR Grant F49620-96-1-0455. The associate editor coordinating the review of this paper and approving it for publication was P. P. Vaidyanathan.

S. Mallat is with the Department of Applied Mathematics, Ecole Polytechnique, Palaiseau, France, and with the Courant Institute, New York University, New York, NY 10012 USA.

F. Falzon was with INRIA, Sophia Antipolis, France. He is now with Alcatel Alsthom Recherche, Marcoussis, France.

Publisher Item Identifier S 1053-587X(98)02558-6.

Let $p(x)$ be the probability density of the random source X . The differential entropy of X is defined by

$$\mathcal{H}_d(X) = - \int_{-\infty}^{+\infty} p(x) \log_2 p(x) dx.$$

A quantizer is said to have a *high resolution* if $p(x)$ is approximately constant on each quantization bin $(y_{k-1}, y_k]$ of size $\Delta_k = y_k - y_{k-1}$. This is the case if the sizes Δ_k are sufficiently small relative to the rate of variation of $p(x)$. The following theorem [5] proves that uniform quantizers are optimal among high-resolution quantizers. It is equivalent to the result of Girsh and Pierce [6], which proves the asymptotic optimality of uniform quantizers when the sizes of the quantization bins tend to zero.

Theorem 1: If Q is a high-resolution quantizer with respect to $p(x)$, then

$$\mathcal{H}(\hat{X}) \geq \mathcal{H}_d(X) - \frac{1}{2} \log_2 (12D). \quad (1)$$

This inequality is an equality if and only if Q is a uniform quantizer, in which case, $D = (\Delta^2/12)$.

For a fixed distortion D , under the high-resolution quantization hypothesis, the minimum average bit rate $R_X = \mathcal{H}(\hat{X})$ is therefore achieved by a uniform quantizer, and

$$R_X = \mathcal{H}_d(X) - \log_2 \Delta. \quad (2)$$

The distortion rate is obtained by taking the inverse of (1)

$$D(R_X) = \frac{1}{12} 2^{2\mathcal{H}_d(X)} 2^{-2R_X}. \quad (3)$$

Farvardin and Modestino [4] proved that even though the high-resolution quantization hypothesis may not hold, for a large class of probability distribution including generalized Gaussians, the uniform quantizer yields a distortion rate that is close to the optimal quantizer if the number of quantization bins K is large enough.

B. Distortion Rate

Let us optimize the transform coding of $Y = \sum_{m=0}^{N-1} A[m] g_m$. The average bit budget to encode $\hat{A}[m] = Q(A[m])$ is $R_m = \mathcal{H}(\hat{A}[m])$. For a high-resolution quantization, Theorem 1 proves that the error $D_m = E\{|A[m] - \hat{A}[m]|^2\}$ is minimized when using a uniform scalar quantization. An optimal bit allocation minimizes the total number of bits $R = \sum_{m=0}^{N-1} R_m$ for a specified total error $D = \sum_{m=0}^{N-1} D_m$. Let $\bar{R} = (R/N)$ be the average number of bits per sample. With Lagrange multipliers, we verify that \bar{R} is minimum if all D_m are equal [7], in which case

$$D(\bar{R}) = \frac{N}{12} 2^{2\bar{\mathcal{H}}_d} 2^{-2\bar{R}} \quad (4)$$

where $\bar{\mathcal{H}}_d$ is the averaged differential entropy

$$\bar{\mathcal{H}}_d = \frac{1}{N} \sum_{m=0}^{N-1} \mathcal{H}_d(A[m]).$$

A mean-square error D is not always a good measurement of visual degradations in images. In particular, we are often less sensitive to high-frequency distortions as opposed to lower

frequencies. A weighted mean-square norm takes into account this sensitivity by emphasizing low-frequency errors. Suppose that g_m is a vector whose Fourier transform is localized in a frequency neighborhood that depends on m , as in a block cosine basis or in a wavelet basis. A weight w_m is adjusted, depending on our visual sensitivity in this frequency range. The resulting weighted distortion is

$$D^w = \sum_{m=0}^{N-1} w_m^2 D_m = \sum_{m=0}^{N-1} D_m^w$$

where $D_m^w = w_m^2 D_m$ is the mean-square quantization error of $A^w[m] = w_m A[m]$. The previous bit allocation result applied to the coefficients $A^w[m]$ proves that D^w is minimized by quantizing uniformly $A^w[m]$ with a bin size Δ , which is equivalent to quantizing $A[m]$ with a bin size $\Delta_m = (\Delta/w_m)$. In the rest of the paper, we choose $w_m = 1$ and, thus, evaluate the error with a standard mean-square norm as opposed to a weighted norm. All calculations are easily extended to any other choice of weights by replacing $A[m]$ by $w_m A[m]$.

The distortion rate (4) depends on the orthonormal basis \mathcal{B} through $\bar{\mathcal{H}}_d$. In general, it is difficult to find \mathcal{B} , which minimizes $\bar{\mathcal{H}}_d$ because the probability density of $A[m] = \langle Y, g_m \rangle$ may depend on g_m in a complicated way. If Y is a Gaussian random vector, then the coefficients $A[m]$ are Gaussian random variables in any basis. In this case, the probability density of $A[m]$ depends only on the variance σ_m^2 , and we can verify that

$$\mathcal{H}_d(A[m]) = \log_2 \sigma_m + \log_2 \sqrt{2\pi e}.$$

Inserting this expression in (4) yields

$$D(\bar{R}) = N \frac{\pi e}{6} \left(\prod_{m=0}^{N-1} \sigma_m^2 \right)^{1/N} 2^{-2\bar{R}}. \quad (5)$$

One can prove that $\prod_{m=0}^{N-1} \sigma_m^2$ is minimum if and only if \mathcal{B} is a Karhunen-Loève basis of Y [5], which means that \mathcal{B} diagonalizes the covariance matrix of Y . The transform coding of a Gaussian process is thus optimized in a Karhunen-Loève basis. When Y is not Gaussian, the Karhunen-Loève basis is *a priori* no longer optimal. This is the case for images that cannot be considered to be realizations of a Gaussian process.

Let us describe a simple wavelet transform coder for images. Separable wavelet bases of images include three wavelets with horizontal, vertical, or diagonal orientations [10], indexed by $1 \leq k \leq 3$. At an orientation k and scale 2^j , the wavelet vector $g_m = \psi_{j,p,q}^k$ is approximately centered at $(2^j p, 2^j q)$, with a square support whose size is proportional to 2^j . At high bit rates, we saw that the distortion rate is optimized by quantizing uniformly all decomposition coefficients. The domains where the image has smooth grey-level variations yield small amplitude wavelet coefficients that are quantized to zero. To improve the efficiency of this transform coding, the wavelet coefficients are scanned in a predefined order, and the position of zero versus nonzero quantized coefficients is recorded with a run-length coding that is entropy encoded. In

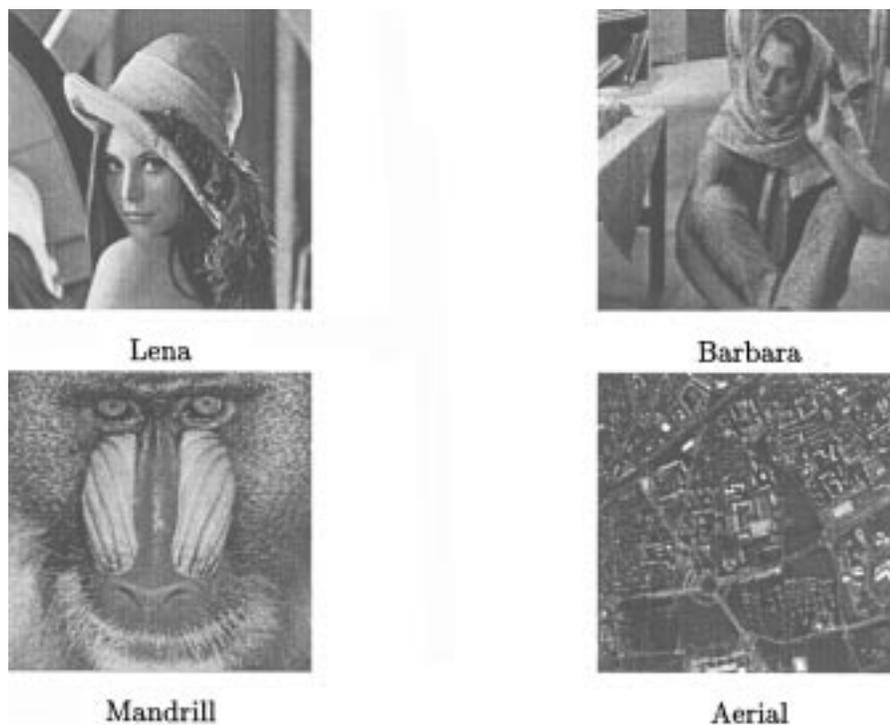


Fig. 1. Four test images for numerical experiments.

the same scanning order, the amplitude of the nonzero quantized coefficients are also entropy encoded with a Huffman or an Arithmetic coding.

Fig. 2 gives the distortion $\log_2 D(\bar{R})$ of this wavelet transform coding for the test images shown in Fig. 1. These numerical experiments are performed with an orthogonal cubic spline Battle-Lemarié wavelet [10]. The distortion rate formula (4) predicts that

$$\log_2 D(\bar{R}) = 2\bar{\mathcal{H}}_d + \log_2 \left(\frac{N}{12} \right) - 2\bar{R}$$

where $\bar{\mathcal{H}}_d$ is the average differential entropy of the wavelet coefficients at all scales and positions. This formula implies that $\log_2 D(\bar{R})$ should decay with a slope of -2 as a function of \bar{R} . This is indeed verified in Fig. 2 for $\bar{R} \geq 1$ but not for $\bar{R} < 1$, where $\log_2 D(\bar{R})$ has a much faster decay. At low bit rates $\bar{R} < 1$, the distortion rate formula (4) is not valid because the high-resolution quantization assumption does not hold. Wavelet transform codings are most often used for $\bar{R} < 1$ because they recover images of nearly perfect visual quality up to $\bar{R} = 0.5$ bits per pixel. The next section studies the distortion rate at these low bit rates.

III. LOW BIT-RATE COMPRESSION

At low bit rates, the decomposition coefficients of an image in an orthonormal basis are coarsely quantized. Since many coefficients are set to zero, the positions of zero versus nonzero quantized coefficients are stored in a binary significance map, which is recorded with a run-length coding or a more sophisticated zero-tree algorithm. The distortion rate theory previously described does not apply for two reasons. First, the high-resolution quantization hypothesis does not hold because the

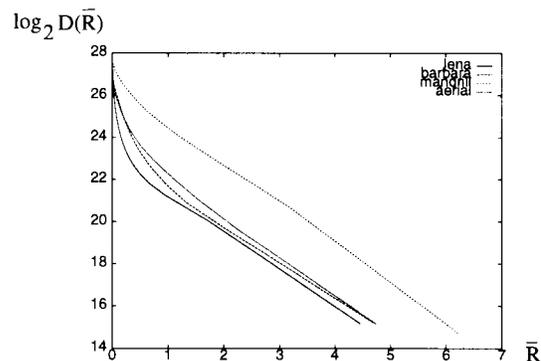


Fig. 2. Log distortion rate curve for the wavelet transform coding of each test image.

quantization bins are large. Second, one cannot treat the total bit budget R as a sum of bits R_m allocated independently to each decomposition coefficient. Indeed, the encoding of the zero quantized coefficients through a significance map is a form of vector quantization, which relates the encoding of different coefficients.

To evaluate the distortion rate, we cannot rely on a precise stochastic model for images. There is, as yet, no model that incorporates the full diversity of image structures, such as nonstationary textures and edges. To avoid this difficulty, we shall consider the signals to be deterministic vectors whose decomposition coefficients in the basis \mathcal{B} have a parameterized decay. The distortion rate is, therefore, not calculated with an ensemble average but for each signal f . The key result shows that this distortion rate depends mostly on the ability to precisely approximate f with a small number of vectors selected from \mathcal{B} . Low bit-rate image compressions in wavelet

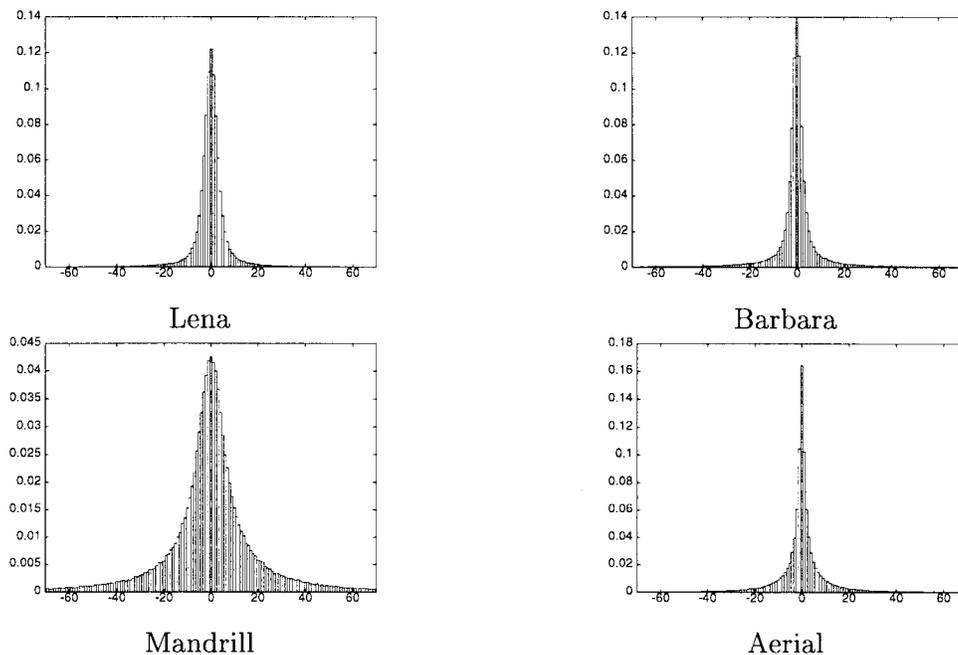


Fig. 3. Normalized histograms of the cubic-spline wavelet coefficients of the test images.

bases and block cosine bases illustrate the distortion rate results.

A. Distortion Rate

Let f be a signal decomposed in an orthonormal basis $\mathcal{B} = \{g_m\}_{0 \leq m < N}$

$$f = \sum_{m=0}^{N-1} a[m] g_m \quad \text{with } a[m] = \langle f, g_m \rangle.$$

The transform coder quantizes all coefficients and reconstructs

$$\hat{f} = \sum_{m=0}^{N-1} Q(a[m]) g_m.$$

The coding error is

$$D = \|f - \hat{f}\|^2 = \sum_{m=0}^{N-1} |a[m] - Q(a[m])|^2. \quad (6)$$

We denote by $h[x]$ the discrete histogram of the N coefficients $a[m]$ normalized so that $\sum_x h[x] = 1$. The values of this histogram are interpolated to define a function $p(x) \geq 0$ for all $x \in \mathbf{R}$ such that $\int_{-\infty}^{+\infty} p(x) dx = 1$. This $p(x)$ is the probability density of a random variable X . We suppose N to be sufficiently large and the histogram sufficiently regular so that for all functions $\phi(x)$ that appear in our calculations, we have

$$\begin{aligned} \frac{1}{N} \sum_{m=0}^{N-1} \phi(a[m]) &= \sum_x \phi(x) h[x] \\ &\approx \int_{-\infty}^{+\infty} \phi(x) p(x) dx = E\{\phi(X)\}. \end{aligned} \quad (7)$$

This hypothesis holds for the histograms of the test images shown in Fig. 3, as well as for most “natural” images. It is

equivalent to the coefficients $a[m]$ being successive values of the random variable X . Applied to $\phi(x) = |x - Q(x)|^2$, (7) yields

$$\frac{D}{N} = \frac{1}{N} \sum_{m=0}^{N-1} |a[m] - Q(a[m])|^2 = E\{|X - Q(X)|^2\}.$$

Let \bar{R} be the average number of bits per coefficient to encode the $Q(a[m])$. If Q is a high-resolution uniform quantizer with step size Δ , then (3) implies a distortion rate formula similar to (4)

$$\frac{D(\bar{R})}{N} = \frac{1}{12} 2^{2\mathcal{H}_d(X)} 2^{-2\bar{R}}.$$

If the basis \mathcal{B} is chosen so that many coefficients $a[m] = \langle f, g_m \rangle$ are close to zero and few have a large amplitude, then $p(x)$ has a sharp high peak at $x = 0$. This is the case for the histograms of wavelet coefficients shown in Fig. 3. If Δ is large, then $p(x)$ has important variations in the zero bin $[-(\Delta/2), (\Delta/2)]$, where coefficients are quantized to zero. Hence, the high-resolution quantization hypothesis does not apply in this zero bin. This explains why $\log_2 D(\bar{R})$, which is shown in Fig. 2, decays like $-2\bar{R}$ only for $\bar{R} \geq 1$.

If $p(x)$ is a Laplacian distribution $p(x) = (1/\sigma\sqrt{2})e^{-\sqrt{2}|x|/\sigma}$, then Sullivan [14] proved that the optimal entropy-constrained scalar quantizer is a nearly uniform quantizer. All the nonzero quantization bins have the same size Δ , but the zero bin $[-T, T]$ is larger, with a ratio $\theta = (T/\Delta)$ that can be calculated. Sullivan’s result does not apply to the quantization of wavelet image coefficients because Section III-B shows that $p(x)$ has a slower decay, which is rational instead of exponential. Yet efficient low bit-rate wavelet transform coders are often implemented with a nearly uniform quantizer, whose zero bin $[-T, T]$ is larger than the other bins. This can be justified with

a high-resolution assumption outside the zero bin. For $|x| > T$, we shall consider that the relative variations of $p(x)$ in each quantization bin is sufficiently small to apply the high-resolution hypothesis. This assumption is only an approximation, but it captures sufficiently precisely the quantizer properties to obtain accurate calculations up to very low bit rates. For $x \in [-T, T]$, the high-resolution quantization hypothesis does not hold because $p(x)$ varies too much. Theorem 1 implies that outside $[-T, T]$, the optimal entropy-constrained quantizer Q has bins of constant size Δ . The ratio $\theta = (T/\Delta)$ is a parameter that must be adjusted to minimize the overall distortion D .

Any coefficient $|a[m]| > T$ that is not quantized to zero is called a *significant coefficient*. Coding the position of nonzero quantized coefficients is equivalent to storing a binary *significance map*, which is defined by

$$b[m] = \begin{cases} 0 & \text{if } |a[m]| \leq T \\ 1 & \text{if } |a[m]| > T \end{cases}. \quad (8)$$

JPEG and the wavelet image coder of Section II-B use a run-length coding to store this significance map. More efficient zero-tree encoding techniques may also be used for wavelet significance maps [9].

Let R_0 be the total number of bits required to encode the significance map. Let M be the number of significant coefficients. There is a proportion $p = (M/N)$ of indices m such that $b[m] = 1$. An upper bound for R_0 is computed by supposing that there is no redundancy in the position of the 0 and the 1 in the significance map. The average number of bits to encode the position of one coefficient is then the entropy of a binary source with a probability $p = (M/N)$ to be equal to 1 and $1 - p$ to be equal to 0

$$\frac{R_0}{N} \leq -p \log_2 p - (1 - p) \log_2 (1 - p).$$

For $x \in (0, 1]$, then $-x \log_2 x \leq (1 - x) \log_2 e$; therefore, the average number of bits per significant coefficient to encode the significance map is

$$r_0 = \frac{R_0}{M} \leq \log_2 \frac{N}{M} + \log_2 e. \quad (9)$$

For wavelet coefficients, when the proportion of significant coefficients M/N is small, a run-length coding yields an average bit rate $r_0 = (R_0/M)$, which is much smaller than the upper bound (9), because of the redundancy in the positions of the zero coefficients. For large classes of images, numerical calculations show that $r_0 = (R_0/M)$ varies slowly relative to M/N .

The amplitude of the M significant coefficients is uniformly quantized with a step Δ , and these quantized values are entropy encoded. Let us compute the total number of bits R_1 of the resulting entropy coding. For $M \gg 1$, the M significant coefficients $a[m]$ above T have a normalized histogram that is interpolated by

$$p_T(x) = \frac{N}{M} p(x) \mathbf{1}_{\{|x| > T\}}.$$

Let X_T be the random variable whose probability density is $p_T(x)$. Since the high-resolution quantization hypothesis

applies to significant coefficients, the average number of bits to encode the amplitude of each quantized significant coefficient, which is denoted r_1 , is calculated from (2)

$$r_1 = \frac{R_1}{M} = \mathcal{H}_d(X_T) - \log_2 \Delta. \quad (10)$$

Overall, the transform coding requires $R = R_0 + R_1$ bits.

To estimate the quantization error [$D = \|f - \hat{f}\|^2$ in (6)], insignificant coefficients quantized to zero are separated from significant coefficients $D = D_0 + D_1$, where

$$D_0 = \sum_{|a[m]| \leq T} |a[m]|^2 \quad (11)$$

is the error due to quantizing insignificant coefficients, and

$$D_1 = \sum_{|a[m]| > T} |a[m] - Q(a[m])|^2 \quad (12)$$

is the error due to quantizing significant coefficients. The average quantization error D_1/M per significant coefficient is calculated with the high-resolution quantization assumption

$$\begin{aligned} \frac{D_1}{M} &= \frac{1}{M} \sum_{|a[m]| > T} |a[m] - Q(a[m])|^2 \\ &= E\{|X_T - Q(X_T)|^2\} = \frac{\Delta^2}{12}. \end{aligned} \quad (13)$$

To compute the error due to quantizing insignificant coefficients, we denote by f_M the approximation of f using the M vectors g_m of \mathcal{B} such that $|a[m]| = |\langle f, g_m \rangle| > T$

$$f_M = \sum_{|a[m]| > T} a[m] g_m.$$

The signal f_M can also be interpreted as an approximation of f from the M vectors of \mathcal{B} , whose inner products with f have the largest amplitude. The distortion D_0 can be rewritten

$$D_0 \left(\frac{M}{N} \right) = \|f - f_M\|^2 = \sum_{|a[m]| \leq T} |a[m]|^2. \quad (14)$$

In approximation theory, D_0 is called a *nonlinear approximation error* because the M vectors are selected depending on f as opposed to linear algorithms that approximate all signals with the same M vectors. Clearly, D_0 decays when M increases, but how fast? This issue is a central question that is studied in approximation theory in relation to particular functional spaces [3], [11].

Let us sort the inner products $a[m]$ by their amplitude. The amplitude of the k th coefficient is written

$$\begin{aligned} x \left(\frac{k}{N} \right) &= |a[m_k]| \leq x \left(\frac{k+1}{N} \right) \\ &= |a[m_{k+1}]| \quad \text{for } 1 \leq k < N. \end{aligned} \quad (15)$$

The approximation error D_0 is the sum of the $N - M$ squared coefficients of smaller amplitude

$$D_0 \left(\frac{M}{N} \right) = \sum_{k=M+1}^N \left| x \left(\frac{k}{N} \right) \right|^2. \quad (16)$$

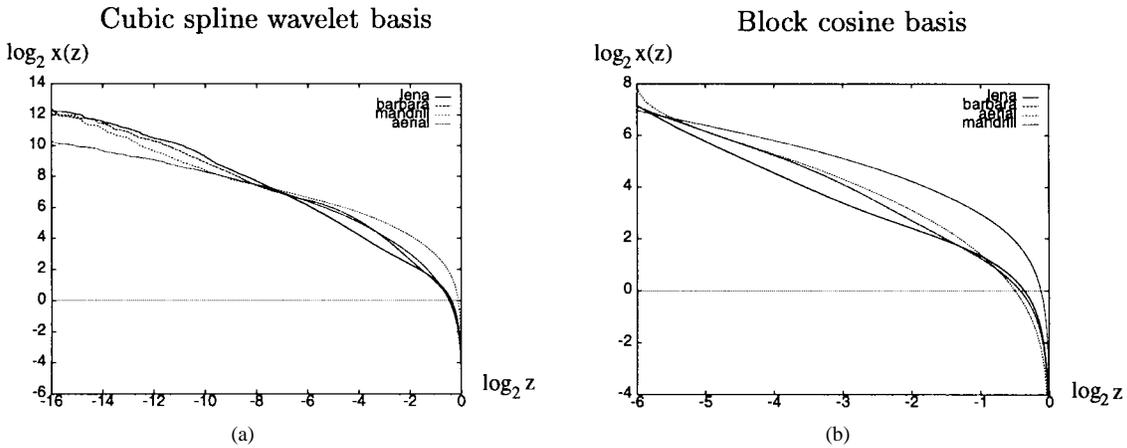


Fig. 4. Amplitude of the sorted decomposition coefficients of the test images in two bases.

The error $D_0(M/N)$ has a fast decay when M/N increases if $x(z)$ decreases quickly when z increases.

Observe that $p(x)$ is related to the inverse $z(x)$ of $x(z)$ by

$$z(x) = 1 - \int_{-x}^x p(u) du. \tag{17}$$

The probability density $p(x)$ is defined as a function of a continuous variable by interpolating the normalized histogram of the decomposition coefficients $a[m]$ of f . This also defines a function $x(z)$ for any $z \in [0, 1]$, which interpolates the values $x(k/N)$.

To estimate the decay of $D_0(z)$ when z increases, a standard approximation theory approach computes the rational decay of the sorted coefficients and, hence, compares $x(z)$ with $z^{-\gamma}$ for some $\gamma > 0$. For functions $f \in [0, N]^2$ decomposed in a wavelet orthonormal basis, the exponent γ characterizes particular functional spaces called Besov spaces [3]. To suppose that $x(z) = Cz^{-\gamma}$ would clearly be too restrictive to model interesting classes of signals. We shall rather suppose that this exponent is slowly varying and define

$$\gamma(z) = -\frac{d \log_2 x(z)}{d \log_2 z}. \tag{18}$$

Fig. 4 plots $\log_2 x(z)$ as a function of $\log_2 z$ for the wavelet coefficients and the block cosine coefficients of the test images. Observe that in both cases, the slope $\gamma(z)$ varies slowly for $z \leq 2^{-1}$. This behavior is further discussed in Section III-B.

To compute the distortion rate, we shall assume that the second-order derivative is bounded by a small $\epsilon > 0$

$$\left| \frac{d^2 \log_2 x(z)}{(d \log_2 z)^2} \right| \leq \epsilon \quad \text{for } z \in \left(0, \frac{1}{2}\right]. \tag{19}$$

We also suppose that

$$\inf_{z \in [0,1]} \gamma(z) > 0 \tag{20}$$

$$\frac{d^2 \log_2 x(z)}{(d \log_2 z)^2} \leq 0 \quad \text{for } z \in (0, 1) \tag{21}$$

and that $p(x)$ is symmetric

$$p(x) = p(-x). \tag{22}$$

The four test images, as well as most natural images, have wavelet and block cosine coefficients that satisfy (19)–(22). The concavity (21) is a technical condition that is used to control corrective terms in distortion rate calculations but is generally satisfied. The symmetry (22) of the probability density is verified in Fig. 3. Assuming (19)–(22), the following theorem relates the distortion rate $D(\bar{R})$ to the approximation error D_0 through parameters that are evaluated as a function of $\gamma(z)$ and the number M of significant coefficients.

Theorem 2: Suppose that $x(z)$ satisfies (19)–(22). Let $\gamma_M = \gamma(M/N) > \frac{1}{2}$. If $(M/N) \leq \epsilon$ and $M \geq (1/\epsilon)$, then

$$D(\bar{R}) = (1 + K)D_0 \left(\frac{\bar{R}}{r_1 + r_0} \right) \tag{23}$$

with

$$K = \frac{D_1}{D_0} = \frac{2\gamma_M - 1}{12\theta^2} [1 + O(\epsilon |\log_2 \epsilon|^2 + \epsilon^{2\gamma_M - 1})] \tag{24}$$

and

$$r_1 = \frac{R_1}{M} = 1 + (1 + \gamma_M) \log_2 e + \log_2 \gamma_M + \log_2 \theta + O(\epsilon). \tag{25}$$

Moreover, the derivative of $D_0(z)$ satisfies

$$\frac{d \log_2 D_0 \left(\frac{M}{N} \right)}{d \log_2 z} = (1 - 2\gamma_M) [1 + O(\epsilon |\log_2 \epsilon|^2 + \epsilon^{2\gamma_M - 1})], \tag{26}$$

The proof is in Appendix A.1. To understand the implications of this theorem, these formulae are simplified with an approximation, and we neglect the corrective terms in ϵ . Since the second-order derivative (19) remains small, the slope

$$\gamma_M = \frac{-d \log_2 x \left(\frac{M}{N} \right)}{d \log_2 z} \tag{27}$$

varies slowly as a function of $\log_2(M/N)$. In the compression range of interest, it will be considered constant $\gamma_M \approx \gamma$. It follows from (24) and (25) that $K = (D_1/D_0)$ and $r_1 = (R_1/M)$ are also constant. We have already mentioned

that this is also the case for $r_0 = (R_0/M)$. Hence, $D(\bar{R})$ is calculated in (23) by scaling and multiplying the nonlinear approximation error $D_0(z)$ by constant factors

$$D(\bar{R}) = (1 + K)D_0\left(\frac{\bar{R}}{r_1 + r_0}\right). \quad (28)$$

Since $\gamma_M \approx \gamma$, (26) implies that $D_0(z) \sim z^{2\gamma-1}$ so that

$$D(\bar{R}) \sim \bar{R}^{1-2\gamma}.$$

This distortion rate decay is very different from the high-resolution formula where $D(\bar{R}) \sim 2^{-2\bar{R}}$.

The distortion D in (28) depends essentially on the approximation error D_0 of f from $M = \bar{R}/(r_1 + r_0)$ vectors selected in the basis \mathcal{B} . To optimize the transform coding, the basis \mathcal{B} must be able to approximate precisely each signal f with a small number of basis vectors. If we consider f to be a realization of a random vector Y , then we may wonder which is the basis that minimizes $E\{D_0(M/N)\}$ over all realizations. This problem is difficult because the M basis vectors are selected depending on each realization f of Y . They are the ones that have the largest inner products with f . The energy compaction theorem [5] proves that the Karhunen-Loève basis is optimal for approximating Y from M vectors chosen once and for all, but it has no optimality property in this nonlinear setting, where the vectors depend on each realization. In some cases, we know how to find bases that minimize the maximum error $D_0(M/N)$ over a whole signal class. For example, wavelet bases are optimal in this min-max sense for piecewise regular signals that belong to Besov spaces [3].

To optimize the quantization, the size of the zero bin $[-T, T]$ must be adjusted with respect to the other quantization bins of size Δ . To minimize the distortion D , we want to find $\theta = (T/\Delta)$ such that for a fixed \bar{R}

$$\frac{\partial D(\bar{R}, \theta)}{\partial \theta} = 0. \quad (29)$$

Appendix A.2 proves the following theorem that gives an analytic formula for θ .

Theorem 3: Suppose that $x(z)$ satisfies (19)–(22) and that $r_0 = (R_0/M)$ is a constant independent of M . Let $\gamma_M = \gamma(M/N) > \frac{1}{2}$. If $(M/N) \leq \epsilon$ and $M \geq (1/\epsilon)$, then the optimal zero bin ratio is

$$\theta = \sqrt{\frac{r_1 + r_0}{6 \log_2 e} - \frac{2\gamma_M - 1}{12}} [1 + \mathcal{O}(\epsilon)]. \quad (30)$$

B. Wavelet Transform Coding

Wavelet bases are known to efficiently approximate piecewise regular functions with a small number of nonzero wavelet coefficients [11]. Since images often include piecewise regular structures, wavelet bases are good candidates for building efficient image transform coders. The central assumption of Theorem 2 is that the sorted decomposition coefficients $x(z)$ of f in the basis \mathcal{B} have a rational decay. In wavelet bases, this is in accordance with asymptotic image models based on Besov spaces introduced for compression by DeVore *et al.* [3]

and further studied in [1]. Let $f \in \mathbf{L}^2[0, N]^2$. If there exists $C > 0$ and $\gamma > \frac{1}{2}$ so that for all $k \geq 0$ the wavelet coefficient of f of rank k is bounded by $Ck^{-\gamma}$, then f belongs to a family of Besov spaces whose indexes depend on γ . Let us consider a piecewise regular image f , which is uniformly regular (Lipschitz $\alpha \geq 1$) inside the regions $\{\Omega_i\}_{1 \leq i \leq K}$ that partition $[0, N]^2$. This image has discontinuities along the borders of the Ω_i , which have a finite length. One can then prove [11] that the sorted wavelet coefficients decay like $Ck^{-\gamma}$ with $\gamma = 1$. The discontinuities create large amplitude wavelet coefficient that are responsible for this decay exponent. This piecewise regular model applies to an image such as Lena because even the fur texture does not create enough large wavelet coefficients to modify the exponent $\gamma = 1$. On the other hand, the Mandrill image is composed of regions with highly irregular textures that create enough high-amplitude wavelet coefficients to reduce the exponent γ . For finite images, $k \leq N$, which is why we renormalize $z = (k/N)$ and compare the decay of $x(z)$ with $z^{-\gamma}$ when z increases in $[0, 1]$. Theorem 2 does not require γ to be a constant, but (19) assumes that it varies slowly as a function of $\log_2 z$.

If $x(z) \sim z^{-\gamma}$ for x large enough, then one can derive from (51) that $p(x) \sim x^{-1-(1/\gamma)}$. However, Fig. 3 shows that $p(x)$ has an exponential decay when x is small. This can be explained by looking at the normalized histograms $p_j(x)$ of the wavelet coefficients $\langle f, \psi_{j,p,q}^k \rangle$ at a fixed scale 2^j for all positions $0 \leq p, q \leq 2^{-j}N$ and orientations $1 \leq k \leq 3$. Such histograms are well modeled by generalized Gaussian distributions [10], which have an exponential decay and a variance that increases with the scale 2^j . It is the aggregation of these histograms that yield a global histogram $p(x)$, which has a rational decay for x sufficiently large [8]. However, the finite image resolution implies that wavelet coefficients are zero for $j < 0$. One can verify that the “border effect” created by the absence of finer scale wavelet coefficients implies that $p(x) \approx p_1(x)$ for x small. This explains the exponential decay of $p(x)$ when x is small, and the rapid variation of the slope γ for $\log_2 z \geq -1$ in Fig. 4(a).

Let us now evaluate numerically the precision of the analytic formula given by Theorem 2. The wavelet transform coder is implemented with a quantizer whose zero bin $[-T, T]$ is twice as large as the other quantization bins, which means that $\theta = (T/\Delta) = 1$. This is a standard choice in most wavelet image transform coders. The significance map is stored with a run-length coding, as explained in Section II-B. The ratio D_1/D_0 was calculated numerically for the wavelet transform coding of Lena. Fig. 5 compares this value with the theoretical estimate (24), where the corrective terms in ϵ are neglected $K = (2\gamma_M - 1)/12\theta^2$. The slope γ_M in (27) is computed numerically from the sorted coefficients $x(z)$ of Lena shown in Fig. 4(a). For $(M/N) \leq 2^{-4}$, Fig. 5 shows that K closely approximates the true value of D_1/D_0 . The error increases for $(M/N) > 2^{-4}$ because the Theorem hypothesis $(M/N) \leq \epsilon$ is not respected.

Fig. 6 displays the value of R_1/M computed numerically from the entropy of the quantized wavelet coefficients of Lena.

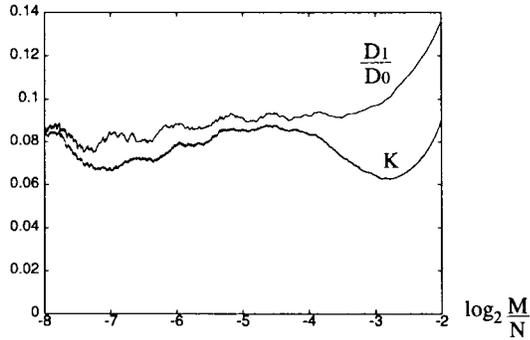


Fig. 5. Comparison of D_1/D_0 with the theoretical estimate $K = (2\gamma_M - 1)/12\theta^2$ for the wavelet coding of Lena.

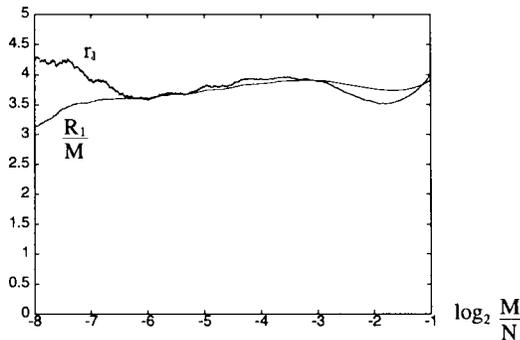


Fig. 6. Comparison of R_1/M with the theoretical estimate $r_1 = 1 + (1 + \gamma_M)\log_2 e + \log_2 \gamma_M + \log_2 \theta$ for the wavelet coding of Lena.

The theoretical estimate (25)

$$r_1 = 1 + (1 + \gamma_M)\log_2 e + \log_2 \gamma_M + \log_2 \theta. \quad (31)$$

is plotted in the same figure. The curves remain close, which verifies the precision of this calculation.

To simplify the expression of the distortion rate, the slope γ_M is approximated by a constant $\gamma_M \approx 1$, which corresponds to piecewise regular image models. Although γ_M can differ from 1 in many images such as Mandrill, this approximation is justified by the small sensitivity of K and $r_1 + r_0$ with respect to fluctuations of γ_M around 1. Since $\theta = 1$, we get $K \approx (1/12)$ and $r_1 \approx 3.9$. Fig. 7 displays

$$\frac{R}{M} = \frac{R_0}{M} + \frac{R_1}{M} = r_0 + r_1$$

, which was computed numerically for the four test images. For $(M/N) \in [2^{-7}, 2^{-1}]_+$, the ratio R/M can be approximated by a constant $r_0 + r_1 \approx 6.5$. The distortion D calculated in (23) is thus approximated by

$$\hat{D}(\bar{R}) = \left(1 + \frac{1}{12}\right) D_0\left(\frac{\bar{R}}{6.5}\right). \quad (32)$$

Fig. 8 compares the peak signal-to-noise ratio

$$P_{\text{SNR}}(D) = 10\log_{10} \frac{N255^2}{D}$$

which was calculated numerically for the four test images, with the approximated $P_{\text{SNR}}(\hat{D})$ derived from (32). Observe that $\hat{D}(\bar{R})$ gives a remarkably precise evaluation of the true

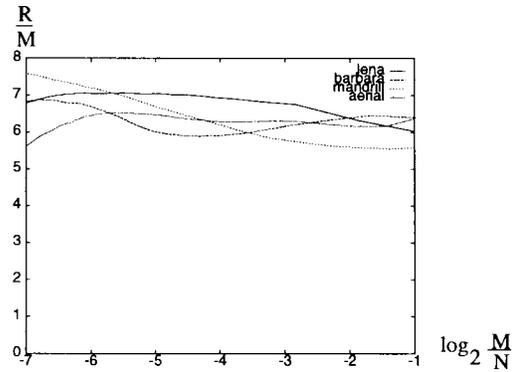


Fig. 7. Variations of R/M for the wavelet transform coding of the test images.

distortion rate $D(\bar{R})$, despite the fact that γ_M is not exactly equal to 1.

The increment of $P_{\text{SNR}}(D)$ for each additional bit is calculated by inserting (26) in (30) while neglecting the variations of K and $r_1 + r_0$

$$\frac{dP_{\text{SNR}}(D)}{d\log_2 \bar{R}} \approx (2\gamma_M - 1)10\log_{10} 2.$$

For Lena, $\gamma_M = 1$ so that $P_{\text{SNR}}(D)$ increases by 3 db for each additional bit, which is indeed verified in Fig. 8. For the three other images, the variations of γ_M cannot be neglected over the whole compression range $\bar{R} \in [2^{-6}, 1]$, and Fig. 8 shows that $P_{\text{SNR}}(D)$ has a slope that varies slowly with $\log_2 \bar{R}$.

Theorem 3 gives an analytical expression that computes $\theta = (T/\Delta)$, which minimizes the distortion D . For $r_1 + r_0 = 6.5$ and $\gamma_M \approx 1$, we get $\theta = 0.81 \approx 1$. This theoretical estimate is precisely the choice that is most often used in wavelet compression softwares after ad-hoc numerical trials.

To optimize the wavelet transform coder, the distortion rate (32) shows that one must choose a wavelet basis that gives precise approximations of images with few wavelet coefficients in order to maintain a small approximation error $D_0(z)$. This essentially depends on the support size and the number of vanishing moments of the wavelets [11]. The optimization of the wavelet basis may depend on the particular class of images to be encoded.

C. JPEG Image Coding

The JPEG image standard decomposes an image in a block cosine basis \mathcal{B} . Images of N^2 pixels are divided in $N^2/64$ blocks of $L = 8$ by 8 pixels. Fig. 4(b) shows that the decay of the sorted local cosine coefficients satisfies the assumptions of Theorem 2 for the four test images.

JPEG uniformly quantizes the block cosine coefficients. In each block of 64 pixels, there is one DC coefficient, which is proportional to the average image value over the block. Instead of directly quantizing this DC coefficient, JPEG quantizes the differences between the DC values of two adjacent blocks. Since the amplitudes of the DC coefficients are not directly quantized, their values are not included in the sorted coefficients shown in Fig. 4(b). A significance map gives the position of zero versus nonzero quantized coefficients.

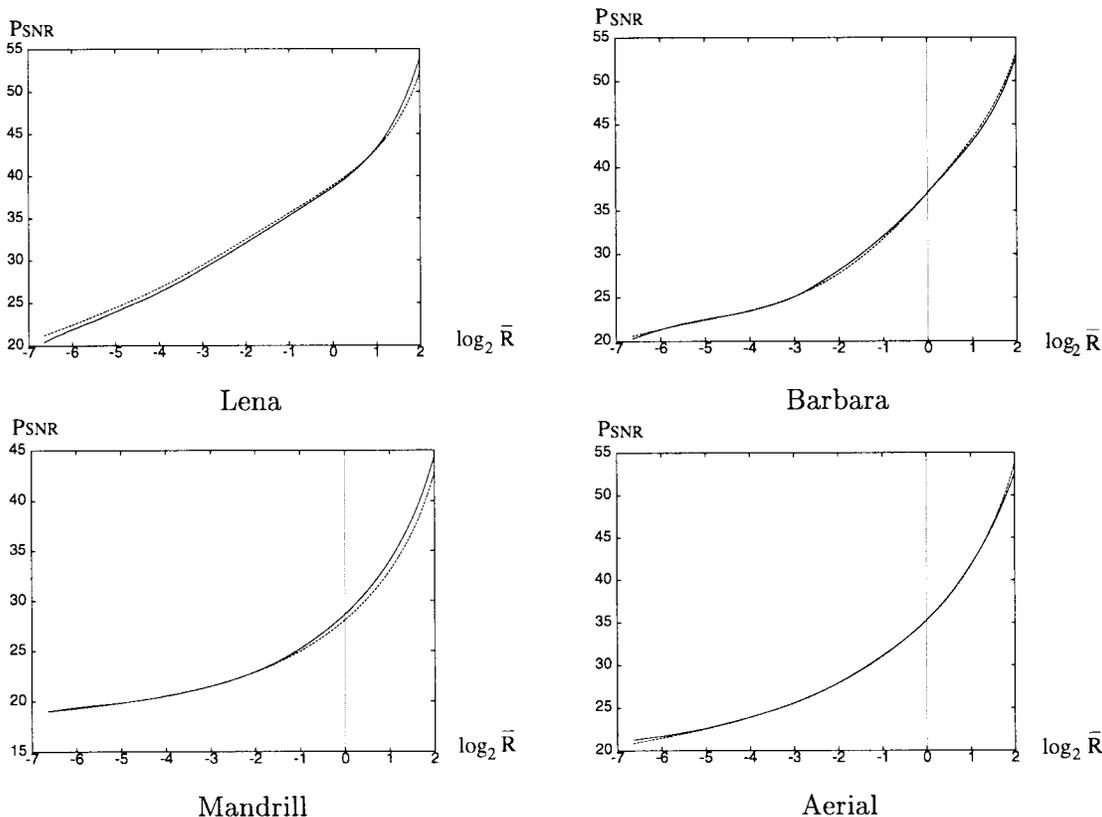


Fig. 8. $PSNR(D)$ (solid line) and $PSNR(\hat{D})$ (dashed line) for a wavelet coder with $\hat{D} = (1 + \frac{1}{12})D_0(\bar{R}/6.5)$.

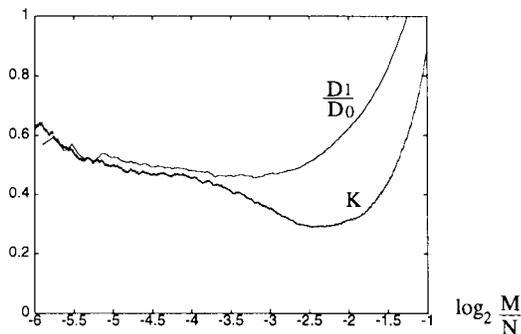


Fig. 9. Comparison of D_1/D_0 with the theoretical estimate $K = (2\gamma_M - 1)/12\theta^2$ for the JPEG coding of Lena.

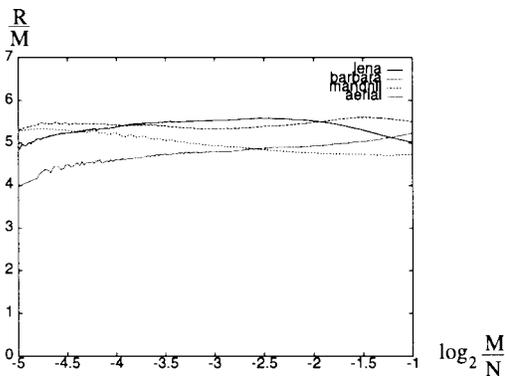


Fig. 10. Variations of R/M for the JPEG coding of the test images.

It is encoded with a run-length coding, which scans the 64 cosine coefficients of each block. The following numerical calculations are done with a baseline JPEG compression coder. Following our explanation in Section II-B, the weights w_m used to improve the visual quality of JPEG images are set to 1, which maintains a uniform quantization with the same bin size Δ for all coefficients.

Since JPEG uniformly quantizes the decomposition coefficients, the zero bin $[-T, T]$ is equal to other bin sizes Δ , and hence, $\theta = (T/\Delta) = \frac{1}{2}$. For Lena, Fig. 9 compares the ratio D_1/D_0 computed numerically and the theoretical estimate derived from Theorem 2

$$K = \frac{2\gamma_M - 1}{12\theta^2} = \frac{2\gamma_M - 1}{3}.$$

The slope γ_M is calculated from the coefficient decay shown in Fig. 4(b). Both curves are very close up to $(M/N) \geq$

2^{-3} . For $(M/N) \geq 2^{-3}$, the estimation error of D_1/D_0 increases because it does not respect the theorem hypothesis that $(M/N) \leq \epsilon$.

JPEG uses a mixed format that encodes the run-length coding for significant coefficients and the amplitude of these significant coefficients together [15]. It is therefore not possible to compute R_0 and R_1 separately. Fig. 10 displays $(R/M) = r_0 + r_1$ as a function of $\log_2(M/N)$ for the four test images. If $(M/N) \in [2^{-5}, 2^{-1}]$, then $r_0 + r_1 \approx 5.5$.

In the compression range of JPEG, the slope γ_M remains slightly above 1. Approximating $\gamma_M \approx \frac{5}{4}$ yields $K \approx \frac{1}{2}$. The distortion rate formula calculated with (23)

$$D(\bar{R}) = (1 + K)D_0\left(\frac{\bar{R}}{r_0 + r_1}\right)$$

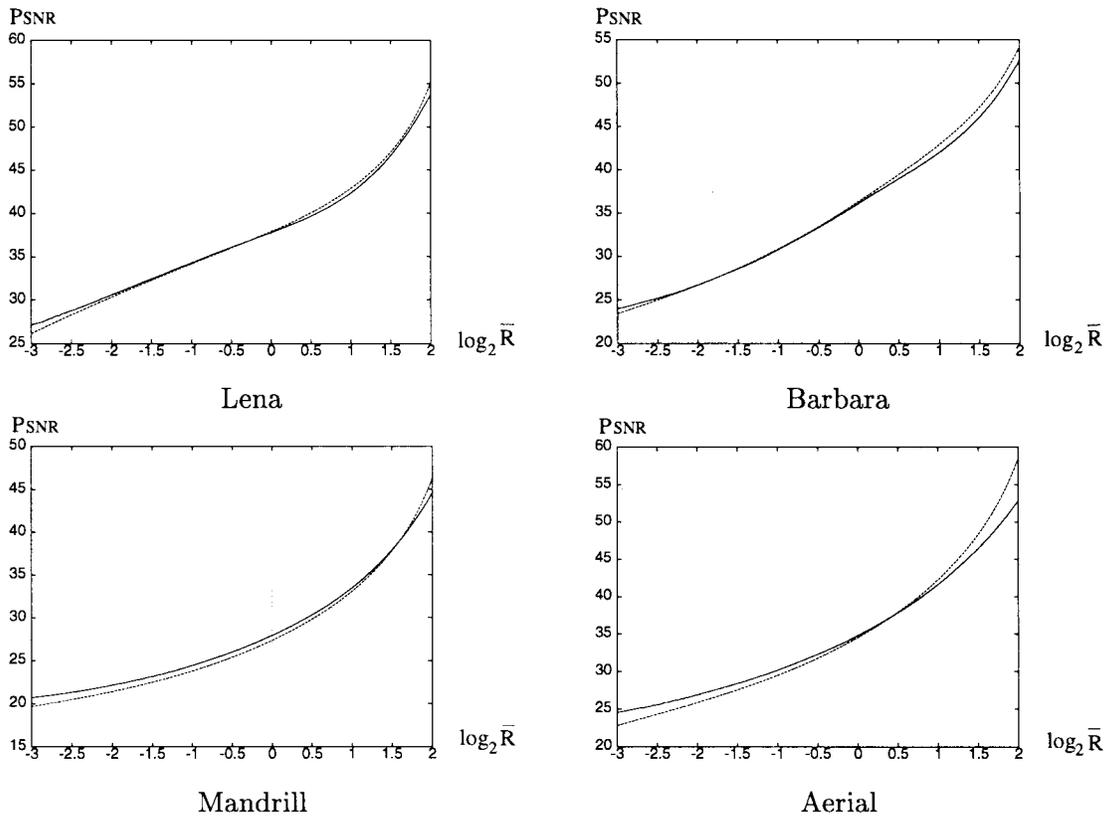


Fig. 11. $P_{\text{SNR}}(D)$ (solid line) and $P_{\text{SNR}}(\hat{D})$ (dashed line) for JPEG, with $\hat{D} = (1 + \frac{1}{2})D_0(\bar{R}/5.5)$.

is thus approximated by

$$\hat{D}(\bar{R}) = \left(1 + \frac{1}{2}\right)D_0\left(\frac{\bar{R}}{5.5}\right). \quad (33)$$

Fig. 11 compares $P_{\text{SNR}}(D)$ calculated numerically with JPEG software and its approximation $P_{\text{SNR}}(\hat{D})$ derived from (33). It shows that the distortion rate formula (33) is a precise approximation over the whole compression range of JPEG and that we can neglect the variations of K and $r_0 + r_1$.

IV. EMBEDDED TRANSFORM CODING

For rapid transmission or fast image browsing from a database, one should provide a coarse signal approximation quickly and then progressively enhance it as more bits are transmitted. Embedded coders offer this flexibility by grouping the bits in order of significance. The decomposition coefficients are sorted, and the first bits of the largest coefficients are sent first. An image approximation can be reconstructed at any time from the bits already transmitted. Embedded coders can take advantage of any prior information on the location of large versus small coefficients. Such prior information is available for natural images decomposed in wavelet bases. As a result, an implementation with zero trees designed by Shapiro [13] yields better compression rates than classical wavelet transform coders.

The decomposition coefficients $a[m] = \langle f, g_m \rangle$ are partially ordered by grouping them in sets \mathcal{S}_k of indexes defined for any $k \in \mathbf{Z}$ by

$$\mathcal{S}_k = \{m: 2^k \leq |a[m]| < 2^{k+1}\}.$$

The set \mathcal{S}_k is encoded with a binary significance map $b_k[m]$

$$b_k[m] = \begin{cases} 0 & \text{if } m \notin \mathcal{S}_k \\ 1 & \text{if } m \in \mathcal{S}_k. \end{cases} \quad (34)$$

An embedded algorithm quantizes $a[m]$ uniformly with a quantization step (bin size) $\Delta = 2^n$ that is progressively reduced. Let $m \in \mathcal{S}_k$ with $k \geq n$. The amplitude $|Q(a[m])|$ of the quantized number is represented in base 2 by a binary string with nonzero digits between the bit k and the bit n . The bit k is necessarily 1 because $2^k \leq |Q(a[m])| < 2^{k+1}$. Hence, $k - n$ bits are sufficient to specify this amplitude to which one bit is added for the sign.

The embedded coding is initiated with the largest quantization step to produce at least one nonzero quantized coefficient. To refine the quantization step from 2^{n+1} to 2^n , the algorithm records the significance map $b_n[m]$ and the sign of $a[m]$ for $m \in \mathcal{S}_n$. This can be done by directly recording the sign of significant coefficients with a variable incorporated into the significance map $b_n[m]$. Afterwards, the code stores the bit n of all amplitudes $|Q(a[m])|$ for $m \in \mathcal{S}_k$ with $k > n$. If necessary, the coding precision is improved by decreasing n and continuing the encoding. The different steps of the algorithm can be summarized as follows [12].

- 1) Store the index n of the first nonempty set \mathcal{S}_n , where $n = \lfloor \sup_m \log_2 |a[m]| \rfloor$.
- 2) Store the significance map $b_n[m]$ and the sign of $a[m]$ for $m \in \mathcal{S}_n$.
- 3) Store the n th bit of all coefficients $|a[m]| > 2^{n+1}$. These are coefficients that belong to some set \mathcal{S}_k for $k > n$,

whose coordinates were already stored. Their n th bit is stored in the order in which their position was recorded in the previous passes.

4) Decrease n by 1, and go to step 2.

This algorithm may be stopped at any time in the loop, providing a code for any specified number of bits. The distortion rate is analyzed when the algorithm is stopped at step 4. All coefficients above $T = 2^n$ are uniformly quantized with a bin size $\Delta = 2^n$. The zero quantization bin $[-T, T]$ is therefore twice as big as the other quantization bins. This quantizer is the same as in the direct transform coding studied in Section III-A for a zero-bin ratio $\theta = T/\Delta = 1$. This value was shown to be nearly optimal for wavelet image coders. The total distortion $D = D_0 + D_1$ is therefore not modified by the embedding strategy.

Once the algorithm stops, we denote by M the number of significant coefficients above $T = 2^n$. The total number of bits of the embedded code is $R = R_0^e + R_1^e$, where R_0^e is the number of bits needed to encode all significance maps $b_k[m]$ for $k \geq n$, and R_1^e is the number of bits used to encode the amplitudes of the quantized significant coefficients $Q(a[m])$, knowing that $m \in \mathcal{S}_k$ for $k > n$.

To appreciate the efficiency of this embedding strategy, the bit budget $R_0^e + R_1^e$ is compared with the number of bits $R_0 + R_1$ used by the direct transform coder of Section III-A. The value R_0 is the number of bits used to encode the overall significance map

$$b[m] = \begin{cases} 0 & \text{if } |a[m]| \leq T \\ 1 & \text{if } |a[m]| > T \end{cases} \quad (35)$$

and R_1 is the number of bits used to encode the quantized significant coefficients.

An embedded strategy encodes $Q(a[m])$ knowing that $m \in \mathcal{S}_k$ and, hence, that $2^k \leq |Q(a[m])| < 2^{k+1}$, whereas a direct transform coding knows only that $|Q(a[m])| > T = 2^n$. Thus, fewer bits are needed for embedded codes: $R_1^e \leq R_1$. This improvement may be offset, however, by the supplement of bits needed to encode the significance maps $\{b_k[m]\}_{k > n}$ of the sets $\{\mathcal{S}_k\}_{k > n}$. A direct transform coder records a single significance map $b[m]$, which specifies $\cup_{k \geq n} \mathcal{S}_k$. It provides less information and is therefore encoded with fewer bits: $R_0^e \geq R_0$. An embedded coder brings an improvement over a direct transform coder if

$$R_0^e + R_1^e \leq R_0 + R_1.$$

This can happen if we have some prior information about the position of large decomposition coefficients $Q(a[m])$ versus smaller ones. It allows us to reduce the number of bits needed to encode the partial sorting of all coefficients provided by the significance maps $\{b_k[m]\}_{k > n}$. The use of such prior information produces an overhead of R_0^e relative to R_0 that is smaller than the gain of R_1^e relative to R_1 . Let $x(z)$ be the sorted amplitudes of the coefficients $a[m]$. The following theorem computes R_1^e with the same hypotheses as Theorem 2, allowing us to compare it with R_1 .

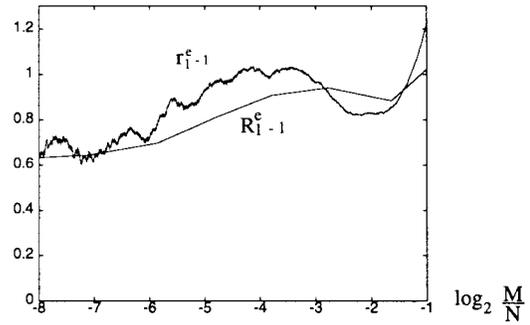


Fig. 12. Comparison of $(R_1^e/M) - 1$ with the theoretical estimate $r_1^e - 1 = 1/(2^{1/\gamma_M} - 1)$ for the embedded coding of Lena.

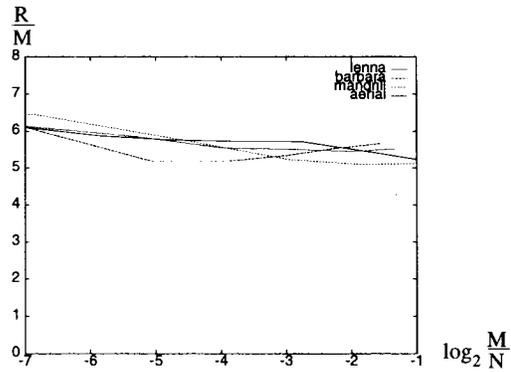


Fig. 13. Variations of R/M for an embedded wavelet coding of the test images.

Theorem 4: Suppose that $x(z)$ satisfies (19)–(22). Let $\gamma_M = \gamma(M/N) > \frac{1}{2}$. If $(M/N) \leq \epsilon$ and $M \geq (1/\epsilon)$, then

$$r_1^e = \frac{R_1^e}{M} = 1 + \frac{1}{2^{1/\gamma_M} - 1} [1 + O(\epsilon \log_2 \epsilon)] \quad (36)$$

and

$$D(\bar{R}) = (1 + K)D_0 \left(\frac{\bar{R}}{r_1^e + r_0^e} \right) \quad (37)$$

where $r_0^e = (R_0^e/M)$ and $K = (D_1/D_0)$ is given by (24).

The proof of this theorem is in Appendix A-3. In the following, we omit the corrective terms to simplify the notation. This theorem proves that R_1^e/M is well approximated by

$$r_1^e = 1 + \frac{1}{2^{1/\gamma_M} - 1}. \quad (38)$$

Fig. 12 verifies that the value of $(R_1^e/M) - 1$ calculated numerically with the embedded wavelet coding software of Said and Pearlman [12] is close to the estimate $r_1^e - 1$ calculated by computing γ_M from Fig. 4(a). We subtract 1 bit because Said and Pearlman do not encode the sign bits with the amplitudes of the significant coefficients but store their values in the significance maps.

Let us compare $r_1^e = (R_1^e/M)$ calculated in (36) with the value $r_1 = (R_1/M)$ estimated in (25) for a direct transform coding, with $\theta = 1$. The bit budget of an embedded coding is smaller than that of a direct transform coding if $r_0^e + r_1^e \leq r_0 + r_1$, and hence

$$\begin{aligned} r_0^e - r_0 &\leq r_1 - r_1^e \\ &= (1 + \gamma_M) \log_2 \exp + \log_2 \gamma_M - \frac{1}{2^{1/\gamma_M} - 1}. \end{aligned} \quad (39)$$

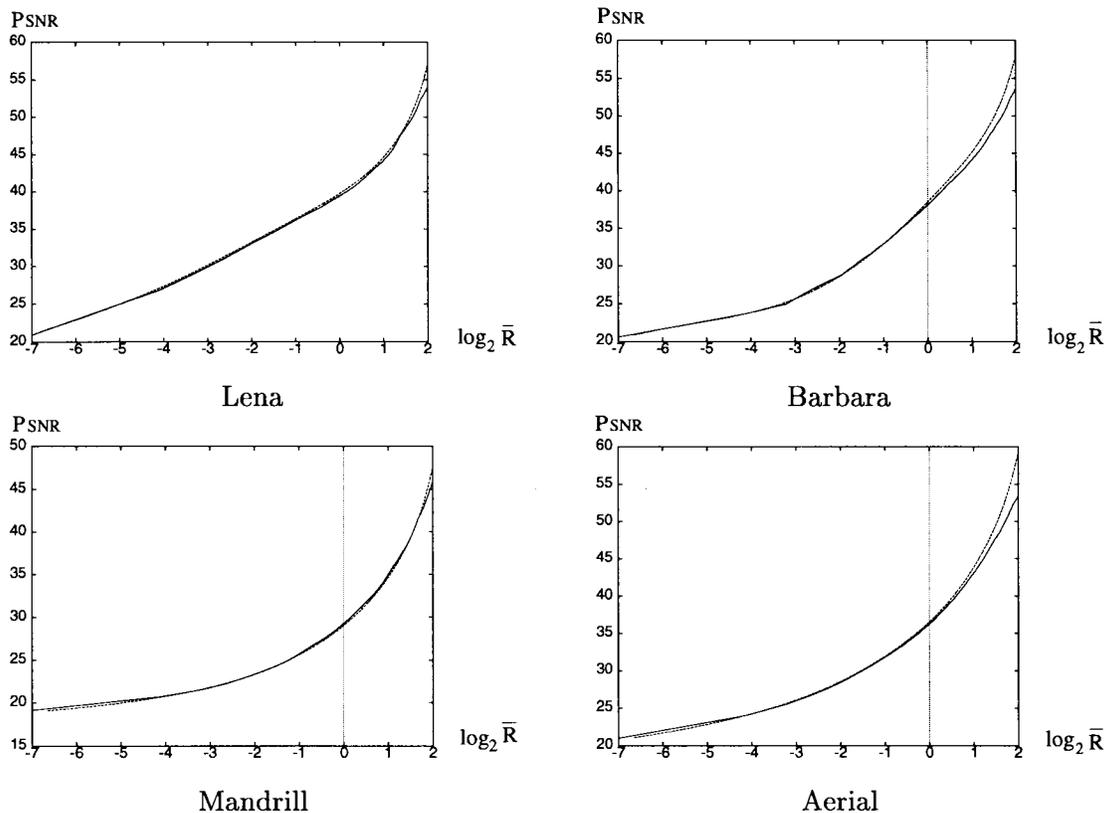


Fig. 14. $P_{\text{SNR}}(D)$ (solid line) and $P_{\text{SNR}}(\hat{D})$ (dashed line) for an embedded wavelet coding, with $\hat{D} = (1 + \frac{1}{12})D_0(\bar{R}/5.5)$.

If $\gamma_M \approx 1$, then $r_1 - r_1^e \approx 1.9$. The inequality (39) is satisfied for embedded transform codings implemented in wavelet bases [13] and in a block cosine basis [17] by taking advantage of prior knowledge of the location of large versus small coefficients using zero trees.

A wavelet coefficient $\langle f, \psi_{j,p,q}^k \rangle$ has a large amplitude where the signal has sharp transitions. If an image f is Lipschitz α in the neighborhood of (x_0, y_0) , then for wavelets $\psi_{j,p,q}^k$ located in this neighborhood, one can prove [11] that there exists $A \geq 0$ such that

$$|\langle f, \psi_{j,p,q}^k \rangle| \leq A 2^{j(\alpha+1)}.$$

The worst singularities are often discontinuities, which means that $\alpha \geq 0$. In the neighborhood of singularities without oscillations, the amplitudes of wavelet coefficients thus decrease when the scale 2^j decreases. This property is not valid for oscillatory patterns. High-frequency oscillations create coefficients at large scales 2^j that are typically smaller than those at the fine scale that matches the period of oscillation. Such oscillatory patterns are not often encountered in images, although they do appear as thin lines in the Barbara image.

Wavelet zero trees, which were introduced by Lewis and Knowles [9], take advantage of the decay of wavelet coefficients by relating these coefficients across scales with quad trees. These zero trees take advantage of a partial self-similarity of the image [2]. Shapiro [13] used this zero-tree structure to encode the embedded significance maps of wavelet coefficients. Numerical experiments were performed with Said and Pearlman's software [12], which improves Shapiro's zero-

tree coder with a set partitioning technique. A cubic spline orthogonal wavelet basis was used. Fig. 13 displays the value of $(R/M) = r_0^e + r_1^e$ as a function of $\log_2 M/N$ for the four test images. These curves have variations centered at 5.5. This graph should be compared with Fig. 7, which shows $r_0 + r_1$ calculated with a direct wavelet image coder. An improvement of approximately 1 bit per significant coefficient is obtained.

The embedded distortion rate function is calculated with (37). Inserting $\gamma_M \approx 1$ in (24) yields $K \approx \frac{1}{12}$. Since $r_0^e + r_1^e \approx 5.5$, we get an approximate distortion rate formula

$$\hat{D}(\bar{R}) = \left(1 + \frac{1}{12}\right) D_0\left(\frac{\bar{R}}{5.5}\right).$$

Fig. 14 compares the $P_{\text{SNR}}(D)$ calculated numerically for the our test images and its theoretical approximation $P_{\text{SNR}}(\hat{D})$. Once more, we verify that the distortion rate essentially depends only upon the approximation error function $D_0(z)$. The variations of the constants K and $r_0^e + r_1^e$ can be neglected. The embedding strategy reduces $r_0 + r_1$ to $r_0^e + r_1^e$, but the variations of the distortion rate still depends on the nonlinear approximation error $D_0(z)$.

V. CONCLUSION

We have shown that at low bit rates, the distortion rate of transform coders can be computed by separating the coefficients quantized to zero from all others. The resulting distortion rate $D(\bar{R})$ depends crucially on the precision of nonlinear image approximations with few nonzero basis coefficients. In wavelet and block cosine bases, we have demonstrated that

if $\bar{R} \leq 1$, then $D(\bar{R})$ decays like $C\bar{R}^{1-2\gamma}$, where γ is an exponent of the order of 1, which varies slowly as a function of $\log_2 \bar{R}$. Embedded transform coders improve the constant C but not the exponent γ , which depends on the image.

APPENDIX
PROOFS OF THEOREMS

A.1. Proof of Theorem 2

The distortion rate formula (23) is derived by observing that $D = D_0 + D_1$ and $\bar{R} = (R_0/N) + (R_1/N)$. By definition, $K = (D_1/D_0)$, $r_0 = (R_0/M)$, and $r_1 = (R_1/M)$. Inserting these variables in the equation $D = D_0 + D_1$ yields (23).

The main difficulty is in computing $K = (D_1/D_0)$ and $r_1 = (R_1/M)$. For this purpose, $x(z)$ is approximated by a function $x_M(z)$ that is tangential to $x(z)$ at $z = M/N$. We first estimate K and r_1 by replacing $x(z)$ with $x_M(z)$ and then evaluate the error introduced by this approximation. We define

$$x_M(z) = T \left(\frac{Nz}{M} \right)^{-\gamma_M}.$$

Since there are exactly M coefficients above T , $x(M/N) = T = x_M(M/N)$, and

$$\frac{d \log_2 x \left(\frac{M}{N} \right)}{d \log_2 z} = \frac{d \log_2 x_M \left(\frac{M}{N} \right)}{d \log_2 z} = -\gamma_M.$$

Both curves are thus tangential at $z = (M/N)$, and the concavity of $\log_2 x(z)$ guarantees that

$$\log_2 x(z) \leq \log_2 x_M(z).$$

Let us now prove that (D_1/D_0) is given by (24). We compute

$$D_0 = \sum_{k=M}^{N-1} x^2 \left(\frac{k}{N} \right) \approx N \int_{M/N}^1 x^2(z) dz. \quad (40)$$

This value is approximated by

$$\begin{aligned} \hat{D}_0 &= N \int_{M/N}^1 x_M^2(z) dz = T^2 N \int_{M/N}^1 \left(\frac{Nz}{M} \right)^{-2\gamma_M} dz \\ &= \frac{T^2 M}{2\gamma_M - 1} \left(1 + \left(\frac{M}{N} \right)^{2\gamma_M - 1} \right). \end{aligned} \quad (41)$$

The concavity of $\log_2 x(z)$ implies that $x(z) \leq x_M(z)$ for $z \geq (M/N)$ and, hence, that $\hat{D}_0 - D_0 \geq 0$. To compute $\hat{D}_0 - D_0$, we write

$$x(z) = T \left(\frac{Nz}{M} \right)^{-\beta(z)} \quad (42)$$

with

$$\beta(z) = \frac{\log_2 x \left(\frac{M}{N} \right) - \log_2 x(z)}{\log_2 z - \log_2 \frac{M}{N}}.$$

The estimation error is

$$\begin{aligned} \hat{D}_0 - D_0 &= NT^2 \int_{M/N}^1 \left(\frac{Nz}{M} \right)^{-2\gamma_M} \\ &\quad \cdot \left[1 - \left(\frac{Nz}{M} \right)^{-2(\beta(z) - \gamma_M)} \right] dz. \end{aligned}$$

We cut this integral in two parts

$$|\hat{D}_0 - D_0| \leq I_1 + I_2 \quad (43)$$

with

$$\begin{aligned} I_1 &= NT^2 \int_{M/N}^{(1/2\epsilon)(M/N)} \left(\frac{Nz}{M} \right)^{-2\gamma_M} \\ &\quad \cdot \left[1 - \left(\frac{Nz}{M} \right)^{-2(\beta(z) - \gamma_M)} \right] dz \end{aligned} \quad (44)$$

and

$$\begin{aligned} I_2 &= NT^2 \int_{(1/2\epsilon)(M/N)}^{+\infty} \left(\frac{Nz}{M} \right)^{-2\gamma_M} dz \\ &= \frac{MT^2 (2\epsilon)^{2\gamma_M - 1}}{2\gamma_M - 1}. \end{aligned} \quad (45)$$

To find an upper bound of I_1 , we compute an upper bound of $\beta(z) - \gamma_M$. Clearly

$$\beta \left(\frac{M}{N} \right) = \gamma \left(\frac{M}{N} \right) = \gamma_M.$$

Since $\log_2 x(z)$ is concave, $|\gamma_M - \beta(z)| \leq |\gamma_M - \gamma(z)|$. The hypothesis (19) guarantees that

$$\forall z \in [0, 2^{-1}], \quad \left| \frac{d\gamma(z)}{d \log_2 z} \right| \leq \epsilon$$

and hence

$$|\gamma_M - \gamma(z)| \leq \epsilon \left| \log_2 z - \log_2 \frac{M}{N} \right| \quad (46)$$

which yields

$$\forall z \in [0, 2^{-1}], \quad |\gamma_M - \beta(z)| \leq \epsilon \left| \log_2 \frac{Nz}{M} \right|. \quad (47)$$

The concavity of $\log_2 x(z)$ also implies that $\beta(z) - \gamma_M \geq 0$ for $z \geq (M/N)$. Since $(1/2\epsilon)(M/N) \leq 2^{-1}$, it follows that

$$\begin{aligned} I_1 &\leq NT^2 \int_{M/N}^{(1/2\epsilon)(M/N)} \left(\frac{Nz}{M} \right)^{-2\gamma_M} \\ &\quad \cdot \left[1 - \left(\frac{Nz}{M} \right)^{-2\epsilon \log_2(Nz/M)} \right] dz \\ &\leq NT^2 \int_{M/N}^{(1/2\epsilon)(M/N)} \left(\frac{Nz}{M} \right)^{-2\gamma_M} \\ &\quad \cdot [1 - 2^{-2\epsilon(\log_2(Nz/M))^2}] dz \\ &\leq MT^2 \int_1^{(1/2\epsilon)} z^{-2\gamma_M} [1 - 2^{-2\epsilon(\log_2 z)^2}] dz. \end{aligned}$$

One can then derive that

$$I_1 \leq \frac{MT^2}{2\gamma_M - 1} \mathcal{O}(\epsilon |\log_2 \epsilon|^2). \quad (48)$$

Inserting (41), (45), and (48), in (43) proves (24).

Let us now prove that (R_1/M) satisfies (25). To compute R_1 , we must calculate the differential entropy

$$\mathcal{H}_d(X_T) = - \int_{-\infty}^{+\infty} p_T(x) \log_2 p_T(x) dx$$

where $p_T(x) = (N/M)p(x)\mathbf{1}_{\{|x|>T\}}(x)$. Since $p(x) = p(-x)$

$$\mathcal{H}_d(X_T) = -2\frac{N}{M} \int_T^{+\infty} p(x) \log_2 p(x) dx. \quad (49)$$

This integral is calculated by relating it to $x(z)$. Since

$$z(x) = 1 - \int_{-x}^x p(u) du \quad (50)$$

it follows that $z'(x) = -p(x) - p(-x) = -2p(x)$, and hence

$$x'(z) = \frac{1}{z'(x)} = \frac{-1}{2p(x)}. \quad (51)$$

Since $p(x) dx = -\frac{1}{2} dz$, the change of variable $z = z(x)$ in (49) yields

$$\mathcal{H}_d(X_T) = 2\frac{N}{M} \int_0^{M/N} \log_2 |2x'(z)| \frac{dz}{2}$$

and with a change of variable

$$\mathcal{H}_d(X_T) = 1 + \int_0^1 \log_2 \left(\frac{M}{N} \left| x' \left(\frac{M}{N} z \right) \right| \right) dz. \quad (52)$$

This integral is first estimated by replacing $x(z)$ by $x_M(z)$

$$\begin{aligned} \hat{\mathcal{H}}_d &= 1 + \int_0^1 \log_2 \left(\frac{M}{N} \left| x'_M \left(\frac{M}{N} z \right) \right| \right) dz \\ &= 1 + \int_0^1 \log_2 (T\gamma_M z^{-\gamma_M-1}) dz \\ &= 1 + \log_2 T + \log_2 \gamma_M + (\gamma_M + 1) \log_e 2. \end{aligned} \quad (53)$$

Let us compute the error

$$|\hat{\mathcal{H}}_d - \mathcal{H}_d(X_T)| \leq \left| \int_0^1 \log_2 \left(\frac{x' \left(\frac{M}{N} z \right)}{x'_M \left(\frac{M}{N} z \right)} \right) dz \right|.$$

Observe that

$$\gamma(z) = -\frac{d \log_2 x(z)}{d \log_2 z} = \frac{-zx'(z)}{x(z)}$$

and $\gamma_M = -(zx'_M(z)/x_M(z))$ so that

$$\begin{aligned} &|\hat{\mathcal{H}}_d - \mathcal{H}_d(X_T)| \\ &\leq \int_0^1 \left| \log_2 \left(\frac{\gamma \left(\frac{M}{N} z \right)}{\gamma_M} \right) + \log_2 \left(\frac{x_M \left(\frac{M}{N} z \right)}{x \left(\frac{M}{N} z \right)} \right) \right| dz \\ &\leq \int_0^1 \left| \log_2 \left(\frac{\gamma \left(\frac{M}{N} z \right)}{\gamma_M} \right) + \left(\gamma_M - \beta \left(\frac{M}{N} z \right) \right) \log_2 z \right| \cdot dz. \end{aligned} \quad (54)$$

Since $\inf_{z \in [0,1]} \gamma(z) = \mu > 0$, it follows that

$$\begin{aligned} \left| \log_2 \left(\frac{\gamma \left(\frac{M}{N} z \right)}{\gamma_M} \right) \right| &= \left| \log_2 \left(1 + \frac{\gamma_M - \gamma \left(\frac{M}{N} z \right)}{\gamma \left(\frac{M}{N} z \right)} \right) \right| \\ &\leq \frac{1}{\mu} \left| \gamma \left(\frac{M}{N} z \right) - \gamma_M \right|. \end{aligned}$$

Since $(M/N)z \leq 2^{-1}$, we proved in (46) and (47) that

$$\begin{aligned} \left| \gamma \left(\frac{M}{N} z \right) - \gamma_M \right| &\leq \epsilon |\log_2 z| \quad \text{and} \\ \left| \gamma_M - \beta \left(\frac{M}{N} z \right) \right| &\leq \epsilon |\log_2 z|. \end{aligned}$$

By inserting these inequalities in (54), one can verify that $|\hat{\mathcal{H}}_d - \mathcal{H}_d(X_T)| = \mathcal{O}(\epsilon)$. From (53) and $T = \theta\Delta$, we thus derive that

$$\begin{aligned} \frac{R_1}{M} &= \mathcal{H}_d(X_T) - \log_2 \Delta \\ &= \log_2 \theta + 1 + (1 + \gamma_M) \log_2 e + \log_2 \gamma_M + \mathcal{O}(\epsilon) \end{aligned}$$

which finishes the proof of (25).

Let us finally prove (26). We must calculate

$$\frac{d \log_2 D_0 \left(\frac{M}{N} \right)}{d \log_2 z} = \frac{d D_0 \left(\frac{M}{N} \right)}{dz} \frac{M}{N D_0 \left(\frac{M}{N} \right)}. \quad (55)$$

We derive from (40) that

$$\frac{d D_0 \left(\frac{M}{N} \right)}{dz} = -N x^2 \left(\frac{M}{N} \right) = -N T^2$$

and (24) yields

$$\frac{1}{D_0} = \frac{2\gamma_M - 1}{M T^2} [1 + \mathcal{O}(\epsilon |\log_2 \epsilon|^2 + \epsilon^{2\gamma_M-1})].$$

Inserting these last two equations in (55) gives (26).

A.2. Proof of Theorem 3

We proved in (24) that

$$D = D_1 \left(1 + \frac{12\theta^2}{2\gamma_M - 1} \right).$$

Since $D_1 = (M\Delta^2/12) = (MT^2/12\theta^2)$ and $T = x(M/N)$

$$\begin{aligned} D &= M T^2 \left(\frac{1}{\theta^2} + \frac{12}{2\gamma_M - 1} \right) \\ &= M x^2 \left(\frac{M}{N} \right) \left(\frac{1}{\theta^2} + \frac{12}{2\gamma_M - 1} \right). \end{aligned}$$

We decompose $(R/M) = (R_0/M) + (R_1/M)$, where $(R_0/M) = r_0$ is a constant, and (25) shows that

$$\begin{aligned} \frac{R_1}{M} &= r_1(M, \theta) = 1 + (1 + \gamma_M) \log_2 e + \log_2 \gamma_M \\ &\quad + \log_2 \theta + \mathcal{O}(\epsilon |\log_2 \epsilon|). \end{aligned} \quad (56)$$

If we neglect the residual terms

$$\frac{R}{M} = 1 + r_0 + (1 + \gamma_M) \log_2 e + \log_2 \gamma_M + \log_2 \theta.$$

The variable M depends on (R, θ) , and we can thus write $(R/M) = \beta(R, \theta)$. Hence

$$D(\bar{R}, \theta) = \frac{N\bar{R}}{\beta(R, \theta)} x^2 \left(\frac{\bar{R}}{\beta(R, \theta)} \right) \left(\frac{1}{\theta^2} + \frac{12}{2\gamma_M - 1} \right). \quad (57)$$

To compute $\partial D(\bar{R}, \theta) / \partial \theta$, we need to calculate $\partial \beta(R, \theta) / \partial \theta$. Observe that

$$\frac{\partial \beta(R, \theta)}{\partial \theta} = \frac{\partial \frac{R}{M}}{\partial \theta} = -\frac{R}{M^2} \frac{\partial M(R, \theta)}{\partial \theta}.$$

We compute $\partial M(R, \theta) / \partial \theta$ by taking the derivative with respect to θ at R fixed of the equality

$$R = M[1 + r_0 + (1 + \gamma_M) \log_2 e + \log_2 \gamma_M + \log_2 \theta].$$

We get

$$0 = \frac{\partial M(R, \theta)}{\partial \theta} \left[\frac{R}{M} + M \log_2 e \frac{d\gamma_M}{dM} + M \frac{1}{\gamma_M} \frac{d\gamma_M}{dM} \right] + M \frac{\log_2 e}{\theta}.$$

The slow variation condition (19) on the slope imposes that

$$\left| \frac{d\gamma_M}{d \log_2 \left(\frac{M}{N} \right)} \right| \leq \epsilon.$$

Inserting this in the previous equation proves that

$$\frac{\partial M(R, \theta)}{\partial \theta} \left[\frac{R}{M} + \mathcal{O}(\epsilon) \right] = -M \frac{\log_2 e}{\theta}.$$

Hence

$$\begin{aligned} \frac{\partial \beta(R, \theta)}{\partial \theta} &= -\frac{R}{M^2} \frac{\partial M(R, \theta)}{\partial \theta} \\ &= \frac{\log_2 e}{\theta} [1 + \mathcal{O}(\epsilon)] = \beta'. \end{aligned} \quad (58)$$

Let us now compute a derivative from (57)

$$\begin{aligned} \frac{\partial \log_e D(R, \theta)}{\partial \theta} &= -\frac{\beta'}{\beta(R, \theta)} + 2 \left(-\frac{R\beta'}{\beta^2(R, \theta)} \right) \\ &\quad \cdot \frac{x' \left(\frac{M}{N} \right)}{Nx \left(\frac{M}{N} \right)} - \frac{\frac{2}{\theta^3}}{\frac{1}{\theta^2} + \frac{12}{2\gamma_M - 1}} = 0. \end{aligned} \quad (59)$$

We know that

$$\frac{d \log_2 x \left(\frac{M}{N} \right)}{d \log_2 \left(\frac{M}{N} \right)} = \frac{Mx' \left(\frac{M}{N} \right)}{Nx \left(\frac{M}{N} \right)} = -\gamma_M. \quad (60)$$

Inserting (58) and (60) with $\beta(R, \theta) = (R/M) = r_1 + r_0$ in (59) proves (30).

A.3. Proof of Theorem 4

Suppose that the algorithm stops at n . It performs a quantization with intervals of size $T = 2^n$. For each coefficient $2^k \leq |a[m]| < 2^{k+1}$, we saw that the number of bits required to specify the quantized value of $|a[m]|$ is

$$p = k - n = \left\lfloor \log_2 \frac{|a[m]|}{T} \right\rfloor$$

and we need one bit for the sign of $a[m]$. The coefficients above 2^n are the M coefficients whose amplitudes are given by $x(z)$. Therefore

$$R_1^\epsilon = \sum_{m=1}^M \left\lfloor \log_2 \frac{x \left(\frac{m}{N} \right)}{T} \right\rfloor + M = I_1 + I_2 + M \quad (61)$$

where I_1 corresponds to the sum for $m \geq \epsilon M \geq 1$

$$I_1 = \sum_{m=M\epsilon}^M \left\lfloor \log_2 \left(\frac{M}{m} \right)^{\beta(m/N)} \right\rfloor$$

and I_2 is calculated from $m < \epsilon M$

$$I_2 = \sum_{m=1}^{M\epsilon-1} \left\lfloor \log_2 \left(\frac{M}{m} \right)^{\beta(m/N)} \right\rfloor.$$

We proved in (47) that

$$\left| \gamma_M - \beta \left(\frac{m}{N} \right) \right| \leq \epsilon \left| \log_2 \left(\frac{m}{M} \right) \right|. \quad (62)$$

As a consequence, for $m \in [\epsilon M, M]$

$$\gamma_M - \epsilon \left| \log_2 \epsilon \right| = \beta_0 \leq \beta \left(\frac{m}{N} \right) \leq \beta_1 = \gamma_M + \epsilon \left| \log_2 \epsilon \right|. \quad (63)$$

Let us define

$$I(\beta) = \sum_{m=M\epsilon}^M \left\lfloor \log_2 \left(\frac{M}{m} \right)^\beta \right\rfloor.$$

The inequalities (63) imply that

$$I(\beta_0) \leq I_1 \leq I(\beta_1). \quad (64)$$

Let us calculate $I(\beta)$ for any β . We decompose the sum in slices where

$$2^i \leq \left(\frac{M}{m} \right)^\beta < 2^{i+1}$$

for which the floor of the \log_2 is equal to i . It corresponds to

$$M2^{-(i+1)/\beta} < m \leq M2^{-i/\beta}.$$

Let $a = \lfloor \log_2 \epsilon \rfloor$. We obtain

$$I(\beta) = \sum_{i=0}^{a\beta} i M(2^{-i/\beta} - 2^{-(i+1)/\beta}) \quad (65)$$

so that

$$\frac{I(\beta)}{M} = (1 - 2^{-(1/\beta)}) \sum_{i=0}^{a\beta} i 2^{-(i/\beta)}. \quad (66)$$

We can verify that

$$\sum_{i=0}^{+\infty} i2^{-(i/\beta)} = \frac{2^{-(1/\beta)}}{(1-2^{-(1/\beta)})^2}. \quad (67)$$

Approximating the finite sum (65) by this infinite sum yields

$$\frac{I(\beta)}{M} = \frac{1}{2^{1/\beta} - 1} + \mathcal{O}(\epsilon |\log_2 \epsilon|). \quad (68)$$

The inequality (64) can thus be rewritten

$$\frac{1}{2^{1/\beta_0} - 1} - \mathcal{O}(\epsilon |\log_2 \epsilon|) \leq \frac{I_1}{M} \leq \frac{1}{2^{1/\beta_1} - 1} + \mathcal{O}(\epsilon |\log_2 \epsilon|).$$

Since $\beta_1 - \gamma_M = \gamma_M - \beta_0 = \epsilon |\log_2 \epsilon|$, it follows that

$$\left| \frac{I_1}{M} - \frac{1}{2^{(1/\gamma_M)} - 1} \right| \leq \mathcal{O}(\epsilon |\log_2 \epsilon|). \quad (69)$$

To compute the discrete sum I_2 , we use (62), which proves that

$$\beta\left(\frac{m}{N}\right) \leq \gamma_M + \epsilon \left| \log_2 \frac{m}{M} \right|.$$

Hence

$$\begin{aligned} I_2 &= \sum_{m=1}^{M\epsilon-1} \left[\log_2 \left(\frac{M}{m} \right)^{\gamma_M + \epsilon |\log_2 (m/M)|} \right] \\ &\leq M \int_0^\epsilon (\gamma_M + \epsilon |\log_2 z|) |\log_2 z| dz \\ &= M \mathcal{O}(\epsilon |\log_2 \epsilon|). \end{aligned}$$

We thus derive from (61) and (69) that

$$\left| \frac{R_1^\epsilon}{M} - \frac{1}{2^{(1/\gamma_M)} - 1} - 1 \right| = \mathcal{O}(\epsilon |\log_2 \epsilon|).$$

REFERENCES

- [1] A. Cohen, I. Daubechies, O. Guleryuz, and M. Orchard, "On the importance of combining wavelet-based nonlinear approximation with coding strategies," 1997, preprint.
- [2] G. M. Davis, "A wavelet-based analysis of fractal image compression," *IEEE Trans. Image Processing*, vol. 6, 1997.
- [3] R. A. DeVore, B. Jawerth, and B. J. Lucier, "Image compression through wavelet transform coding," *IEEE Trans. Inform. Theory*, vol. 38, pp. 719–746, Mar. 1992.
- [4] N. Farvardin and J. Modestino, "Optimum quantizer performances for a class of non-Gaussian memoryless sources," *IEEE Trans. Inform. Theory*, vol. IT-30, pp. 485–497, May 1984.
- [5] A. Gersho and R. M. Gray, *Vector Quantization and Signal Compression*. Boston, MA: Kluwer, 1992, ch. 2.
- [6] H. Gish and J. N. Pierce, "Asymptotically efficient quantizing," *IEEE Trans. Inform. Theory*, vol. IT-14, pp. 676–683, Sept. 1968.
- [7] N. J. Jayant and G. Noll, *Digital Coding of Waveforms*. Englewood Cliffs, NJ: Prentice-Hall, 1984.
- [8] E. LePennec, "Modélisation d'images par ondelettes, Dec. 1997," Mémoire de DEA, CMAP, Ecole Polytechnique, Paris, France.
- [9] A. S. Lewis and G. Knowles, "Image compression using the 2-D wavelet transform," *IEEE Trans. Image Processing*, vol. 1, pp. 244–250, Apr. 1992.

- [10] S. Mallat, "A theory for multiresolution signal decomposition: The wavelet representation," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 11, pp. 674–693, July 1989.
- [11] ———, *A Wavelet Tour of Signal Processing*. Boston, MA: Academic, 1998.
- [12] A. Said and W. A. Pearlman, "A new, fast, and efficient image codec based on set partitioning in hierarchical trees," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 6, pp. 243–250, June 1996.
- [13] J. M. Shapiro, "Embedded image coding using zerotrees of wavelet coefficients," *IEEE Trans. Signal Processing*, vol. 41, pp. 3445–3462, Dec. 1993.
- [14] G. Sullivan, "Efficient scalar quantization of exponential and laplacian random variables," *IEEE Trans. Inform. Theory*, vol. 42, pp. 1365–1374, Sept. 1996.
- [15] G. K. Wallace, "The JPEG still picture compression standard," *Commun. ACM*, vol. 34, no. 4, pp. 30–44, Apr. 1991.
- [16] I. Witten, R. Neal, and J. Cleary, "Arithmetic coding for data compression," *Commun. ACM*, vol. 30, no. 6, pp. 519–540, 1987.
- [17] Z. X. Xiong, O. Guleryuz, and M. T. Orchard, "Embedded image coding based on DCT," in *Proc. VCIP Euro. Image Processing Conf.*, 1997.



Stéphane Mallat (M'91) was born in Paris, France. He graduated from Ecole Polytechnique, Palaiseau, France, in 1984 and from Ecole Nationale Supérieure des Télécommunications, Paris, in 1985. He received the Ph.D. degree in electrical engineering from the University of Pennsylvania, Philadelphia, in 1988.

In 1988, he joined the Computer Science Department, the Courant Institute of Mathematical Sciences, New York University, New York, NY, and became Associate Professor in 1993. In the fall 1994, he was a Visiting Professor with the Department of Electrical Engineering, Massachusetts Institute of Technology, Cambridge, and, in the spring 1994, with the Department of Applied Mathematics, University of Tel Aviv, Tel Aviv, Israel. Since 1995, he has been a Professor with the Department of Applied Mathematics, Ecole Polytechnique. His research interest include computer vision, signal processing, and diverse applications of wavelet transforms. He is the author of the book *A Wavelet Tour of Signal Processing* (Boston, MA: Academic, 1998).

Dr. Mallat received the 1990 IEEE Signal Processing Society's Paper Award, the 1993 Alfred Sloan Fellowship in Mathematics, the 1997 Outstanding Achievement Award from the SPIE Optical Engineering Society, and the 1997 Blaise Pascal Prize in applied mathematics from the French Academy of Sciences.



Frédéric Falzon received the diplôme d'ingénieur from Ecole Supérieure, Paris, France, in sciences informatiques and the Diplôme d'Etudes Approfondies in Automatics and Signal Processing from the University of Nice-Sophia Antipolis, France, in 1990. He worked on Vision Systems and Image Processing at INRIA, Sophia Antipolis, and received the Ph.D degree from the University of Nice-Sophia Antipolis in 1994.

From 1995 to 1996, he was an expert engineer at INRIA, Sophia Antipolis, and studied problems such as deconvolution and low bit-rate compression of high-resolution satellite images. He is currently a research engineer at Alcatel Alsthom Recherche, Paris, and his area of research encompasses image segmentation and classification, texture analysis, and, more generally, image processing.