

DCT-DOMAIN IMAGE WATERMARKING AND GENERALIZED GAUSSIAN MODELS

J. R. Hernández, F. Pérez-González and M. Amado

Dept. Tecnologías de las Comunicaciones, ETSI Telecom., Universidad de Vigo, 36200 Vigo, Spain
email: jhernan@tsc.uvigo.es, fperez@tsc.uvigo.es

ABSTRACT

A novel DCT-domain watermark extraction procedure for still images that does not require the original image is presented. This method is based on a generalized Gaussian model, which includes as a special case the cross-correlation-based watermark detector, used so far in the literature. The optimal maximum likelihood (ML) detector is given, which allows to analytically assess the performance of watermarking methods in the DCT domain within a statistical framework. These original theoretical results are validated with experiments that show a considerable improvement over the existing watermark decoders. The perceptual model used in the tests is also described.

1. INTRODUCTION

In recent years we have witnessed a striking proliferation of techniques for representation, storage and distribution of digital multimedia information. Unfortunately, these developments have also opened the gate to unauthorized copying, distribution and manipulation of data, mostly images. Specialized and costly hardware may alleviate the problem of images duplication, at the price of a dramatic reduction in marketing possibilities –this is the cryptographic approach taken by pay TV channels, not foreseeable for scenarios such as Internet–. Watermarking techniques can at least ensure that ownership information is invisibly embedded into the image, thus preventing or deterring users from illegal uses.

Although many watermarking methods have sprouted over the few past years, even with commercial products available, the results up to date are quite discouraging, since there are freely available programs (e.g., unZign, Stirmark) that have succeeded in wiping the watermark away with little impact on the quality of the resulting image. Parallel to this, the lack of theoretical analyses in most of the available literature makes it difficult to know the actual limits in the performance of the various methods and to provide well-founded solutions which are the only way to eventually turn digital copyright protection into a mature discipline. In this

paper we make a contribution in this direction by showing how watermarking in the DCT domain (the most commonly used) can be dramatically improved by carefully modeling the problem and designing the proper watermark detector. We will assume throughout the paper that the original image is not known. While knowledge of the original image greatly simplifies the extraction procedure [1] it also narrows the range of possible applications.

2. DEFINITIONS

Let $x[\mathbf{n}]$ be a two-dimensional sequence representing the luminance of the original image, where $\mathbf{n} = (n_1, n_2)$. For the sake of readability, we will use in the sequel this vector notation to represent two-dimensional discrete indexes. Let $X[\mathbf{k}]$ be the result of applying a DCT transform to $x[\mathbf{n}]$ in a 8×8 pixels block basis. For copyright protection purposes, a watermark $W[\mathbf{k}]$ carrying some hidden information (owner and image identification number, transaction date, etc.) is added to the original image in the DCT domain, obtaining as a result the watermarked version $Y[\mathbf{k}] \triangleq X[\mathbf{k}] + W[\mathbf{k}]$.

In the watermarking technique we analyze in this paper, the watermark $W[\mathbf{k}]$ is generated in the DCT domain employing a 2-dimensional multipulse amplitude modulation scheme [2, 3]. In other words, $W[\mathbf{k}]$ can be expressed as the sum of N orthogonal pulses $\{P_i[\mathbf{k}]\}_{i=1}^N$

$$W[\mathbf{k}] = \sum_{i=1}^N b_i P_i[\mathbf{k}], \quad (1)$$

where the coefficients $\mathbf{b} \triangleq (b_1, \dots, b_N)$ are used to encode the hidden message. The modulation pulses $\{P_i[\mathbf{k}]\}_{i=1}^N$ are generated as a function of a secret key K , only known by the copyright owner. They are expressed as

$$P_i[\mathbf{k}] = \begin{cases} \alpha[\mathbf{k}]s[\mathbf{k}], & \mathbf{k} \in \mathcal{S}_i \\ 0, & \text{otherwise,} \end{cases} \quad (2)$$

where $s[\mathbf{k}]$ is a key-dependent pseudorandom sequence such that $s[\mathbf{k}] \in \{-1, 1\}$, $\forall \mathbf{k}$, and the sets of indexes $\mathcal{T} \triangleq \{\mathcal{S}_i\}_{i=1}^N$ are also key-dependent and determine the spatial shape of the pulses. The sequence $\alpha[\mathbf{k}]$ is called the *perceptual mask*

and indicates the maximum allowable magnitude of the alteration that the coefficient $X[\mathbf{k}]$ may suffer without achieving noticeable distortions. The sets $\{\mathcal{S}_i\}_{i=1}^N$ are assumed to be non-overlapping, i.e. $\mathcal{S}_i \cap \mathcal{S}_j = \emptyset, \forall i \neq j$, and sparsely spread over the whole image in a pseudorandom fashion to provide security and robustness against cropping [2, 3].

Given a watermarked image $Y[\mathbf{k}]$ and the secret key K , the *watermark extraction* procedure obtains an estimate of the secret message \mathbf{b} . Although intimately related to our problem, the so-called *watermark detection test* which answers the ownership question, is out of the scope of this paper. In addition, we will assume that no attacks aimed at desynchronizing the watermark are performed. However, both the synchronization and watermark detection problems can be tackled within the statistical framework presented in Section 2 and analyzed in Section 3. The perceptual model used in our particular watermarking scheme is given in Section 4. Section 5 is devoted to experimental results, while Section 6 presents our conclusions and future lines of research.

3. DETECTOR STRUCTURE

Detector structures usually proposed for hidden information decoding in DCT-domain spread spectrum data hiding techniques are based on the crosscorrelation between the watermarked image $Y[\mathbf{k}]$ and the pseudorandom sequence $s[\mathbf{k}]$. This scheme would be appropriate if noise –in watermarking, the original image– followed a Gaussian distribution. However, the Gaussian assumption is inaccurate for DCT coefficients of common images. Some authors have proposed the generalized Gaussian probability density function (pdf)

$$f_x(x) = A e^{-|\beta x|^c}. \quad (3)$$

as an alternative leading to improved statistical models [4]. Note that the Gaussian and the Laplacian pdf's are just special cases of this expression, given by $c = 2$ and $c = 1$, respectively. Previous works in this field show that DCT coefficients at low frequencies are reasonably well-modeled by a generalized Gaussian distribution with $c = 1/2$. Coefficients at high frequencies are better approximated by a Gaussian distribution and sometimes by a Laplacian distribution.

The parameters A and β in Eq. (3) can be expressed as

$$\beta = \frac{1}{\sigma} \left(\frac{\Gamma(3/c)}{\Gamma(1/c)} \right)^{1/2}; \quad A = \frac{\beta c}{2\Gamma(1/c)},$$

where σ is the standard deviation. Hence, the pdf is completely specified by c and σ . Let us define the sequence

$$C_{i,j}[k_1, k_2] \triangleq X[8k_1 + i, 8k_2 + j], \quad i, j \in \{0, \dots, 7\},$$

which results if we take the (i, j) -th DCT coefficient of every block. We will model each of these 64 sequences as the output of an two-dimensional i.i.d. random process whose marginal distribution follows Eq. (3), with parameters $c(i, j)$ and $\sigma(i, j)$. Let us also define the sequences $c[\mathbf{k}]$ as

$$c[\mathbf{k}] \triangleq c(k_1 \bmod 8, k_2 \bmod 8)$$

and $\sigma[\mathbf{k}]$ in a similar fashion. Thus, these two sequences indicate the parameters c and σ associated with each sample $X[\mathbf{k}]$. Let us assume that M possible different messages can be encoded with the vector $\mathbf{b} = (b_1, \dots, b_N)$ and let $\mathbf{b}_l, l \in \{1, \dots, M\}$ indicate the codeword associated to one of those messages. Also, let $W_l[\mathbf{k}], l \in \{1, \dots, M\}$ be the watermark obtained from $\mathbf{b}_l = (b_{l,1}, \dots, b_{l,N})$ using Eq. (1). Then, assuming the i.i.d. generalized Gaussian model for $X[\mathbf{k}]$, it can be easily shown that the optimum detector in the ML sense is the one that chooses the index $l \in \{1, \dots, M\}$ verifying

$$\sum_{\mathbf{k}} \frac{|Y[\mathbf{k}] - W_m[\mathbf{k}]|^{c[\mathbf{k}]} - |Y[\mathbf{k}] - W_l[\mathbf{k}]|^{c[\mathbf{k}]}}{\sigma[\mathbf{k}]^{c[\mathbf{k}]}} > 0 \\ \forall m \neq l.$$

Assuming that $b_{l,i} \in \{-1, 1\}, \forall l \in \{1, \dots, M\}, i \in \{1, \dots, N\}$, this optimization problem is equivalent to finding the codeword \mathbf{b}_l which maximizes the expression $\sum_{i=1}^N b_{l,i} r_i$, where the coefficients r_i are sufficient statistics for the detection problem and are defined as

$$r_i \triangleq \sum_{\mathbf{k} \in \mathcal{S}_i} \frac{|Y[\mathbf{k}] + \alpha[\mathbf{k}]s[\mathbf{k}]|^{c[\mathbf{k}]} - |Y[\mathbf{k}] - \alpha[\mathbf{k}]s[\mathbf{k}]|^{c[\mathbf{k}]}}{\sigma[\mathbf{k}]^{c[\mathbf{k}]}}.$$

When a binary antipodal constellation is used to encode $M = 2^N$ possible messages, the ML detector structure is equivalent to a bit-by-bit hard decisor, so the outputs of the decoder are

$$\hat{b}_i = \text{sgn}(r_i), \quad i \in \{1, \dots, N\}.$$

4. PERFORMANCE ANALYSIS

Now let us analyze the performance of the watermark decoding process in terms of the *probability of bit error* P_b . Obviously, performance results strongly depend on image characteristics, so we will obtain P_b conditioned to a given original image $X[\mathbf{k}]$ or, in other words, the probability of getting a bit error when a secret key is taken at random and is applied in both the watermarking and decoding processes. In this context, $X[\mathbf{k}]$ will be regarded as a deterministic signal while the sequence $s[\mathbf{k}]$ and the sets $\mathcal{T} = \{\mathcal{S}_i\}_{i=1}^N$ will be modeled statistically.

If the pseudorandom sequence $s[\mathbf{n}]$ is modeled as an i.i.d. two-dimensional random process with marginal pdf $f_s(s)$, then, each sufficient statistic r_i is the sum of $|\mathcal{S}_i|$ statistically independent contributions ($|\mathcal{S}_i|$ is the cardinality of the set $\{\mathbf{k}, P_i[\mathbf{k}] \neq 0\}$). Hence, by central limit theorem arguments, $\mathbf{r} \triangleq (r_1, \dots, r_N)$ can be accurately approximated as the output of a vector Gaussian channel. Therefore, the probability of error conditioned to $X[\mathbf{k}]$ can be expressed as a function of the first and second order moments of r_1, \dots, r_N . Let us define the two-dimensional sequence

$$r[\mathbf{k}] \triangleq |Y[\mathbf{k}] + \alpha[\mathbf{k}]s[\mathbf{k}]|^{c[\mathbf{k}]} - |Y[\mathbf{k}] - \alpha[\mathbf{k}]s[\mathbf{k}]|^{c[\mathbf{k}]},$$

extracted from Eq. (5). Since the elements of $r[\mathbf{k}]$ are statistically independent, the mean and variance of r_i conditioned to a certain tiling $\mathcal{T} = \{\mathcal{S}_i\}_{i=1}^N$ are

$$E[r_i | \mathcal{T}] = \sum_{\mathbf{k} \in \mathcal{S}_i} \frac{E[r[\mathbf{k}]]}{\sigma[\mathbf{k}]^{c[\mathbf{k}]}}; \quad \text{Var}(r_i | \mathcal{T}) = \sum_{\mathbf{k} \in \mathcal{S}_i} \frac{\text{Var}(r[\mathbf{k}])}{\sigma[\mathbf{k}]^{2c[\mathbf{k}]}}$$

Now, since the tiling \mathcal{T} is also key-dependent, the overall first and second order moments of r_i can be obtained using the expressions

$$E[r_i] = E_{\mathcal{T}}[E[r_i | \mathcal{T}]] \quad (8)$$

$$\text{Var}(r_i) = E_{\mathcal{T}}[\text{Var}(r_i | \mathcal{T})] + \text{Var}_{\mathcal{T}}(E[r_i | \mathcal{T}]), \quad (9)$$

where the expectations are now evaluated over the set of possible tilings \mathcal{T} . If each index $\mathbf{k} \in \mathbb{N}^2$ belongs to \mathcal{S}_i with probability $1/N$ for all $i \in \{1, \dots, N\}$ and assignments of indices to sets are performed independently, i.e. $\Pr\{\mathbf{k} \in \mathcal{S}_i, \mathbf{m} \in \mathcal{S}_j\} = \Pr\{\mathbf{k} \in \mathcal{S}_i\}\Pr\{\mathbf{m} \in \mathcal{S}_j\}, \forall \mathbf{k} \neq \mathbf{m}, i, j \in \{1, \dots, N\}$, then after some algebraic manipulations it can be proven that

$$E[r_i] = \frac{1}{N} \sum_{\mathbf{k}} \frac{E[r[\mathbf{k}]]}{\sigma[\mathbf{k}]^{c[\mathbf{k}]}}. \quad (10)$$

$$\text{Var}(r_i) = \frac{1}{N} \sum_{\mathbf{k}} \frac{\text{Var}(r[\mathbf{k}])}{\sigma[\mathbf{k}]^{2c[\mathbf{k}]}} + \frac{N-1}{N^2} \sum_{\mathbf{k}} \frac{E^2[r[\mathbf{k}]]}{\sigma[\mathbf{k}]^{2c[\mathbf{k}]}}. \quad (11)$$

Assume that $b_i = 1$. Then, $Y[\mathbf{k}] = X[\mathbf{k}] + \alpha[\mathbf{k}]s[\mathbf{k}], \forall \mathbf{k} \in \mathcal{S}_i$, and, as a consequence,

$$r[\mathbf{k}] = |X[\mathbf{k}] + 2\alpha[\mathbf{k}]s[\mathbf{k}]|^{c[\mathbf{k}]} - |X[\mathbf{k}]|^{c[\mathbf{k}]}.$$

If $s[\mathbf{k}]$ follows a discrete uniform two-level distribution, it can be easily shown that the mean and variance of $r[\mathbf{k}]$ are

$$\begin{aligned} E[r[\mathbf{k}]] &= \frac{1}{2} \left[\left(|X[\mathbf{k}]| + 2\alpha[\mathbf{k}] \right)^{c[\mathbf{k}]} \right. \\ &\quad \left. + \left| |X[\mathbf{k}]| - 2\alpha[\mathbf{k}] \right|^{c[\mathbf{k}]} \right] - |X[\mathbf{k}]|^{c[\mathbf{k}]} \\ \text{Var}(r[\mathbf{k}]) &= \frac{1}{4} \left[\left(|X[\mathbf{k}]| + 2\alpha[\mathbf{k}] \right)^{c[\mathbf{k}]} \right. \\ &\quad \left. - \left| |X[\mathbf{k}]| - 2\alpha[\mathbf{k}] \right|^{c[\mathbf{k}]} \right]^2. \end{aligned}$$

These expressions can be applied in equations (10) and (11) to compute the moments of r_i . When $b_i = -1$, it can be verified that $\text{Var}(r_i)$ is given by Eq. (11) and $E[r_i]$ is negative and its absolute value is given by Eq. (10). When a binary antipodal constellation with $M = 2^N$ is used to encode the hidden message, the probability of bit error P_b of the ML decoder (a bit-by-bit hard decisor) is $P_b = Q(\text{SNR})$, where $Q(x) \triangleq \frac{1}{\sqrt{2\pi}} \int_x^\infty e^{-t^2/2} dt$ and the signal to noise ratio SNR is defined as

$$\text{SNR} \triangleq \frac{E[r_i]}{\sqrt{\text{Var}(r_i)}}. \quad (13)$$

5. PERCEPTUAL MODEL

In Eqs. (1,2) the watermark $W[\mathbf{k}]$ depends on a perceptual mask $\alpha[\mathbf{k}]$ that multiplies the pseudorandom sequence $s[\mathbf{k}]$. This perceptual mask determines the maximum amplitude distortion that each coefficient of the original image may suffer while satisfying the invisibility constraint. A good psychovisual model in the DCT-domain (with 8x8 blocks) is capital to render the sequence $\alpha[\mathbf{k}]$. For our work we have followed the model proposed in [5, 6] that has been also applied to derive adaptive quantization matrices for the JPEG algorithm [7]. This model has been here simplified by disregarding the so-called *contrast-masking effect* for which the perceptual mask at a certain coefficient depends on the amplitude of the coefficient itself. Consideration of this effect constitutes a future line of research. On the other hand, the *background intensity effect*, for which the mask depends on the magnitude of the DC coefficient (i.e., the background), has been taken into account.

The so-called *visibility threshold* $T(i, j), i \in \{0, \dots, 7\}, j \in \{0, \dots, 7\}$, determines the maximum allowable magnitude of an invisible alteration of the (i, j) -th DCT coefficient and can be approximated in logarithmic units by the following quadratic function with parameter K

$$\begin{aligned} \log T(i, j) &= \log \left(\frac{T_{\min}(f_{i,0}^2 + f_{0,j}^2)^2}{(f_{i,0}^2 + f_{0,j}^2)^2 - 4(1-r)f_{i,0}^2 f_{0,j}^2} \right) \\ &\quad + K \left(\log \sqrt{f_{i,0}^2 + f_{0,j}^2} - \log f_{\min} \right)^2, \end{aligned}$$

where $f_{i,0}$ and $f_{0,j}$ are respectively the vertical and horizontal spatial frequencies (in cycles/degree) of the DCT-basis functions, T_{\min} is the minimum value of $T(i, j)$, associated to the spatial frequency f_{\min} , and r is taken as 0.7 following [5]. The threshold $T(i, j)$ can be corrected for each block by considering the DC coefficient $X_{0,0}$ and the average luminance of the screen $\bar{X}_{0,0}$ (1024 for an 8-bit image) in the following way

$$T'(i, j) = T(i, j) \left(\frac{X_{0,0}}{\bar{X}_{0,0}} \right)^{\alpha_T}.$$

Note that the actual dependence of $X_{0,0}$ on the block indices has been dropped in the notation for conciseness. Following [5], the parameters used in our scheme have been set to $a_T = 0.649$, $f_{min} = 3.68$ cycles/degree, $T_{min} = 1.1548$ and $K = 1.728$. Once the corrected threshold value $T'_{i,j}$ has been obtained, the perceptual mask is calculated as

$$\alpha[k_1, k_2] = 4 \left(1 + \frac{\delta(l_1)}{\sqrt{2} + 1} \right) \left(1 + \frac{\delta(l_2)}{\sqrt{2} + 1} \right) \gamma \cdot T'(l_1, l_2) \quad (15)$$

where $l_1 = k_1 \bmod 8$, $l_2 = k_2 \bmod 8$ and $\gamma < 1$ is a scaling factor that allows to introduce a certain degree of conservativeness in the watermark due to those effects that have been overlooked (e.g., spatial masking in the frequency domain [8]). The remaining factors in (15) allow to express the corrected threshold in terms of DCT coefficients instead of luminances.

6. EXPERIMENTAL RESULTS

In order to validate the theoretical analysis presented in previous sections, we have watermarked the well-known image ‘Lena’ (256 x 256 pixels), shown in Fig. 1, following the method described in sections 1 and 4, modifying only the mid-frequency coefficients (low frequency coefficients have very low capacity, i.e., slight modifications become quite visible; high frequency coefficients can be easily erased by compression algorithms). We run 100 experiments with different keys for a fixed number of pixels per information bit and computed the resulting bit error rate (BER). Figure 2 shows the resulting watermark for one of such experiments. In Figure 3 both empirical and theoretical results for different values of the generalized Gaussian parameter c are plotted. Note that the parameter γ in Eq. (15) has been set to $1/5$ —so the watermark is well below the visibility level—in order to produce statistically significant results. The actual performance is substantially better, but the qualitative conclusions remain the same. As can be inferred from Figure 3 and also from Figure 4, where the SNR in Eq. (13) is plotted for different values of c , good results are obtained in the range $1/2 \leq c \leq 1$. Interestingly enough, the performance for $c = 2$, corresponding to the cross-correlation-based detector used so far in the literature [9], suffers a severe deterioration, corresponding to a drop of more than 6 dB in the SNR (cf. Fig. 4).

Although not directly discussed here, our analysis can be somewhat straightforwardly extended to the case of JPEG compression. Figure 5 shows the theoretical BER obtained when image ‘Lena’ is watermarked (with a 100 bits hidden message) and later compressed with JPEG to a percentage of its original quality. As it can be seen, in this case, the Laplacian detector ($c = 1$) performs slightly better than that with $c = 1/2$. The Gaussian (cross-correlation) detector, not shown in the Figure, produces a much higher BER.



Figure 1: Original image used in the tests.

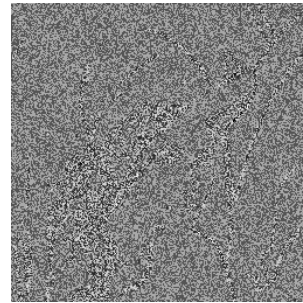


Figure 2: One of the watermarks used in the tests.

The curves labeled as ‘Optimum’ correspond to a detector specifically designed for the JPEG compression attack.

7. CONCLUSIONS AND FURTHER WORK

A new structure based on the use of generalized Gaussian models has been proposed for the ML detection of DCT-domain watermarks embedded in still images. By considering these models, we have been able to dramatically improve the performance of the cross-correlation-based detector, that has been used up to date. In any case, we also have presented a theoretical analysis that allows to assess the performance of DCT-based methods, measured in terms of the bit error rate for a given hidden message length and a given image. The Gaussian detector is simply a particular case of the generalized model, so the analytical results given here are directly applicable. One immediate extension of our analysis is the consideration of channel codes which have been already shown to considerably improve on spatial-domain watermarking methods [10]. Another future line of research consists in fitting a generalized Gaussian model (from a discrete set of parameters c and σ) to the DCT coefficients histogram so as to decide upon and apply the optimal detector structure.

8. REFERENCES

- [1] I. J. Cox, J. Kilian, F. T. Leighton, and T. Shamoan, “Secure spread spectrum watermarking for multimedia,” *IEEE*

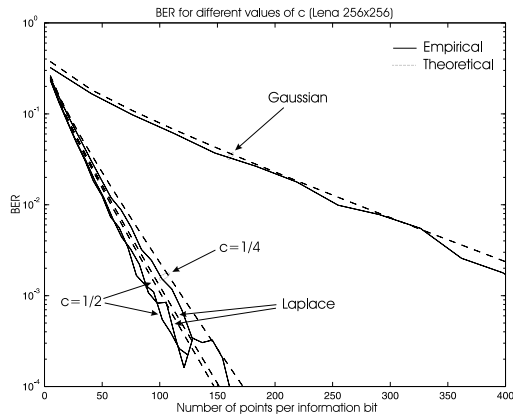


Figure 3: BER as a function of the pulse size for Lena (256 × 256).

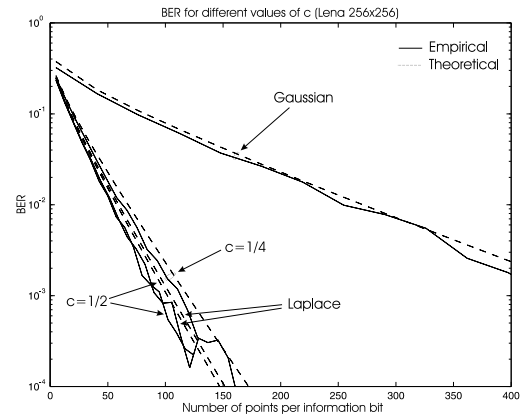


Figure 5: BER as a function of JPEG final quality for Lena.

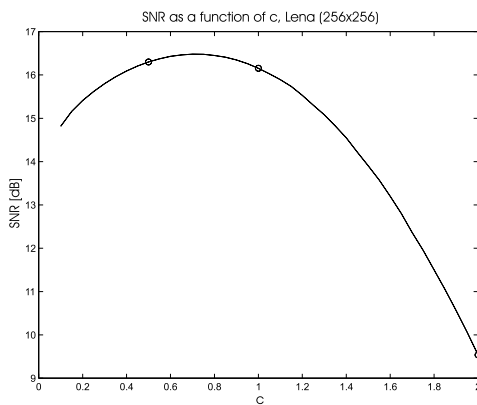


Figure 4: SNR as a function of c for Lena (256 × 256).

Transactions on Image Processing, vol. 6, pp. 1673–1687, December 1997.

- [2] J. R. Hernández, F. Pérez-González, J. M. Rodríguez, and G. Nieto, “Performance analysis of a 2d-multipulse amplitude modulation scheme for data hiding and watermarking of still images,” *IEEE J. Select. Areas Commun.*, vol. 16, pp. 510–524, May 1998.
- [3] J. R. Hernández, F. Pérez-González, and J. M. Rodríguez, “The impact of channel coding on the performance of spatial watermarking for copyright protection,” in *Proc. ICASSP’98*, vol. 5, (Seattle, Washington (USA)), pp. 2973–2976, May 1998.
- [4] R. J. Clarke, *Transform Coding of Images*. Academic Press, 1985.
- [5] A. J. Ahumada and H. A. Peterson, “Luminance-model-based DCT quantization for color image compression,” in *Human Vision, Visual Processing, and Digital Display III (Proc. of the SPIE)* (B. E. Rogowitz, ed.), 1992.
- [6] J. A. Solomon, A. B. Watson, and A. J. Ahumada, “Visibility of DCT basis functions: Effects of contrast masking,” in *Proc. Data Compression Conference*, (Snowbird, Utah (USA)), pp. 361–370, IEEE Computer Society Press, 1994.

- [7] A. B. Watson, “Visual optimization of DCT quantization matrices for individual images,” in *Proc., AIAA Computing in Aerospace 9*, (San Diego, California (USA)), pp. 286–291, American Institute of Aeronautics and Astronautics, 1993.
- [8] A. N. Netravali and B. G. Haskell, *Digital Pictures. Representation, Compression and Standards*. New York: Plenum Press, 1995.
- [9] M. Barni, F. Bartolini, V. Cappellini, and A. Piva, “A DCT-domain system for robust image watermarking,” *Signal Processing*, vol. 66, pp. 357–372, May 1998.
- [10] J. R. Hernández, J. M. Rodríguez, and F. Pérez-González, “Improving the performance of spatial watermarking of images using channel coding.” Accepted for publication in *Signal Processing*, Elsevier.