

# EFFECT OF JPEG COMPRESSION ON IMAGE WATERMARK DETECTION

*Minghui Xia, Bede Liu*

Information Sciences and Systems, Department of Electrical Engineering  
Princeton University, Princeton, NJ 08544  
Contact: {mxia,liu}@ee.princeton.edu

## ABSTRACT

Compression on a watermarked image can significantly affect the detection of the embedded watermark. To show that a particular watermarking scheme is robust against compression, simulation is often relied upon. In this paper, we investigate this problem analytically. We first characterize the noise introduced by JPEG compression. We then propose a maximum a posteriori (MAP) detector, analyze its performance, and demonstrate that it has superior performance over the correlation detector.

## 1. INTRODUCTION

The detection of the presence or absence of a watermark embedded in an image is often affected if the watermarked image has undergone compression. Compression can also be considered as an attack on watermarked images [1]. Partly because it is difficult to analyze the effect of compression on the detection of watermarked image, past investigations of this problem have relied heavily on simulation. In this paper, we analyze the effect of JPEG compression on a well known watermark scheme. Noise is introduced by JPEG compression when the block-wise DCT coefficients are quantized. The statistics of this quantization noise generally depend on the statistics of the input source, and cannot be modeled merely as additive white Gaussian noise. Modeling of the quantization noise is presented in Section 3 and used in a maximum a posteriori (MAP) detector proposed in Section 4. The performance of the MAP detector is analyzed and, as expected, is shown to have better performance than the commonly used correlation detector. Experimental results are presented in Section 6.

We assume the image is of size  $N \times N$  pixels. (The case of non-square arrays are handled similarly.) We need to deal with two kinds of DCT's. The first is that of the entire image; it is also a 2D array of size  $N \times N$ . The second is a block-wise DCT or BDCT, which is obtained by first dividing the image into blocks of  $M \times M$  pixels ( $M=8$  in JPEG), and then taking the DCT of each block. The result is also a 2D array

of size  $N \times N$ . To study the statistical properties of images and its transforms, it is convenient sometimes to convert the 2D array into a 1D vector. This can be done by taking the rows in successive order. Another way is to divide the 2D array into blocks, arrange all blocks in a 1D array, and convert each block to a vector. When the vector  $\xi$  is the realization of a random variable  $\xi$ , we sometimes use  $\xi$  to represent either case for simplicity. When not all entries in a vector  $\xi$  are involved in a computation, it is convenient to make use of a select operation  $S$  on the vector to pick out only those element of interest:  $S \xi$ . For convenience, we may at times use the same symbol to denote the different representations of the same object, if the meaning is clear from the context. Thus,  $y$  may denote the  $N \times N$  array of pixels, or a 1D vector of all the pixels. Similarly,  $\xi$  may denote a vector or a smaller vector  $S \xi$ .

## 2. SPREAD SPECTRUM EMBEDDING IN FREQUENCY DOMAIN

In this section, we investigate a specific watermark scheme proposed in [2]. Similar analysis can be carried out for other schemes. In this scheme, the DCT of the entire host image is taken and the largest  $K$  coefficients, denoted by  $x$ , are modified to embed a watermark  $w$ . An often used watermark is a spreading sequence of normalized Gaussian distributed random numbers, generated from a key. It is embedded according to:

$$y = (1 + \alpha \cdot w) x. \quad (1)$$

That is, the original  $K$  largest coefficients  $x$  is replaced by the  $K$  modified coefficients  $y$ . The inverse DCT of the entire  $N \times N$  array is the marked image. The masking factor  $\alpha$  is adjustable, to make the watermark invisible. This method inserts a single watermark in the entire image, similar to using spread spectrum for transmitting one bit information. The detector decides whether or not a test image contains the specific watermark. This is done by subtracting the original image from the test image, taking the DCT of the difference, and selecting the prescribed set of  $K$  coefficients to form the test watermark, denoted by  $\hat{w}$ . It is

---

This research is supported by an NJ State R&D excellence grant.

compared with the original watermark  $w$  using correlation:

$$\text{sim}(w, \hat{w}) = \frac{\langle \hat{w} \cdot w \rangle}{\sqrt{\langle \hat{w} \cdot \hat{w} \rangle}}, \quad (2)$$

where  $\langle \cdot \rangle$  denotes inner product.

### 3. NOISE DUE TO JPEG COMPRESSION

The block diagram of Figure 1 depicts watermark embedding, followed by JPEG compression and watermark detection. A watermark  $w$  is embedded in the DCT of the image to obtain  $y$ , which is then transformed into image domain by inverse DCT. JPEG compression takes the 8x8 block DCT of the marked image  $z$ , quantize the coefficients  $f$ , and produce the distorted coefficients  $\tilde{f}$ . An 8x8 block inverse DCT then gives the compressed and watermarked image,  $\tilde{z}$ . To detect the watermark, the DCT of  $\tilde{z}$  is taken and subtracted from it the DCT of the original image  $x$ .

As the blocks labeled with IDCT and 8x8 BDCT are linear operations,  $f$  can be related to  $y$  via

$$f = H y, \quad (3)$$

where  $H$  is the linear transformation matrix. Similarly,

$$\tilde{y} = H^{-1} \tilde{f} \quad (4)$$

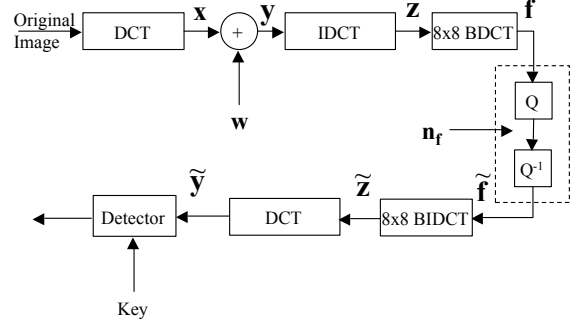
The major source of error is due to quantization of the 8x8 block DCT coefficients in JPEG compression. The error introduced can be approximated as an additive noise  $n_f$ , whose statistics depend in general on that of the input source and the quantization step size  $\Delta$ . The value of  $\Delta$  is position dependent in a 8x8 block. DCT coefficients of images without watermark can be modeled by Generalized Gaussian or Laplacian distribution [3]. When an image independent watermark is embedded in the image, the marked image can be viewed as dithered signal added before the quantizer. Dithered quantizer was analyzed in [4]. However, the watermark inserted according to Eq. 1, is dependent of the original image, making the analysis of [4] invalid.

#### 3.1. Watermark absent

When there is no watermark, the first and second order statistics of the quantization noise  $n_f$  are given by

$$\begin{aligned} E\{n_{f_0}\} &= -j \sum_{\substack{n=-\infty \\ n \neq 0}}^{\infty} \frac{(-1)^n}{u} \Phi_f(u) \big|_{u=\frac{2n\pi}{\Delta}}, \quad (5) \\ E\{n_{f_0}^2\} &= \frac{\Delta^2}{12} + \sum_{\substack{n=-\infty \\ n \neq 0}}^{\infty} \frac{2(-1)^n}{u^2} \Phi_f(u) \big|_{u=\frac{2n\pi}{\Delta}}, \end{aligned}$$

where  $\Phi_f(u)$  is the characteristic function of the block DCT coefficients  $f$ . Since Laplacian or generalized Gaussian distribution will be used to calculate the characteristic function, the quantization error in this case has zero mean.



**Fig. 1.** Watermarked image going through JPEG compression-decompression

#### 3.2. Watermark present

When watermark is present, we keep the notation  $f$  to represent the block DCT coefficients of the original image. The quantity to be quantized is now

$$\begin{aligned} f_w &= H y = H (x (1 + \alpha w)) \\ &= H x + H (\alpha w x) = H x + w_H, \end{aligned} \quad (6)$$

where  $w_H = H (\alpha w x)$ .

Note that  $w$  is Gaussian with zero mean. So for a given image,  $f_w$  is Gaussian with mean  $H x$  and certain variance that can be computed from  $w_H$ . We do so by modeling the blocks individually. As in the no-watermark case, the first and second moments of the quantization error can be calculated from

$$E\{n_{f_1}\} = -j \sum_{\substack{n=-\infty \\ n \neq 0}}^{\infty} \frac{(-1)^n}{u} \Phi_{f_w}(u) \big|_{u=\frac{2n\pi}{\Delta}}, \quad (7)$$

$$E\{n_{f_1}^2\} = \frac{\Delta^2}{12} + \sum_{\substack{n=-\infty \\ n \neq 0}}^{\infty} \frac{2(-1)^n}{u^2} \Phi_{f_w}(u) \big|_{u=\frac{2n\pi}{\Delta}},$$

$$\text{Var}\{n_{f_1}\} = E\{n_{f_1}^2\} - E\{n_{f_1}\}^2,$$

where  $\Phi_{f_w}(u)$  is the characteristic function of  $f_w$ , which can be obtained by computing the mean  $H x$  and the variance of  $w_H$ .

Note that  $n_{f_1}$  and  $w_H$  are dependent. It will be seen later in Section 5 that this dependency will affect the performance of the MAP detector. The correlation between  $n_{f_1}$  and  $w_H$  can be obtained by

$$E(n_{f_1} \cdot w_H) = \frac{d^2}{dudv} \Phi_{\tilde{f}_w f_w} \big|_{u=0, v=0} - E\{f_w^2\} - f E\{n_{f_1}\}, \quad (8)$$

where  $f$  is the block DCT coefficient of the original image,

$$\Phi_{f_w f_w} = \sum_{n=-\infty}^{\infty} \Phi_{f_w f_w} \left( u + \frac{2n\pi}{\Delta} v \right) \frac{\sin \frac{\Delta}{2} \left( u + \frac{2n\pi}{\Delta} \right)}{\frac{\Delta}{2} \left( u + \frac{2n\pi}{\Delta} \right)}. \quad (9)$$

$\Phi_{f_w f_w}(u, v)$  can be calculated from the auto-correlation function of  $\mathbf{f}_w$ .

The quantization noise  $\mathbf{n}_f$  propagates to the input of the detector via linear transformation  $\mathbf{H}$ . Therefore, if the noise at the input of the detector  $\mathbf{n}_y$  is  $\tilde{\mathbf{y}} - \mathbf{y}$ , its covariance matrix  $\mathbf{C}_{n_y}$  can be expressed as

$$\mathbf{C}_{n_y} = \mathbf{H}^{-1} \mathbf{C}_{n_f} \mathbf{H}, \quad (10)$$

where  $\mathbf{C}_{n_f}$  is the covariance matrix for  $\mathbf{n}_f$ .

#### 4. MAXIMUM-A-POSTERIOR (MAP) DETECTOR

A MAP detector which uses statistical knowledge of the signal and noise would give superior performance over a correlation detector. The detector tests hypothesis  $H_0$  (watermark not present) against hypothesis  $H_1$  (watermark present):

$$\begin{cases} H_0 : \text{no watermark } (\hat{\mathbf{w}} = \mathbf{n}_0), \\ H_1 : \text{watermark present } (\hat{\mathbf{w}} = \mathbf{w} + \mathbf{n}_1). \end{cases} \quad (11)$$

Here  $\mathbf{n}_0$  is the noise at the detector when there is no watermark and  $\mathbf{n}_1$  is the noise when watermark is present. The detector uses the following log likelihood function

$$l(\hat{\mathbf{w}}) = \ln \frac{p_1(\hat{\mathbf{w}} | H_1)}{p_0(\hat{\mathbf{w}} | H_0)}, \quad (12)$$

where  $p_1$  and  $p_0$  are respectively the probability density functions of the test watermark  $(\hat{\mathbf{w}})$  under the hypotheses  $H_0$  and  $H_1$ . If the noises  $\mathbf{n}_0$  and  $\mathbf{n}_1$  are assumed to follow Gaussian distributions with covariance matrix  $\mathbf{C}_0$  and  $\mathbf{C}_1$  respectively, Eq. 12 reduces to

$$l(\hat{\mathbf{w}}) = \hat{\mathbf{w}}^T \mathbf{C}_0^{-1} \hat{\mathbf{w}} - (\hat{\mathbf{w}} - \mathbf{w})^T \mathbf{C}_1^{-1} (\hat{\mathbf{w}} - \mathbf{w}). \quad (13)$$

The matrices  $\mathbf{C}_0$  and  $\mathbf{C}_1$  are quite large. For images of size  $256 \times 256$ ,  $\mathbf{C}_0$  and  $\mathbf{C}_1$  are  $(256)^2 \times (256)^2$ , which would require 32GB storage for double-precision computation. To simplify subsequent calculations, we shall use the approximation that  $\mathbf{C}_0 = \mathbf{C}_1 = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_M^2)$ , where  $M$  is the number of rows/columns in  $\mathbf{C}_0$  and  $\mathbf{C}_1$ .

That  $\mathbf{C}_0$  and  $\mathbf{C}_1$  are diagonal can be justified by observing that the noises  $\mathbf{n}_0$  and  $\mathbf{n}_1$  were generated by quantization and passed through the DCT step to reach the input of the detector and that DCT has the effect of de-correlation. That  $\mathbf{C}_0 = \mathbf{C}_1$  can be justified because invisible watermark is much smaller than the host signal so it is not ex-

pected to change the variance of the quantization error significantly. With this approximation, Eq 13 reduces to

$$l(\hat{\mathbf{w}}) = \sum_{i=1}^K \frac{\hat{w}_i w_i}{\sigma_i^2} / \sum_{i=1}^K \frac{w_i^2}{\sigma_i^2}, \quad (14)$$

where  $K$  is the length of the inserted watermark, and the term  $w_i^2/\sigma_i^2$  is due to normalization.

Given a JPEG image with known compression quality factor (which determines the quantization matrix), the overall detector function can be computed as follows. First extract the watermark  $\hat{\mathbf{w}}$  in DCT domain by subtracting the original image. Then we calculate the covariance matrix for  $\mathbf{n}_f$  as in Figure 1, using Eqs. 5 and 7. The covariance matrices of the noise at the detector input can then be obtained by Eq. 10.  $\sigma_i^2$  in Eq. 14 are just the diagonal elements of the resulting matrix. In fact, the intermediate computation can be simplified since we only need the diagonal elements of the final covariance matrices.

#### 5. MAP DETECTOR PERFORMANCE ANALYSIS

When watermark is present, we have

$$\tilde{w}_i = w_i + n_i. \quad (15)$$

Eq. 14 becomes

$$l(\hat{\mathbf{w}}) = 1 + \sum_{i=1}^N \frac{n_i w_i}{\sigma_i^2} / \sum_{i=1}^N \frac{w_i^2}{\sigma_i^2}. \quad (16)$$

This is a sum of a large number of nearly independent variables, it may be approximated by a Gaussian variable,  $l(\hat{\mathbf{w}}) \sim N(\mu_{l_1}, \sigma_{l_1}^2)$ . Normalizing the variance of inserted watermark to 1,  $\mu_{l_1}$  and  $\sigma_{l_1}^2$  can be calculated by

$$\begin{aligned} \mu_{l_1} &= 1 + \sum_{i=1}^N \frac{E\{n_i w_i\}}{\sigma_i^2} / \sum_{i=1}^N \frac{1}{\sigma_i^2}, \\ \sigma_{l_1}^2 &= 1 / \sum_{i=1}^N \frac{1}{\sigma_i^2}, \end{aligned} \quad (17)$$

where  $E\{n_i w_i\}$  can be obtained by first computing the noise-watermark correlation at the quantizer output from Eq. 8, then calculating the detector input noise-watermark covariance matrix by using Eq. 10.

Similarly, Under hypothesis  $H_0$ ,

$$\begin{aligned} \mu_{l_0} &= 0, \\ \sigma_{l_0}^2 &= 1 / \sum_{i=1}^N \frac{1}{\sigma_i^2}. \end{aligned} \quad (18)$$

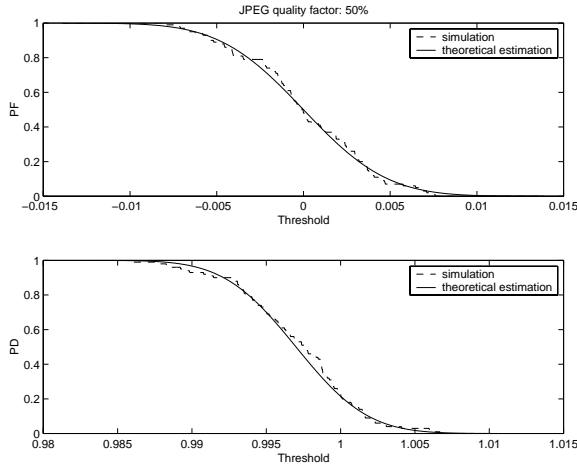
As a result,  $P_D$  and  $P_F$  can be expressed by

$$\begin{aligned} P_D &= Q\left(\frac{\eta - \mu_{l_1}}{\sigma_{l_1}}\right), \\ P_F &= Q\left(\frac{\eta}{\sigma_{l_0}}\right). \end{aligned} \quad (19)$$

where  $Q(x) \triangleq \frac{1}{2\pi} \int_x^\infty e^{-t^2/2} dt$ .

## 6. RESULTS

A gray level Lena image of size  $256 \times 256$  is used to test our theoretical analysis. The JPEG compression quality factor is set at 50%. Figure 2 shows the theoretical estimation of  $P_D$  and  $P_F$  of the MAP detector, as well as the simulation result. Note that the mean of the detector likelihood function  $l(\hat{w})$  is the threshold value at which  $P_D = 0.5$ . The ideal mean  $\mu_{l_1}$  should be 1 because of normalization. However, it can be seen from the graph that the mean is ‘shifted’ from 1. This is due to the correlation between the quantization noise and the inserted watermark.

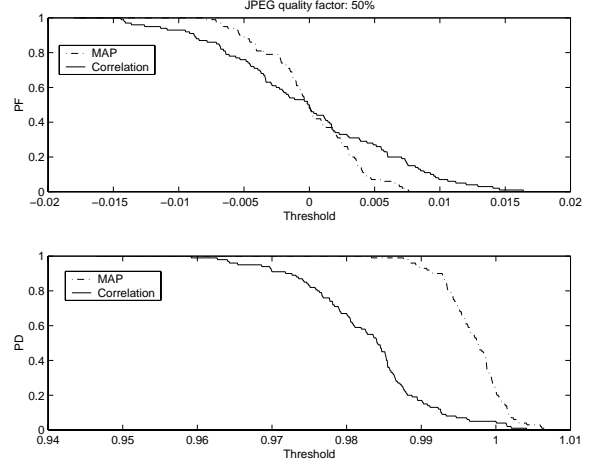


**Fig. 2.** MAP detector performance after JPEG: theoretical estimation vs. simulation

Figure 3 compares the simulation result of the false alarm  $P_F$  and detection probabilities  $P_D$  as functions of detector threshold for the MAP detector and the commonly used correlation detector. It is clear that the MAP detector has higher  $P_D$  and has lower  $P_F$  over the range where  $P_F$  is small.

## 7. CONCLUSIONS

We presented statistical modeling for JPEG compression noise in image watermarking. MAP detector was proposed to improve detector performance. Theoretical estimation of the detector performance was also given. Similar analysis



**Fig. 3.** Detector performance comparison with JPEG compression: MAP vs. Correlation

approach can be applied to other watermarking schemes, such as multi-bit embedding.

## 8. REFERENCES

- [1] F. A. P. Petitcolas, R. J. Anderson, and M. G. Kuhn, “Attacks on copyright marking systems,” in *2nd workshop on information hiding*, Portland, Oregon, Apr. 1998, pp. 218–238.
- [2] I. J. Cox, J. K. F. T. Leighton, and T. Shamoan, “Secure spread spectrum watermarking for multimedia,” *IEEE Trans. on Image Processing*, vol. 6, no. 12, pp. 1673–1686, 1997.
- [3] Julia Minguillon and Jaume Pujol, “Uniform quantization error for laplacian sources with applications to jpeg standard,” in *SPIE Conference on Mathematics of Data/Image Coding, Compression, and Encryption*, San Diego, California, Jul. 1998, pp. 77–88.
- [4] L. Schuchman, “Dither signals and their effect on quantization noise,” *IEEE Trans. on Communication Technology*, no. 12, pp. 162–165, 1964.