

# Efficient Algorithms in Determining JPEG-Effective Watermark Coefficients

Chih-Wei Tang and Hsueh-Ming Hang

Department of Electronics Engineering,  
National Chiao Tung University,  
Hsinchu 30050, Taiwan.

chihwei.ee88g@nctu.edu.tw, hmhang@mail.nctu.edu.tw

**Abstract.** A set of coefficient selection rules is proposed for efficiently determining the effective DCT watermarking coefficients of an image for JPEG attack. These rules are simple in computation but they are derived from the theoretically optimized data set with the aid of the parametric classifiers. They improve the watermark robustness (correctly decoding) and, in the mean time, decrease the error detection probability (correct detection). The frequency versus watermark strength space is used in constructing the selection rules. Simulation results show that the computational complexity is significantly reduced compared to our previous theory-based optimization work, but still the selected coefficients can achieve nearly the same performance as the original scheme.

## 1 Introduction

Many digital watermarking schemes have been proposed for copyright protection, data hiding and other purposes. In our previous work, we focus on the tradeoffs between the achievable watermarking data payload, allowable distortion for information hiding, and robustness against attacks [1]. Although many methods have been developed to improve the watermark data payload and robustness while maintaining reliable detection and visual fidelity [2]-[5], few researchers have proposed techniques to identify the exact coefficient locations for watermarking. Thus, we suggested a generic approach for selecting the most effective coefficients for watermark embedding. Using this set of coefficients improves the watermark robustness and reliability while it maintains the watermark visual transparency. To a certain extent, we try to find the performance limit of invisible watermarking for a given natural image under the assumptions of known attack and non-blind detection for DCT-domain watermarking. The non-blind detection can be used in applications such as transaction-tracking. The synchronization attack is not considered as a problem due to non-blind detection. Since digital images are often compressed for efficient storage and transmission, we use JPEG and JPEG2000 as the examples of attacking sources in the design phase.

Although the coefficient selection procedure performs rather well, its computational complexity is very high. Therefore, in this paper, we develop a fast algorithm with nearly no performance loss. Due to the limited space, only the

simplified rules for JPEG compression attacking source is presented. Note that the methodology of the coefficient selection procedure in [1] and the simplified algorithm proposed in this paper both can easily be extended to the other types of attacks. Section 2 briefly describes our previous work theory-based optimal coefficient selection. Section 3 describes the newly proposed coefficient selection rules. Simulation results are summarized in Section 4 and Section 5 concludes this presentation.

## 2 Our Previous Optimization Algorithm

Two optimization stages are proposed in [1] for selecting effective coefficients. One is the robust and imperceptible coefficient selection stage (Stage One), and the other is the detection reliability improvement stage (Stage Two). Stage One conducts a deterministic analysis on the transform coefficients, and then the proper coefficients and the associated watermark strength are determined so that the coefficients after a specified attack can still bear the valid marks. The additive embedding is adopted in the DCT domain, where  $x[i]$  is the watermark strength of the  $i$ th AC coefficient and  $w[i]$  is the watermark bit. All AC coefficients are watermarked. For an attack in either the spatial or other transform domains, the watermarked image is converted back to the spatial domain and the attack is applied. We decode the watermark bits in the DCT-domain. Several different watermark patterns are tested. If all watermark bits associated with a certain DCT coefficient are correctly decoded, this coefficient is retained in the Stage One candidate set. We examine the all-positive and all-negative watermark patterns. When the attack is not applied to individual coefficients in the DCT-domain, we also test the alternate polarity pattern in which the odd-index watermark bits (in zigzag scan order) are +1 and the even-index ones are -1. This is because the attack distortion on a DCT coefficient also depends on its neighboring watermark bits. Our experiments indicated that we can identify robust coefficients with rather high probability by only 4 patterns. The Watson's visual model is adopted for contrast masking threshold computation and the parameter values are taken from the Checkmark package [6].

Some robust coefficients may produce higher detection error probability. Thus, Stage Two calculates the statistical measures on images and attacks, and it discards the weak coefficients. An iterative procedure is proposed and only one coefficient is discarded in each iteration. At the beginning of one iteration, if  $N$  coefficients remain,  $N$  candidate sets are formed by deleting one coefficient alternatively in this  $N$ -coefficient set. That is, there are  $N-1$  coefficients in each candidate set. Then, the watermark detection statistics based on signal dependent channel distortion model [7] and the Bayes' decision rule for each candidate set is calculated for each candidate set. The error detection probability is the average of the false positive probability and false negative probability. Then, the set with the lowest detection error probability is chosen if the average error probability decreases from the previous iteration. The coefficient discarding process is repeated until the overall error probability cannot be further reduced. If

there are  $N$  selected coefficients at the beginning of Stage Two, and  $K$  dropped coefficients in the process, the execution time of Stage Two will be  $O(KN^2)$ . Thus, a fast algorithm is very desirable.

### 3 Efficient Robust and Reliable Coefficient Selection Rules

Our goal is finding simplified rules to separate the selected coefficients and dropped coefficients for a given input image based on the theoretically optimized data set derived from [1]. We adopt a parametric linear classifier for classification [8]. For a parametric approach, most of the estimated expressions are functions of expected vectors and covariance matrices. Although linear classifiers are not optimum, we use it due to its simplicity. The classifier (linear discriminate function) is

$$h(X) = V^T X + v_0, \tag{1}$$

where  $X$  is the given input data vector which distributions are not limited,  $V = [v_1 v_2 \dots]^T$  is the coefficient vector, and  $v_0$  is a threshold value. To find the optimal  $V^T$  and  $v_0$  for a given distribution, the criterion  $g$  is maximized, which measures the between-class scatter normalized by the within-class scatter,

$$g = \frac{P_1 \eta_1^2 + P_2 \eta_2^2}{P_1 \sigma_1^2 + P_2 \sigma_2^2}, \tag{2}$$

where  $P_i$ ,  $\eta_i$ , and  $\sigma_i$  are the priori probability, expected value of  $h(X)$ , and variance of  $h(X)$  for class  $i$ , respectively. As a result,

$$V = [P_1 \Sigma_1 + P_2 \Sigma_2]^{-1} (M_2 - M_1), \tag{3}$$

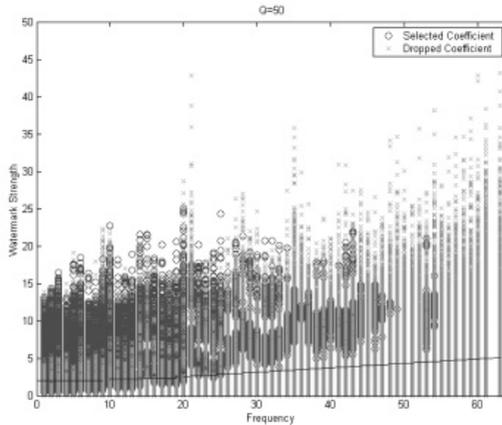
$$v_0 = -V^T [P_1 \Sigma_1 + P_2 \Sigma_2], \tag{4}$$

where  $\Sigma_i$  is the covariance matrix for a given expected vector  $X$ . Here, the well known fisher criterion is not adopted since it cannot determine the optimum  $v_0$ . The features in our problem are frequency  $f$ , amplitude  $x$  and admissible watermark strength  $\alpha$ . Our target is to find a piece-wise linear classifier (discriminator) that separate the selected coefficients from the dropped ones. We have looked at the case that uses all three features  $(f, x, \alpha)$  (3-D domain). To simplify calculations, we also search for a 2-D feature space with smallest average misclassification rate. Our experiments show that the "optimal" average misclassification rate in the 2-D space " $(f, \alpha)$ " is only 1% lower than that of the 3-D domain classifier. There are three 2-D domain candidates:  $(f, x)$ ,  $(f, \alpha)$ , and  $(x, \alpha)$ . Let  $D_{fx}$ ,  $D_{f\alpha}$  and  $D_{x\alpha}$  be the misclassification rate due to the selected coefficients are misclassified as dropped coefficients in the aforementioned three candidate spaces, respectively,  $S_{fx}$  and  $S_{f\alpha}$ , and  $S_{x\alpha}$  be the misclassification rate due to the dropped coefficients are misclassified as selected coefficients. To decrease  $S_{fx}$  and  $S_{f\alpha}$ , and  $S_{x\alpha}$ , we set  $P_1 = 0.4$  and  $P_2 = 0.6$ . For further improving the classification accuracy, we divide a space into three subspaces,

and design one linear classifier for each subspace. For  $(f, x)$  and  $(f, \alpha)$  spaces, the separation is based on  $f=0-9$ ,  $f=10-19$  and  $f=20-63$ . For space  $(x, \alpha)$ , they are  $x=0-49$ ,  $x=50-99$  and  $x=100-\infty$ . Our image data base contains 30 natural images. The training set is generated using the method described in Sect. 2. Four JPEG quality factors ranging from 50 to 80 are used. We adopt the definition of JPEG quantization step size defined in [9]. The misclassification rates in all cases (2D domains) are listed in Table 1. Because the best 2-D  $(f, \alpha)$  space is 1% worse than the 3-D  $(f, x, \alpha)$  classifier, the former is adopted for a much lower computation complexity.

**Table 1.** Misclassification rates in three 2-D feature spaces

Design Phase	$S_{fx}$	$S_{f\alpha}$	$S_{x\alpha}$	$D_{fx}$	$D_{f\alpha}$	$D_{x\alpha}$
JPEG50	0.31	<b>0.18</b>	0.66	0.07	<b>0.10</b>	0.40
JPEG60	0.29	<b>0.18</b>	0.65	0.06	<b>0.09</b>	0.38
JPEG70	0.27	<b>0.18</b>	0.63	0.06	<b>0.09</b>	0.35
JPEG80	0.27	<b>0.16</b>	0.57	0.05	<b>0.08</b>	0.30



**Fig. 1.** The classifier at JPEG quality factor 50 with coefficients from 30 natural images

Thus, we can now select effective watermarking coefficients with the simplified rules. Figure 1 shows the classifier (coefficient selecting rules) for the JPEG quality factor 50 in the design phase. Although these rules eliminate a number of poor candidate coefficients, the remaining coefficients do not necessarily have the required robustness. Therefore, we apply the original Stage One process to the retained coefficients for further removing weak coefficients.

## 4 Simulation Results

To examine the performance of the proposed rules, we test images which are not used in training. Limited by space, only the results for pictures Lena and Baboon are included. For the JPEG quality factor 50 in the design phase, the PSNR values between the original and the watermarked images are 45.2 dB and 39.98 dB for Lena and Baboon, respectively. And, they are 42.9 dB and 36.82 dB for JPEG quality factor 80 in the design phase. The embedded watermarks are invisible as we inspect them visually. The comparisons between the original and the simplified schemes are shown in Tables 2 and 3. Let the overlapped percentage be the number of coefficients selected by both the original Stage One and the simplified scheme divided by the number of selected coefficients by the original Stage One. We find that the overlapped percentage is higher than 70%. The detection error probability using the simplified scheme is still very small (all less than  $10^{-135}$  for Lena). Practically these rules are as good as the original massive iteration scheme. In the case of Baboon image, the overlapped percentage is over 85% and the detection error probability is all less than  $10^{-245}$ . The data shown in Fig. 2 is each averaged over 5000 watermarked images with different random watermark sequences. Also, the same 5000 watermark sequences are correlated with the unmarked but JPEG compressed image and the results are averaged in Fig. 3. Figure 3 shows that the selected coefficient survives JPEG compression at higher quality factors may not survive JPEG compression at lower quality factors. To verify the designed false negative and positive error probabilities, the mean, variance, minimum and maximum values of the normalized correlation sum after the JPEG attacks are computed. (The normalization is normalized against the embedded watermark power as discussed in [1], and thus is not bounded to  $[-1, 1]$ .) Due to the limited space, only the mean of the normalized correlation sum for watermarked images is shown in Fig. 2. The mean value of the normalized correlation sum  $C$  is computed by

$$C = \frac{1}{M} \sum_{i=1}^M c[i] = \frac{1}{M} \sum_{i=1}^M \frac{y[i] \times (w[i] \times \alpha[i])}{\sigma_d^2}, \quad (5)$$

where  $y[i]$  is the difference between the DCT coefficients of the received image and the original image,  $w[i]$  is the watermark signature and  $M$  is the number of selected coefficients. For a watermark sequence,  $C$  is compared against the detection threshold which is approximately the average of the mean values of the normalized correlation sum of the watermarked  $E\{c|H_1\}$  and unmarked images  $E\{c|H_0\}[1]$ . The presence of the watermark is declared if  $H_1$  is favored. In all cases, there is no failure for either watermarked or unmarked 5000 images. Finally, small variance implies lower error detection probability. The variance values  $Var\{c|H_0\}$  and  $Var\{c|H_1\}$  are all smaller than 0.0018 after JPEG attacks with different quality factors for both watermarked and unmarked cases. We also test the JPEG-robust watermark against several other signal processing attacks by as shown in Fig. 4 and the data are obtained by averaging over 100 different random watermark sequences. The  $E\{c|H_1\}$  is over 0.8 after JPEG2000

**Table 2.** The comparisons of the selected coefficients for Lena

Design Phase	No. of Selected Coeff. by Org. Stage 1	No. of Selected Coeff. by Org. Stage 2	Estimated $P_{error}$ after Org. Stage 2	No. of Selected Coeff. by Fast Scheme	Estimated $P_{error}$ by Fast Scheme
JPEG50	4738	4019	5.505e-299	3609	2.097e-136
JPEG60	6007	5082	0.000e+000	4516	6.803e-181
JPEG70	8041	6587	0.000e+000	5911	2.320e-253
JPEG80	111473	9439	0.000e+000	8166	0.000e+000

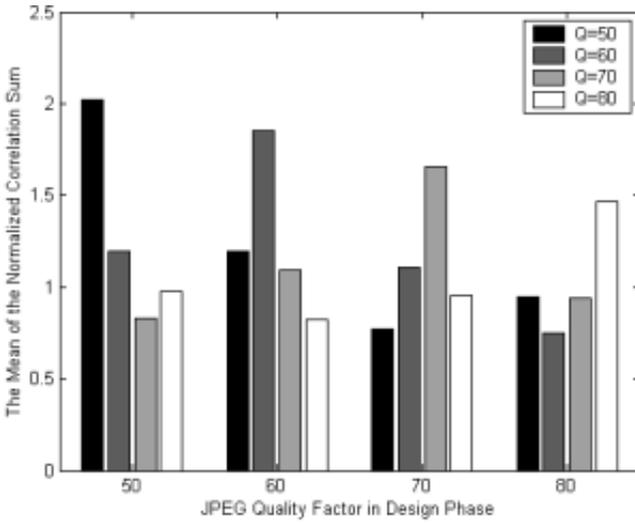
**Table 3.** The comparisons of the selected coefficients for Baboon

Design Phase	No. of Selected Coeff. by Org. Stage 1	No. of Selected Coeff. by Org. Stage 2	Estimated $P_{error}$ after Org. Stage 2	No. of Selected Coeff. by Fast Scheme	Estimated $P_{error}$ by Fast Scheme
JPEG50	9270	7359	0.000e+000	7877	3.099e-246
JPEG60	11743	8972	0.000e+000	10130	0.000e+000
JPEG70	15912	13105	0.000e+000	13708	0.000e+000
JPEG80	22931	18885	0.000e+000	20120	0.000e+000

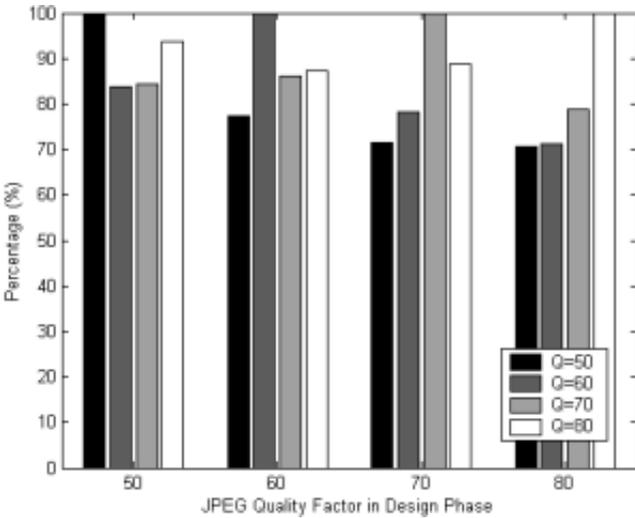
**Table 4.** The comparisons of the selected coefficients for Baboon

Design Phase	No. of Processed Coeff. by Org. Stage 1	No. of Processed Coeff. by Org. Stage 2	No. of Processed Coeff. by Fast Scheme
JPEG50	64512	3152520	73629
JPEG60	64512	5134207	74108
JPEG70	64512	10641870	75139
JPEG80	64512	21277960	76366

attacks at bit rates 0.125 bpp and 0.0625 bpp. We also compare the computational complexity between the original and the simplified stages as shown in Table 4. The computational complexity is expressed by the number of processed DCT coefficients. For image Lena at JPEG quality factor 80, the simplified scheme requires roughly  $\frac{1}{266}$  of the computations of the original scheme (Stage One + Stage Two) for large candidate sets. The simplified scheme does greatly reduce the computational complexity.



**Fig. 2.** The mean of the normalized correlation sum after JPEG attacks at different quality factors for watermarked Lena



**Fig. 3.** The percentage of correctly decoded coefficients at the detector after JPEG attacks for Lena

## 5 Conclusions

In this paper, we propose an efficient algorithm for selecting JPEG-effective watermark coefficients. In most cases, the new scheme uses only  $\frac{1}{100}$  of the computation needed in the original scheme in [2]. The methodology of both the original coefficient selection procedure in [2] and the simplified algorithm proposed here can be easily extended to the other types of attacks.

**Acknowledgements.** This work is partially supported by the Lee and MTI Center for Networking Research at National Chiao Tung University, Taiwan.

## References

1. C.-W. Tang and H.-M. Hang: Exploring Effective Transform-Coefficients in Perceptual Watermarking. Proc. SPIE Security and Watermarking of Multimedia Contents IV. **4675** (2003) 572–583
2. P. Moulin and J. A. O’Sullivan: Information-Theoretic Analysis of Information Hiding. IEEE Trans. Information Theory. **49** (2003) 563–593
3. Q. Cheng, Y. Wang and T. S. Huang: How to Design Efficient Watermarks? IEEE International Conference on Acoustics, Speech, and Signal Processing. **3** (2003) 49–52
4. M. Barni, F. Bartolini, A. D. Rosa and A. Piva: Optimum Decoding and Detection of Multiplicative Watermarks. IEEE Trans. Signal Processing. **51** (2003) 1118–1123
5. A. Giannoula, A. Tefas, N. Nikolaidis and I. Pitas: Improving the Detection Reliability of Correlation-Based Watermarking Techniques. ICME. **1** (2003) 209–212
6. <http://watermarking.unige.ch/Checkmark/index.html>
7. J. J. Eggers and B. Girod: Quantization Effects on Digital Watermarks. Signal Processing. **831** (2000) 239–253
8. K. Fukunaga: Introduction to Statistical Pattern Recognition. Academic Press. (1990)
9. <http://www.cl.cam.ac.uk/fapp2/watermarking/stirmark/>