# Chapter 4

# Exploring Effective Coefficients in DCT-Domain Perceptual Watermarking

## 4.1 Introduction

Several works studying the watermark data payload issue have been published using the theoretical analysis approach [3][46][47]. Clearly, there are tradeoffs between the achievable watermarking rate, allowable distortion for information hiding, and robustness against attacks [3]. It has been reported that the transform–domain watermarking techniques can offer a higher capacity under specific attacks (such as compression) [46][48]. For our targeting applications, we are especially interested in the watermark payload and robustness under the combined criteria of reliable detection and visual fidelity. The perceptual watermark payload in different transform domains has been analyzed in [49]. In [47], the capacity constrained by reliable statistical detection is calculated. In [50], the minimum number of coefficients in discrete wavelet domain with spread spectrum watermark embedding is theoretically analyzed using the human visual model and a probabilistic detection model.

The previous research work estimates the theoretical watermark capacity bound of, for example, transform-domain watermarking, but the exact locations of the coefficients for watermark embedding are not identified. In this paper, we develop a procedure that identifies these effective coefficients for natural images. The goal is to achieve both high detection reliability and watermark invisibility. Since the achievable rate (data payload) in the case of

blind detection is upper-bounded by that of non-blind detection [3], we employ non-blind detection and assume the attack source is known to explore the best achievable performance under the above assumptions. To a certain extent, we are exploring the "performance limit" of DCT-domain watermarking for a given specific attack. In general, this method can be extended to other watermarking schemes and/or with multiple attacks. In practice, the non-blind detection watermarks can be used in the applications such as transaction-tracking [1].

Since digital images are often compressed for efficient storage and transmission, in this study, JPEG and JPEG2000 compression schemes are adopted as the attacking sources. In general, other attacks can be used in the design (training) phase, because in our approach, the watermarking algorithm and the attacking process are unrelated.

A two–stage procedure is developed for choosing appropriate coefficients. In the first stage, deterministic analysis is exercised to pick up the proper coefficients so that the attacked coefficients can still hold the retrievable mark and in the meanwhile the distortion due to watermark embedding is lower than the visual threshold. In the second stage we calculate the statistical properties of watermarked images and discard the coefficients that reduce detection reliability (or equivalently, increase error probability). In Section 4.2, the robust and imperceptible coefficient selection process (stage 1) is developed. Section 4.3 describes the human visual masking model used in our experiment. Section 4.4 contains the description of the detection reliability improvement process (stage 2). Section 4.5 and 4.6 cover the details of the watermark embedding and detection procedures. Simulation results summarized in Section 4.7 will show the performance of our scheme. At the end, Section 4.8 concludes this presentation.

# 4.2 Robust and Imperceptible Coefficient Selection

Our goal is to achieve the maximum detection robustness while the watermark imperceptible property can still be maintained. Several factors affect the watermark detection ability. In the case of DCT-domain watermarking, intuitively one may want to use more coefficients. However, some coefficients with low energy, for example, may be inappropriate for carrying watermarks. Similar to the signal design problem in digital transmission over noisy channels, signals (now DCT coefficients) have to be carefully selected to achieve the robustness goal. Increasing the magnitude of watermark generally increases the watermark robustness. But on the other hand, large-magnitude changes on coefficients may be perceptually visible. Also different types (and amount) of attacks produce different-levels of damages on the watermarks. The coefficients that can tolerate a specific type of attack can be identified with the aid of damage analysis due to potential attacks [68]. For example, if JPEG compression is used for image distribution and thus may be viewed as an attack to the embedded watermark. In this case, we can include JPEG compression in the design phase in choosing the most robust watermark parameters.

We first assume both the attacks and the watermark embedding method are known. In our first experiment, JPEG compression is used as the attacking means. And as said earlier, we adopt the DCT-domain watermarking embedding technique. The robustness of watermark (correctness of decoded watermark bits) can be increased by either selecting proper coefficients and/or adjusting watermark embedding parameters based on the input images.

We use the DCT-domain additive watermark embedding, $x'[i] = x[i] + \alpha[i] \cdot w[i]$, where $\alpha[i]$ is the watermark strength for the DCT coefficient $x[i]$, and its value is decided by the visual threshold (details to be described in Section 4.3). A DCT coefficient $x[i]$ is positively watermarked if the watermark bit $w[i]$ is +1, and it is negatively watermarked if $w[i]$ is –1.

There are two stages in our proposed scheme. In the first robustness and imperceptibility coefficient selection stage (Robustness stage in short), the attack effect on the watermarked coefficients are checked for robustness. In the JPEG attack case, a DCT coefficient is declared robust and is selected if both its positive and negative watermark embedding can survive the attack. On the other hand, because the human eyes are rather sensitive to low-frequency coefficient variations, we do not embed the watermark on the DC coefficients. The robustness and imperceptibility of all the AC coefficients are examined.

The embedded watermark strength $\alpha[i]$ of the $i$th AC coefficient $x[i]$ is set to be the visual masking threshold $hvs[i]$ of this coefficient. Let the quantization step size of JPEG compression with quality factor $q$ be $\Delta_q$. The JPEG compression is applied to both the positively and negatively watermarked values of the same DCT coefficient. As a result, the distortion between the unwatermarked AC coefficient $x[i]$ and its quantized positively watermarked coefficient $x'[i]$ is $e_{posw}[i]$; that is,

$$
e_{posw}[i] = \begin{cases} \Delta_q \cdot Round\left\lfloor \dfrac{\Delta_q/2 + (x[i] + hvs[i])}{\Delta_q} \right\rfloor - x[i], \text{ if } (x[i] + hvs[i]) \geq 0 \\ \Delta_q \cdot Round\left\lfloor \dfrac{-\Delta_q/2 + (x[i] + hvs[i])}{\Delta_q} \right\rfloor - x[i], \text{ if } (x[i] + hvs[i]) < 0 \end{cases}, \qquad (4.1)
$$

and the distortion between the original AC coefficient $x[i]$ and its quantized negatively watermarked coefficient $x'[i]$ is $e_{negw}[i]$; that is,

$$
e_{negw}[i] = \begin{cases} \Delta_q \cdot Round\left\lfloor \dfrac{\Delta_q/2 + (x[i] - hvs[i])}{\Delta_q} \right\rfloor - x[i], \text{ if } (x[i] - hvs[i]) \geq 0 \\ \Delta_q \cdot Round\left\lfloor \dfrac{-\Delta_q/2 + (x[i] - hvs[i])}{\Delta_q} \right\rfloor - x[i], \text{ if } (x[i] - hvs[i]) < 0 \end{cases}. \qquad (4.2)
$$

A DCT coefficient is retained in the candidate set if the sign of $e_{posw}[i]$ is +1 and the sign of

$e_{negw}[i]$ is −1.

The sets of robust coefficients under JPEG compression corresponding to different quality factors are separately produced. Since $\Delta_q$ decreases as the quality number of the JPEG compression increases, a coefficient that survives the JPEG compression with a low quality factor (e.g., 50) usually also survives the JPEG compression with a higher quality factor (e.g., 80). Therefore, the number of the coefficients passing a low-quality-factor JPEG attack is smaller than that passing a higher quality-factor attack. These selected coefficients will be further screened in Section 4.4.

The preceding approach can be extended to cover the other types of attacks. An original image is embedded with the watermark in the DCT-domain. 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 then applied. We decode the watermark bits in the DCT-domain. If the watermark bit associated with a certain DCT coefficient is correctly decoded, this coefficient is declared to be robust.

However, there is a major difference when the attack is not applied to individual coefficients in the DCT-domain. In the case of JPEG, where each DCT coefficient is separately quantized, the watermark robustness of a DCT coefficient is easily examined as we did in (1) and (2). In general, the amount of distortion on a single DCT-coefficient due to an attack depends on the entire watermark pattern. That is, the attack distortion on a certain DCT coefficient depends on not only the watermark bit added to that coefficient but also its neighboring watermark bits. If the watermark pattern changes, the same attack may produce different distortion on a particular coefficient. Therefore, to make sure a coefficient can survive for any embedded watermark pattern, we should examine all possible watermark patterns. Clearly, this is not possible when the watermark bits are large (say, more than 20 or 30 bits). Looking for "best" test patterns can be a research topic. Heuristically, we first examine the all-positive and all-negative patterns. Then, we 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. Reversing the polarity of the previous pattern produces another alternate polarity pattern. Additional patterns included in the robustness checking process will improve the accuracy of picking up robust coefficients at the cost of computational complexity. Our experiments indicate that even with only 4 patterns, we can identify robust coefficients with rather high probability.

If the target is to design a watermark that survives a specific attack, then the preceding procedure is sufficient and no further refinement is needed. However, in the probabilistic sense (and so that it may be practically useful), we hope the embedded watermark can survive "similar" attacks (with the same probabilistic model) with high reliability. Also, we like to reduce the false alarm probability in which a watermark is detected although none is really embedded. Therefore, a second stage of increasing reliability is added.

# 4.3 Visual Masking Effect in DCT Domain

For the purpose of imperceptible watermark design, we are particularly interested in the masking properties of the human visual system (HVS) [69]. Masking effect means that the visibility of one signal (image) is changed due to the existence of the other (image) signal. Several visual masking effects have been identified such as spatial masking, luminance masking, and contrast masking.

The inclusion of human perceptual characteristics into the watermarking design process helps maintaining the watermark imperceptibility. Another advantage of this approach is that if the watermarked image spectra is similar in shape to the spectrum of the original image, then the attackers cannot easily identify the embedded watermark by using some prior knowledge on the image's statistics [36][44][45].

In the following procedure, we use the visual masking model in the DCT domain since our watermark embedding process is conducted in this domain. The visual masking thresholds

are calculated only for AC coefficients because the DC coefficients are not marked. The popular Watson's DCT-based visual model is employed in calculating the contrast masking threshold $e_{mnk}$ of the AC coefficient $x[i]$ at the 2-D frequency index $(m,n)$ of block $k$. The visual threshold $hvs[i]$ is thus set to $e_{mnk}$ and it is used to adjust the watermark embedding strength $\alpha[i]$ as described in Section 4.2. The contrast masking effect often has the strongest impact on the subjective visual quality.

In our experiment, the parameter values used in contrast masking threshold calculation are the same as those used in the Checkmark package [70]. This set of setting is decided through subjective tests and is widely adopted in image research. More details can be found in [70] and [71]. Here, we only briefly describe its computational steps as below.

1. Set $W_x$ to $(180/\pi) \times (v/(344 \times d_v))$, $W_y$ to $(180/\pi) \times (u/(342 \times d_v))$, where $v$ is the vertical screen image size, $u$ is the horizontal screen image size, and $d_v$ is the viewing distance. In our experiments, $v$ is 8.8, $u$ is 9.4, and $d_v$ is 72.

2. Calculate $f_{mn} = \frac{1}{16}\sqrt{(m/W_x)^2 + (n/W_y)^2}$, where $(m,n)$ is the 2-D frequency index of an $8 \times 8$ DCT block.

3. Calculate $r_{mn} = r + (1-r)\cos^2\theta_{mn}$, where $\theta_{mn} = \sin^{-1}\frac{2f_{m0}f_{0n}}{f_{mn}^2}$ and $r = 0.7$.

4. Calculate $T_{\min} = \begin{cases} \dfrac{L}{S_0} & \text{if } L > L_T \\ \dfrac{L}{S_0}(\dfrac{L_T}{L})^{1-a_t} & \text{if } L \le L_T \end{cases}$, where $L = 10^{7/((255 \times 128)-4)}$, $L_T = 13.45$, $S_0 = 94.7$, and $a_t = 0.649$.

5. Calculate $k = \begin{cases} k_0 & \text{if } L > L_k \\ k_0(\dfrac{L}{L_k})^{a_k} & \text{if } L \le L_k \end{cases}$, where $L_k = 300$, $k_0 = 3.125$, and $a_k = 0.0706$.

6. Calculate $f_{\min} = \begin{cases} f_0 & \text{if } L > L_f \\ f_0(\dfrac{L}{L_f})^{a_f} & \text{if } L \le L_f \end{cases}$, where $L_f = 300$, $a_f = 0.182$, and $f_0 = 6.78$.

7. Calculate $\log_{10} t_{mn} = \log_{10}\dfrac{T_{\min}}{r_{mn}} + k(\log_{10} f_{mn} - \log_{10} f_{\min})^2$.

8. Calculate the luminance masking threshold at frequency index $(m,n)$ for block $k$ :

$$t_{mnk} = t_{mn}(\frac{c_{00k}}{c_{00}})^{a_t}, \text{ where } c_{00} = 128 \times 8, \text{ and } c_{00k} \text{ is the DC coefficient of block } k.$$

9. Calculate the contrast masking threshold: $e_{mnk} = \max[t_{mnk}, |c_{mnk}|^{w_{mn}} \cdot t_{mnk}^{1-w_{mn}}]$, where

$w_{mn}$ is chosen experimentally. We set $w_{00} = 0$ and $w_{mn} = 0.7$ for $(m,n) \neq (0,0)$. Note

that $c_{mnk}$ is the $(m,n)$ AC coefficient of block $k$.

# 4.4 Detection Reliability Improvement

A watermarking system can be viewed as a communication system with, possibly, side information [72]. If the watermark detector is known and the type of attacks is also known in advance, the coefficients that have higher detection error probability can be pre-estimated and dropped to improve the overall detection reliability. There are two types of error probability. The false positive probability that an unmarked image is wrongly declared watermarked by the detector is $P_{FP}$. On the other hand, the probability of undetected watermark is false negative error probability, $P_{FN}$. The average error probability is $P_{error} = (P_{FP} + P_{FN})/2$ if we assume an image is equally likely marked or unmarked. Let $H_0$ denote the state that an image is marked and $H_1$ denote the watermarked state.

In the second stage, DCT coefficients are further screened to enhance the detection reliability. For this purpose, we like to know how the detection error probability of a particular coefficient is affected by a specific attack distortion model (JPEG compression, say). We first collect statistics from the real image data by running the (JPEG) attack on the watermarked images. Then, the error probability is estimated based on a statistical model of the distorted watermarked coefficients. In [73][74], a theoretical model for additive watermarks under JPEG quantization effect is proposed based on the dither quantization theory [75]. The pseudo-noise watermark and the original image are assumed to be statistically independent. It was shown that the JPEG quantization distortion on individual

coefficient cannot be approximated by an AWGN channel model and such distortion should be signal dependent. In particular, the distributions of the fine and coarse quantization errors are different. Therefore, we take the approach based on the central limit theorem [73] instead of applying the normal distribution model to the individual coefficient. That is, the mean value of the normalized correlation sum can be approximated by the normal distribution. This model can be also extended to other attacking sources.

The candidate coefficients that have passed the robustness and imperceptibility stage (stage 1) in Section 4.2 are further examined against the reliability test at the reliability improvement stage (stage 2). We propose an iterative procedure to discard the "poor" coefficients. Only one coefficient is discarded in each iteration. This process continues until the overall error probability cannot be further reduced. At the beginning of one iteration, if there are $N$ coefficients, $N$ candidate sets are formed by deleting one coefficient alternatively in this $N$-coefficient set. Consequently, there are $N$-1 coefficients in each candidate set. Then, the statistics of each candidate set based on the Bayes' decision rule for watermark detection is calculated separately. The set with the lowest error probability is retained if the overall error probability decreases monotonically. Clearly, our proposal is one special type of searching algorithms. There are other searching algorithms that may be employed in this stage. The procedure of our algorithm is described below.

The detection error probability is calculated based on the watermark detection rule. Here, the watermark detection rule is designed to minimize the average cost using the Bayes' rule. The binary hypotheses of watermark detection for a received image are:

$$
\begin{aligned}
&\mathrm{H}_0 : y[i] = (x[i] + e_{\mathrm{H}_0}[i]) - x[i] = e_{\mathrm{H}_0}[i] \\
&\mathrm{H}_1 : y[i] = (x[i] + d[i] + e_{\mathrm{H}_1}[i]) - x[i] = d[i] + e_{\mathrm{H}_1}[i]
\end{aligned}
$$

,

where $y[i]$ is the difference between the received coefficient and the unwatermarked

coefficient $x[i]$, $d[i]$ is the embedded watermark, $e_{H_0}[i]$ is the distortion due to attack on the

original coefficient, and $e_{H_1}[i]$ is the attack distortion on the watermarked coefficient.

Let $c[i]$ be the normalized correlation value between $y[i]$ and $d[i]$ and $C$ be the mean

value of the normalized correlation sum. Let $c_{10}$ be the cost of the false positive decision,

$c_{01}$ be the cost of the false negative decision, $c_{00}$ be the cost of detecting watermark

correctly, and $c_{11}$ be the cost of detecting the absence of watermark correctly. Then, the

Bayes' decision rule for minimum cost implies that $H_1$ is chosen if [76]

$$H_1 : \quad \frac{P(C\,|\,H_1)}{P(C\,|\,H_0)} > K = \frac{(c_{10} - c_{00}) \cdot P(H_0)}{(c_{01} - c_{11}) \cdot P(H_1)} \quad , \tag{4.3}$$

where $C$ is an estimated value of $E\{c\}$,

$$E\{c\} \approx C = \frac{1}{M}\sum_{i=1}^{M} c[i] = \frac{1}{M}\sum_{i=1}^{M} \frac{y[i] \cdot d[i]}{\sigma_d^2} \quad \text{and} \quad d[i] = w[i] \cdot \alpha[i] \;. \tag{4.4}$$

As described earlier that $w[i]$ is the watermark signature with antipodal signaling $\{-1,1\}$, and

$\alpha[i]$ is the adjustable watermark embedding strength of $x[i]$ (described in Section 4.2). In

(4.4), $M$ is the number of the watermarked coefficients, and $\sigma_d^2$ is the variance of embedded

watermark signals. Since $\sigma_d^2$, rather than $\dfrac{1}{M}\sqrt{\sum\limits_{i=1}^{M} y^2[i] \cdot \sum\limits_{i=1}^{M} d^2[i]}$, is in use in the denominator, $C$

is not bounded to [-1, 1]. When $M$ is sufficiently large, the probability distribution of $C$ can be

approximated by the Gaussian distribution according to the central limit theorem.

The variance of $C$ is

$$\text{Var}\{C\} = \frac{1}{M}\text{Var}\{c\} \;. \tag{4.5}$$

Therefore, the left hand side of the decision rule (4.3) becomes

$$H_1 : \quad \frac{(2\pi\,\text{Var}\{c\,|\,H_1\}/M)^{-1/2}\exp(-\dfrac{(C - E\{c\,|\,H_1\})^2}{2\text{Var}\{c\,|\,H_1\}/M})}{(2\pi\,\text{Var}\{c\,|\,H_0\}/M)^{-1/2}\exp(-\dfrac{(C - E\{c\,|\,H_0\})^2}{2\text{Var}\{c\,|\,H_0\}/M})} > K \quad . \tag{4.6}$$

Equivalently,

$$H_1 : \begin{array}{l} 2\log K + \log(\dfrac{\mathrm{Var}\{c \mid H_1\}}{\mathrm{Var}\{c \mid H_0\}}) + \dfrac{E^2\{c \mid H_1\}}{\mathrm{Var}\{c \mid H_1\}/M} - \dfrac{E^2\{c \mid H_0\}}{\mathrm{Var}\{c \mid H_0\}/M} \\[4mm] < (\dfrac{1}{\mathrm{Var}\{c \mid H_0\}/M} - \dfrac{1}{\mathrm{Var}\{c \mid H_1\}/M}) \cdot C^2 + 2(\dfrac{E\{c \mid H_1\}}{\mathrm{Var}\{c \mid H_1\}/M} - \dfrac{E\{c \mid H_0)}{\mathrm{Var}\{c \mid H_0\}/M}) \cdot C \end{array} \qquad (4.7)$$

Finally, the maximum-likelihood (ML) detector is obtained with $K=1$ in (4.3) assuming that (i) $c_{10} - c_{00} = c_{01} - c_{11}$ (symmetric cost function), and (ii) $P(H_0) = P(H_1)$. Then, (4.7) can be simplified and expressed as $(C - x_1)(C - x_2) > 0$, where $x_1$ and $x_2$ are two constants determined by $E\{c \mid H_0\}$, $E\{c \mid H_1\}$, $\mathrm{Var}\{c \mid H_0\}$, and $\mathrm{Var}\{c \mid H_1\}$. The detection threshold $x_c$ is either $x_1$ or $x_2$ as its value should locate between $E\{c \mid H_0\}$ and $E\{c \mid H_1\}$. As a result, the image is declared watermarked if $C > x_c$. Consequently,

$$P_{FP} = \int_{x_c}^{\infty} \frac{\exp(-\dfrac{(x - E\{c \mid H_0\})^2}{2\,\mathrm{Var}\{c \mid H_0\}/M})}{\sqrt{2\pi \mathrm{Var}\{c \mid H_0\}/M}} dx = \frac{1}{2} erfc(\sqrt{M}\,\frac{x_c - E\{c \mid H_0\}}{\sqrt{2\,\mathrm{Var}\{c \mid H_0\}}}), \qquad (4.8)$$

$$P_{FN} = \int_{-\infty}^{x_c} \frac{\exp(-\dfrac{(E\{c \mid H_1\} - x)^2}{2\,\mathrm{Var}\{c \mid H_1\}/M})}{\sqrt{2\pi \mathrm{Var}\{c \mid H_1\}/M}} dx = \frac{1}{2} erfc(\sqrt{M}\,\frac{E\{c \mid H_1\} - x_c}{\sqrt{2\,\mathrm{Var}\{c \mid H_1\}}}). \qquad (4.9)$$

To estimate $P_{FP}$ and $P_{FN}$ in (4.8) and (4.9), the statistics $E\{c \mid H_0\}$, $E\{c \mid H_1\}$, $\mathrm{Var}\{c \mid H_0\}$, and $\mathrm{Var}\{c \mid H_1\}$ are derived from the image data. We assume a coefficient is equal likely being positively or negatively watermarked. Then, $E\{c \mid H_0\}$, $E\{c \mid H_1\}$, $\mathrm{Var}\{c \mid H_0\}$, and $\mathrm{Var}\{c \mid H_1\}$ are calculated by the following equations:

$$E\{d\} = \frac{1}{M} \sum_{i=1}^{M} d[i], \text{ where } d[i] = w[i] \cdot \alpha[i] \qquad (4.10)$$

$$\sigma_d^2 = \mathrm{Var}\{d\} = \frac{1}{M-1} \cdot \sum_{i=1}^{M} (d[i] - E\{d\})^2 = (\frac{1}{M-1} \sum_{i=1}^{M} d^2[i]) - (\frac{M}{M-1} E^2\{d\}) \qquad (4.11)$$

$$E\{c \mid H_0\} = \frac{1}{M} \sum_{i=1}^{M} c_{H_0}[i] = \frac{1}{M} \sum_{i=1}^{M} \frac{y_{H_0}[i] \cdot d[i]}{\sigma_d^2} = \frac{1}{M} \sum_{i=1}^{M} \frac{e_{H_0}[i] \cdot d[i]}{\sigma_d^2} \qquad (4.12)$$

$$E\{c \mid H_1\} = \frac{1}{M} \sum_{i=1}^{M} c_{H_1}[i] = \frac{1}{M} \sum_{i=1}^{M} \frac{y_{H_1}[i] \cdot d[i]}{\sigma_d^2} = \frac{1}{M} \sum_{i=1}^{M} \frac{(d[i] + e_{H_1}[i]) \cdot d[i]}{\sigma_d^2} \qquad (4.13)$$

$$\begin{array}{l} \mathrm{Var}\{c \mid H_0\} = \dfrac{1}{M-1} \cdot \sum_{i=1}^{M} (c_{H_0}[i] - E\{c \mid H_0\})^2 \\[4mm] \qquad = \dfrac{1}{M-1} \sum_{i=1}^{M} (\dfrac{e_{H_0}[i] \cdot d[i]}{\sigma_d^2})^2 - (\dfrac{M}{M-1} E^2\{c \mid H_0\}) \end{array} \qquad (4.14)$$

$$\text{Var}\{c \mid \text{H}_1\} = \frac{1}{M-1} \cdot \sum_{i=1}^{M} (c_{\text{H}_1}[i] - \text{E}\{c \mid \text{H}_1\})^2$$

$$= \frac{1}{M-1} \sum_{i=1}^{M} \left( \frac{(d[i] + e_{\text{H}_1}[i]) \cdot d[i]}{\sigma_d^2} \right)^2 - \left( \frac{M}{M-1} \text{E}^2\{c \mid \text{H}_1\} \right) \tag{4.15}$$

In each iteration, the average error probability of (4.8) and (4.9) is computed for every candidate set. The minimum error set is singled out and therefore one coefficient is removed. This iterative process repeats until the average error does not decrease any further. Note that the coefficients sets associated with different JPEG compression quality numbers are different as discussed in Section 4.2. And thus, $M$ is the number of total selected coefficients for the entire image associated with a JPEG quality factor.

A slight variation of the above scheme is formed for selecting a pre-determined number of coefficients. The iteration procedure is similar as before; that is, "drop the least reliable coefficient". However, the stopping rule is changed to "stop when the number of coefficients reaches the pre-determined number". We can also use the same framework for picking up the largest set of coefficients that meet a pre-selected error probability.

# 4.5 Watermark Embedding Scheme

The watermark embedding process is described as follows. First, an original image is transformed by $8 \times 8$ non-overlapped 2-D DCT and the contrast masking thresholds of all AC coefficients are calculated. Then the locations of robust coefficients are determined by the Robustness and the Reliability stages as described before. After the robust coefficients have all been selected, watermarks are embedded on the selected DCT coefficients. If we group the DCT coefficients with the same 2-D frequency index together to form a sub-channel, there are in total 63 sub-channels. Typically, the AC coefficients belonging to a sub-channel can be modeled as a generalized Gaussian distribution source [77]. The watermark bits are inserted to the coefficients belonging to the same sub-channel in the raster-scan order, and then from the lower frequency sub-channels to the higher frequency. In fact, the order of embedding is

not important in our scheme because the watermark detection is judged based on the entire sequence and is irrelevant to the coefficient order. The watermarked coefficient is generated by $x'[i] = x[i] + \alpha[i] \cdot w[i]$. The watermark strength $\alpha[i]$ of this watermark bit is determined by visual threshold as shown in Section 4.3. At the end, the watermarked $8 \times 8$ blocks of coefficients are converted back to the spatial domain by 2-D IDCT.

## 4.6 Watermark Detection Scheme

Figure 4.1 shows the block diagram of our watermark detection scheme. The selected coefficients of the original image are identified or pre-recorded. Since the locations of watermarked coefficients are image-dependent, if they are not pre-recorded they have to be found with the aid of the original image during watermark detection. The Robustness and the Reliability stages are the same as those at watermark embedding. So are the visual masking thresholds. Finally, watermark sequences are extracted from the received image and they are correlated with the original watermark for binary hypothesis testing and decision.



Fig. 4.1. Watermark detection scheme.

For a watermark sequence, the hypothesis decision is made based on the Bayes' decision rule. The mean value of the normalized correlation sum $C$ is computed by

$$C = \frac{1}{M}\sum_{i=1}^{M} c[i] = \frac{1}{M}\sum_{i=1}^{M} \frac{y[i] \cdot d[i]}{\sigma_d^2} \text{ with } d[i] = w[i] \cdot \alpha[i], \qquad (4.16)$$

where $y[i]$ is the difference between the DCT coefficients of the received image and the original image, $\alpha[i]$ is the watermark strength of coefficient $x[i]$, and $w[i]$ is the watermark signature. Then, $C$ is compared against the threshold $x_c$ derived from (4.7). The presence of the watermark is declared if $H_1$ is favored.

## 4.7 Simulation Results

We test the proposed watermarking scheme on two $256 \times 256$ images, Lena and Baboon. Due to the limited space, the experimental results listed below are mainly coming from the image Lena. Sets of robust coefficients are generated corresponding to JPEG compression quality factors ranging from 50 to 80 with a step of 10. Smaller quality numbers imply higher distortion. Two examples of the difference images between the original images and the watermarked images are shown in Figs. 4.2 (a)-(d). They are magnified by a factor of 15 so that we can see the differences visually. For the JPEG quality factor 50 in the design phase, the PSNR values between the original and the watermarked images are 46.1 dB and 42.9 dB for Lena and Baboon, respectively. And, they are 45.4 dB and 39.2 dB for JPEG quality factor 80 in the design phase. The watermark mainly spreads over the visually significant areas as we expect. Because the human visual model is used to control the watermark strength, subjectively we cannot see the distortion caused by watermarking. The Baboon image offers higher watermark capacity than Lena due to its highly textured content.

Some statistics of the selected coefficients corresponding to different JPEG compression quality factors after two processing stages are shown in Table 4.1 for image Lena. The number of dropped coefficients and the error detection probability for higher JPEG

compression quality factors are usually larger than those for lower JPEG compression quality factors. This is partially due to the fact that there are more candidate coefficients surviving JPEG compression with higher quality factors in the design phase. The estimated error probability of the selected robust coefficients is very small because only one binary decision is made on the entire image (and watermark) and the attacking source is assumed to be known in the design phase. If the number of coefficients is pre-selected to be, say, 200 and 1000, then the Reliability process stops only when the desired number is reached.

The estimated statistics for the selected coefficients after detection reliability improvement stage is shown in Table 4.2 for image Lena. They are calculated as follows. For a set of selected coefficients corresponding to a certain JPEG compression quality factor (e.g., 50) in the design phase, these statistics are estimated based on the unmarked and the watermarked images after JPEG compression at the same quality factor (e.g., 50). The experiment results show that the mean of the embedded watermark strength $E\{d\,|\,H_1\}$, the variance of the embedded watermark strength $Var\{d\,|\,H_1\}$ and the variances of normalized correlation sum $Var\{c\,|\,H_0\}$ and $Var\{c\,|\,H_1\}$ are larger for lower JPEG compression quality factors. This is due to the fact that the attack produced by a lower JPEG compression quality factor generates higher quantization distortion. Typically, the estimated $E\{d\,|\,H_1\}$ is around 0 in the design phase. The reason is that the number of the selected coefficients is large enough such that the numbers of the positive and the negative watermarks are about equal. We also notice that $Var\{c\,|\,H_0\}$ is not equal to $Var\{c\,|\,H_1\}$ for real images. The calculated $E\{c\,|\,H_0\}$ is not identical to zero but it is usually very close to zero. The value of $E\{c\,|\,H_1\}$ is distortion (attack) dependent and it is approximately 1.1 in JPEG case. Finally, the detection threshold $x_c$ computed from (4.7) is roughly near the average value of $E\{c\,|\,H_0\}$ and $E\{c\,|\,H_1\}$ as we expect.

StirMark 3.1[78] is used to test the robustness of our watermark. As shown in Figs. 4.3 (a) and (b) for images Lena and Baboon, respectively, our scheme can survive JPEG compression

at different quality factors. The data shown in Fig. 4.3 are each averaged over 5000 watermarked images with different pseudo random watermark sequences. Since the quantization step size decreases as the JPEG compression quality factor increases, a selected coefficient survives JPEG compression at higher quality factors may not survive JPEG compression at lower quality factors. This can be seen from Fig. 4.4(a). For example, in one experiment, among the selected coefficients in the design phase for JPEG quality 60 for image Lena, 75 percentage of coefficients can survive JPEG compression attack with quality 50, and 81 percentage of coefficients can survive JPEG compression attack with quality 70. On the other hand, coefficients designed for lower JPEG quality has a better surviving probability under higher quality attacks as one may expect.

Although our scheme is originally designed to be only JPEG-robust, we test its robustness against several other signal processing attacks including color reduction, Gaussian filtering, frequency-mode Laplacian removal (FMLR) and JPEG2000. The mean values of the correlation sum are shown in Fig. 4.5 and the data are obtained by averaging over 300 different pseudo random watermark sequences. The percentages of correctly detected images are shown in Fig. 4.6. The same detection thresholds designed for JPEG compression attacks are used in these experiments. As shown in Fig. 4.5, for the combined attacks of JPEG with either FMLR attacks or Gaussian filtering, the mean values of the correlation sum are often larger. Although the embedded watermark can survive most attacks, the combined attack of JPEG with $4 \times 4$ median filtering fails our scheme. The reason may be due to the strong lowpass filtering characteristics of the median filtering and a large percentage of our JPEG-robust watermark are embedded in the middle frequency band. In fact, the attacked images are so strongly lowpass filtered that the image quality degradation is visible. Finally, the simulation also shows that our watermark survives the JPEG2000 attack at bit rates 0.125 and 0.0625 bpp, and there is no detection failure in all cases.

To verify the designed false negative and positive error probabilities in the experiments, the mean, variance, minimum and maximum values of the normalized correlation sum $C$ are calculated for both the watermarked and unwatermarked image Lena. Again, 5000 watermarked Lena are generated with 5000 different pseudo random watermark sequences and the data shown in Figs. 4.3(a) and 4.7 are each averaged over these test data. To test the false positive (false alarm) case, 5000 different pseudo random watermark sequences are correlated with the unwatermarked but JPEG compressed image, and the results are averaged and shown in Figs. 4.8 and 4.9.

As shown in Fig. 4.7, for watermarked images, the variances of $C$ are all smaller than 0.0002 after JPEG compression attacks with different quality factors. The estimated (designed) and measured (simulated) average $E\{c \mid H_1\}$ are shown in Table 4.2 and their differences are shown in Table 4.3. The differences are very small and the measured $C$ values are fairly stable through out different sets of test watermarks. Our simulation matches the design target quite well. Similar false alarm analysis is conducted for the unwatermarked image cases. As shown in Fig. 4.8(b), the variances of $C$ in the design phase are smaller than 0.0005 after JPEG compression at different quality factors in the attack phase. The absolute differences between the estimated and measured $E\{c \mid H_0\}$ in Table 4.2 and Fig. 4.8(a) are shown in Table 4.3. Again, as we expect, the measured values of $C$ are close to the designed values and no failure cases occur in 5000 runs.

Regarding the computation complexity, the two stages in our algorithm behave quite differently. The second reliability improvement stage is the most time consuming. Given $N$ candidate coefficients, at the beginning of the iteration process, there are $N$ different sets of coefficients formed and their statistics have to be all estimated. Therefore, if there are $N$ selected coefficients after the first robustness stage, and there are totally $K$ coefficients dropped in the second reliability improvement stage, then the total number (times in calculation) of processed transform coefficients, $R$, is

$$R = N \cdot (N-1) + (N-1) \cdot (N-2) + (N-2) \cdot (N-3) + \ldots + (N-K) \cdot (N-K-1).$$

If $K \ll N$, the above expression indicates that $R$ is something around $\mathrm{O}(KN^2)$. The computation complexity can be very large, particularly, for an image with deep texture (and thus many coefficients are "robust"). As the statistics shown in Table 4.4, the value of $KN^2$ at quality 80 is about 6.6 times of that at quality 60 for image Lena and its associated computational complexity $R$ is about 4.1 times of that at quality 60. So, the true computational complexity is somewhat smaller than $\mathrm{O}(KN^2)$.

In the examples of pre-fixed target number of coefficients, Figs. 4.10 and 4.11 give the detection performance corresponding to 200, 1000 and 4019 selected coefficients for JPEG quality factor 50 in the design phase. Due to the space limit, the figures of correctly detected images are not shown. However, our simulation results show that there is neither missing nor false alarm cases for images with 5000 different pseudo random watermark sequences. The trade-off between the capacity and reliability is shown clearly here. For example, the error variance for 4019 coefficients is around $10^{-4}$ (Fig. 4.11(a)), the variance is increased to for 200 coefficients.

For comparison purpose, a middle-frequency DCT watermarking scheme [15] is also simulated. For a fair comparison with our proposed scheme, the watermark strength is chosen to be the visual masking threshold described in Section 4.3. The detection thresholds designed for JPEG attacks are still in use where the quality factor 50 is assumed. The range of middle frequency coefficients in [15] covers the frequency indices from 14 to 33 in a zigzag-scan manner. In our simulations, we narrow down the frequency range to 15 to 20 for stronger JPEG-robust coefficients. For a $256 \times 256$ image divided into $8 \times 8$ blocks, we randomly select the watermarking coefficients from the $1024 \times 6$ middle-frequency coefficients. A chosen coefficient is denoted by $x_{i,j}[m]$, where $i$ and $j$ are the horizontal and the vertical block coordinates, respectively, $0 \leq i, j \leq 63$, and $m$ is the frequency index and $15 \leq m \leq 20$.

We pick up one coefficient by randomly choosing the value of *i*, *j* and *m* independently is their specified ranges. Three experiments with 200, 1000 and 4019 DCT coefficients are chosen to match our designated coefficient capacity for the JPEG quality factor 50 in the design phase.

To calculate the false negative and positive error probabilities, the mean of the normalized correlation sum *C* and the percentage of the correctly detected images are computed for both the watermarked and unwatermarked image Lena. Again, for calculating the missing (false negative) probability, 5000 watermarked Lena are generated with 5000 different pseudo random watermark sequences and the results shown in Figs. 4.12 and 4.14(a) are each averaged over these test data. To test the false positive (false alarm) case, 5000 different pseudo random watermark sequences are correlated with the unwatermarked but JPEG compressed image, and the results are averaged and shown in Figs. 4.13 and 4.14(b). Fig. 4.12(a) indicates that the mean values of all the normalized correlation sums are greater than 0.9 for totally 4019 watermarked coefficients, and therefore, there is no failure case in Fig. 4.12(b). However, there are failure cases for fewer watermarked coefficients (e.g., 200). Also, the variance of the normalized correlation sum for the middle-frequency watermarking scheme (Fig. 4.14(a)) is about 40 times larger than that of our proposed scheme (Fig. 4.11(a)). Therefore, the error probability of the middle-frequency watermarking scheme is much higher. In the unwatermarked image case, Fig. 4.13(b), the false alarm cases appear for the sets with 200 marked coefficients. One may recall that neither false positive nor false negative cases occur in 5000 runs of our proposed scheme. Also, the variance of the normalized correlation sum for the middle-frequency watermarking scheme (Fig. 4.14(b)) is about 10 times larger than that of our scheme (Fig. 4.11(b)). It is clear that our scheme has much higher detection reliability.

As another example of attack, the JPEG2000 compression (two quality values) [79] is used in the design phase for picking up the DCT watermarking coefficients. The four

watermark patterns described in Section 4.2 are used in the training process. The training results for image Lena are shown in Tables 4.1 and 4.2. We again generate 5000 watermarked Lena embedded with 5000 pseudo random watermark sequences and the detection performance is shown in Fig. 4.15. Also, to test the false alarm case, 5000 different pseudo random watermark sequences are correlated with the unwatermarked but JPEG2000 compressed image, and the results are averaged and shown in Fig. 4.16. In both watermarked and unwatermarked cases, the variances are very small – an indication of very small missing and false alarm probability. Indeed in deciding the existence of correct watermark, there is neither a missing nor false alarm case. Note that the percentage of correctly decoded coefficients shown in Fig. 4.15(b) is around 80% because only 4 watermark patterns are used in the design phase. Comparing Figs. 4.6(a) and 4.15(b), we found that the JPEG-robust coefficients can usually survive the JPEG2000 compressions while the JPEG2000-robust coefficients are more sensitive to the JPEG attacks. Including more watermark patterns in the training process may help selecting the more robust coefficients.

# 4.8 Summary

In this chapter, a selection procedure is designed to identify the most effective DCT coefficients for watermarking purpose. The target is to explore the watermarking performance limit (capacity and detection reliability) of a give picture under a specified attack. There are essentially two steps in the design phase. Candidate coefficients and watermark signal (strength) are first chosen to achieve both robustness against the attack and perceptual invisibility. Then, we examine the error probability of using these candidate coefficients. The weak ones that lower the detection probability are discarded. Finally, we obtain a set of robust coefficients, which are both picture and attack dependent. Our simulations show that the proposed watermarking scheme performs very well in achieving high detection   probability

while maintaining the transparency of the embedded watermark. The methodology presented here for finding the most effective coefficients can be extended to the other types of attacks and/or watermarking techniques.

Table 4.1. The statistics of the selected coefficients at different JPEG and JPEG 2000 compression settings after the Robustness stage (stage 1) and the Reliability stage (stage 2) for image Lena.

| JPEG Quality Factor/ JPEG2000 Rate in Design Phase | No. of Selected Coefficients after Stage 1 | No. of Selected Coefficients after Stage 2 | Estimated $P_{error}$ after Stage 1 | Estimated $P_{error}$ after Stage 2 |
|---|---|---|---|---|
| JPEG 50 | 4738 | 200 | 4.772122e-143 | 8.703533e-017 |
| JPEG 50 | 4738 | 1000 | 4.772122e-143 | 4.092384e-076 |
| JPEG 50 | 4738 | 4019 | 4.772122e-143 | 0.000000e+000 |
| JPEG 60 | 6007 | 5082 | 5.394741e-188 | 0.000000e+000 |
| JPEG 70 | 8041 | 6587 | 1.364150e-246 | 0.000000e+000 |
| JPEG 80 | 11473 | 9439 | 2.307777e-263 | 0.000000e+000 |
| JPEG2000 0.0625bpp | 7120 | 5388 | 0.000000e+000 | 0.000000e+000 |
| JPEG2000 0.125bpp | 17656 | 13765 | 0.000000e+000 | 0.000000e+000 |

Table 4.2. The estimated statistics of the selected coefficients at different JPEG and JPEG2000 compression settings after the reliability improvement stage for image Lena.

| Design Phase | $E\{c \mid H_0\}$ | $Var\{c \mid H_0\}$ | $E\{c \mid H_1\}$ | $Var\{c \mid H_1\}$ | $E\{d \mid H_1\}$ | $Var\{d \mid H_1\}$ | Detection Threshold $x_c$ |
|---|---|---|---|---|---|---|---|
| JPEG 50 | 0 | 0.000446 | 1.654130 | 0.000560 | 0 | 18.282436 | 0.779981 |
| JPEG 60 | 0 | 0.000289 | 1.528975 | 0.000412 | 0 | 16.643265 | 0.696441 |
| JPEG 70 | 0 | 0.000141 | 1.357382 | 0.000244 | 0 | 14.089002 | 0.554339 |
| JPEG 80 | 0 | 0.000085 | 1.296038 | 0.000219 | 0 | 11.564191 | 0.496916 |
| JPEG2000 0.0625bpp | 0 | 0.000066 | 1.072749 | 0.000216 | 0 | 21.103317 | 0.381604 |
| JPEG2000 0.125bpp | 0 | 0.000015 | 1.035836 | 0.000087 | 0 | 12.421462 | 0.306017 |

Table 4.3. The absolute differences between the estimated and measured average $E\{c \mid H_1\}$ and $E\{c \mid H_0\}$, and the differences in ratio for watermarked and unwatermarked Lena images.

| JPEG Quality Factor in Design Phase | Watermarked Images ($H_1$) Absolute Difference Value | Unwatermarked Images ($H_0$) Absolute Difference Value |
|---|---|---|
| 50 | 0.00137 | 0.00046 |
| 60 | 0.00032 | 0.00024 |
| 70 | 0.00022 | 0.00014 |
| 80 | 0.00036 | 0.00009 |

Table 4.4. The processing complexity of the reliability improvement stage for images Lena and Baboon.

| JPEG Quality Factor in Design Phase | Lena | | | Baboon | | |
|---|---|---|---|---|---|---|
| | No. of Selected Coeff. after Stage 1 ($N$) | No. of Dropped Coeff. in Stage 2 ($K$) | Processed Coeff. | No. of Selected Coeff. after Stage 1 ($N$) | No. of Dropped Coeff. in Stage 2 ($K$) | Processed Coeff. |
| 50 | 4738 | 719 | 3152520 | 7359 | 1911 | 15897324 |
| 60 | 6007 | 925 | 5134207 | 11743 | 2771 | 28710990 |
| 70 | 8041 | 1454 | 10641870 | 15912 | 2807 | 40739868 |
| 80 | 11473 | 2034 | 21277960 | 22931 | 4046 | 84614676 |



| (a) | (b) | (c) | (d) |

Fig. 4.2. The (absolute) difference image between the original image and the watermarked image. The magnitude in display is amplified by a factor of 15: (a) Lena, Q=50. (b) Lena, Q=80. (c)Baboon, Q=50. (d) Baboon, Q=80.

(a)



(b)

Fig. 4.3. The mean of the normalized correlation sum $C_{\mathrm{E}\{c\,|\,\mathrm{H}_1\}}$ after the JPEG attacks (at four different quality factors) for watermarked images: (a) Lena and (b)Baboon.



(a)

(b)

Fig. 4.4. The percentage of correctly decoded coefficients at the detector after the JPEG attacks (at four different quality factors) for image Lena: (a) Watermarked and (b) Unwatermarked.



(a)



(b)

Fig. 4.5. The mean of the normalized correlation sum $E\{c \mid H_1\}$ under various signal processing attacks: (a) Lena and (b) Baboon.



(a)



(b)

Fig. 4.6. The percentage of correctly detected watermarks (images) under various signal processing attacks: (a)Lena and (b)Baboon.

Fig. 4.7. The variance of the normalized correlation sum $\mathrm{Var}\{c\,|\,\mathrm{H}_1\}$ after the JPEG attacks (at four different quality factors) for the watermarked image Lena.



(a)



(b)

Fig. 4.8. The mean and variance of the normalized correlation sum after the JPEG attacks (at

four different quality factors) for the unwatermarked image Lena: (a) Mean $E\{c\,|\,H_0\}$ and (b) Variance $Var\{c\,|\,H_0\}$.



(a)



(b)

Fig. 4.9. The maximum and minimum values of the normalized correlation sum $E\{c\,|\,H_0\}$ after the JPEG attacks (at four different quality factors) for the unwatermarked image Lena: (a) Maximum and (b) Minimum.

(a)



(b)

Fig. 4.10. The mean of the normalized correlation sum due to JPEG attacks for image Lena of designated capacity: (a) Watermarked and (b) Unwatermarked images.



(a)

(b)

Fig. 4.11. The variance of the normalized correlation sum due to JPEG attacks for image Lena of designated capacity: (a) Watermarked and (b) Unwatermarked images.



(a)



(b)

Fig. 4.12. The detection performance due to JPEG attacks for the middle-frequency watermarked image Lena: (a) The mean of the normalized correlation sum and (b) The percentage of correctly detected images.
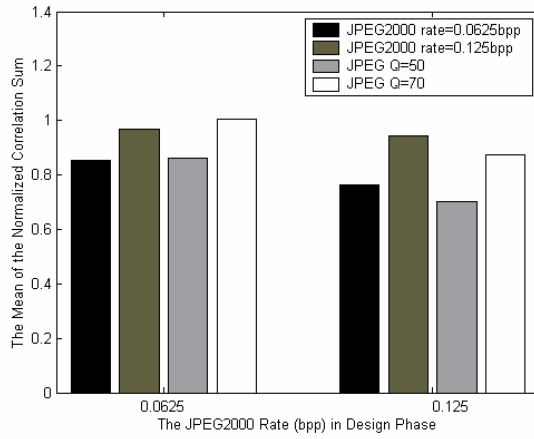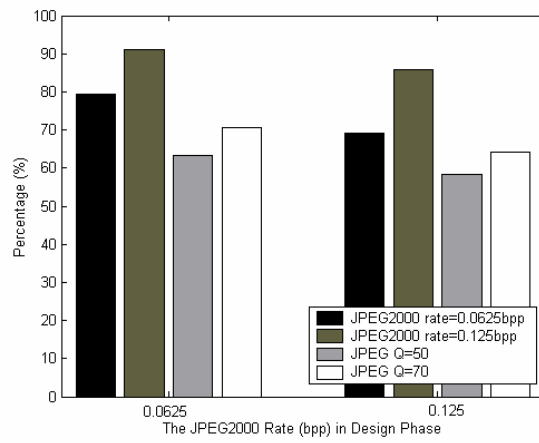


(a)

(b)

Fig. 4.13. The detection performance due to JPEG attacks for the middle-frequency unwatermarked image Lena: (a) The mean of the normalized correlation sum and (b) The percentage of correctly detected images.
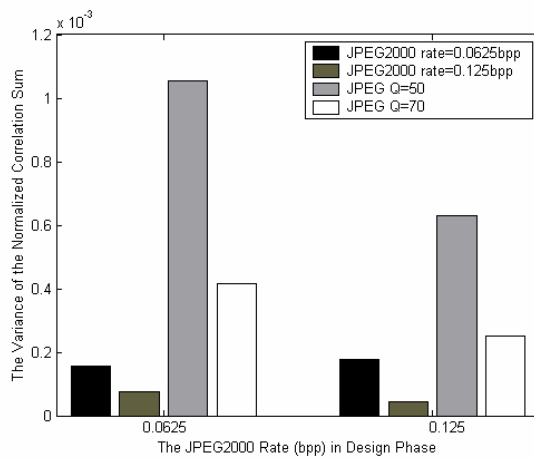
(a)



(b)

Fig. 4.14. The variance of the normalized correlation sum due to JPEG attacks for the middle-frequency watermarking for image Lena: (a) Watermarked images and (b) Unwatermarked images.
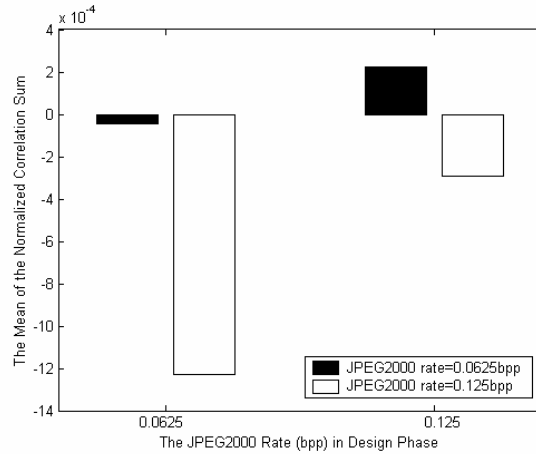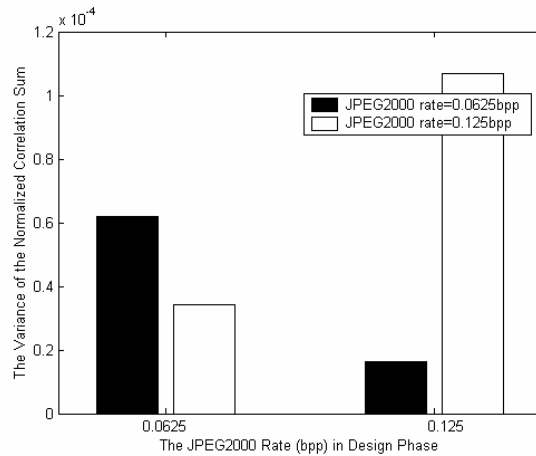
(a)



(b)



(c)

Fig. 4.15. The detection performance for JPEG2000-trained image Lena under JPEG and JPEG2000 attacks: (a) The mean of the normalized correlation sum, (b) The percentage of correctly decoded coefficients, and (c) The variance of the normalized correlation sum.



(a)



(b)

Fig. 4.16. The detection performance for unwatermarked image Lena under the JPEG2000 and JPEG attacks: (a) The mean of the normalized correlation sum, and (b) The variance of the normalized correlation sum.