

REVERSIBLE IMAGE AUTHENTICATION BASED ON WATERMARKING

Sang-Kwang Lee¹, Young-Ho Suh¹, and Yo-Sung Ho²

¹Electronics and Telecommunications Research Institute (ETRI)

161 Gajeong-Dong, Yuseung-Gu, Daejeon 305-700, KOREA

²Gwangju Institute of Science and Technology (GIST)

1 Oryong-Dong, Buk-Gu, Gwangju 500-712, KOREA

{sklee, syh}@etri.re.kr, hoyo@gist.ac.kr

ABSTRACT

In this paper, we propose a new reversible image authentication technique based on watermarking where if the image is authentic, the distortion due to embedding can be completely removed from the watermarked image after the hidden data has been extracted. This technique utilizes histogram characteristics of the difference image and modifies pixel values slightly to embed more data than other lossless data hiding algorithm. We show that the lower bound of the PSNR (peak-signal-to-noise-ratio) values of watermarked images are 51.14 dB. Moreover, the proposed scheme is quite simple and the execution time is rather short. Experimental results demonstrate that the proposed scheme can detect any modifications of the watermarked image.

1. INTRODUCTION

Digital representation of multimedia content offers various advantages, such as easy and wide distribution of multiple and perfect replications of the original contents. In many circumstances, alternations to content serve legitimate purposes. However, in other cases, the changes may be intentionally malicious or may inadvertently affect the interpretation of the content. For example, an inadvertent change to an X-ray image might result in a misdiagnosis, whereas malicious tampering of photographic evidence in a criminal trial can result in either a wrong conviction or acquittal. Thus, in applications for which we must be certain a content has not been altered, there is a need for verification or authentication of the integrity of the content [1].

In most authentication techniques based on watermarking, the original image is inevitably distorted due to the authentication itself. Typically, this distortion cannot be removed completely due to quantization, bit-replacement, or truncation at the grayscale 0 and 255. Although the distortion is often quite small, it may be unacceptable for medical or legal imagery or images with a high strategic importance in certain military applications. Thus, it is desired to undo the changes introduced by authentication if the image is verified as au-

thentic. Data embedding techniques satisfying this requirement, are referred to as *reversible* (also referred as *lossless*) image authentication techniques.

Ni *et al.* [2] proposed a lossless data embedding technique, which utilizes the zero or the minimum point of the image histogram. It can embed a large amount of data and the PSNR values of watermarked images are always higher than 48.13 dB. However, gray level values of the zero point and the peak point should be transmitted to the receiving side for data retrieval.

In this paper, we propose a new reversible authentication technique for images, which can embed a significant amount of data while keeping high visual quality. In order to verify the integrity of the image, we use a cryptographic hash function such as the MD5. The hash code is combined with a binary logo image by a bit-wise exclusive OR (XOR) and then embedded in the histogram of the difference image from the original image. In the worst case, a half the number of pixels of the image are added or subtracted by 1. Thus, the PSNR is guaranteed to be higher than 51.14 dB.

2. PROPOSED AUTHENTICATION SCHEME

2.1. Watermark Embedding

The watermark embedding procedure of the proposed authentication scheme is shown in Fig. 1.

For a grayscale image $I(i, j)$ of size $M \times N$ pixels, we form the difference image $D(i, j)$ of size $M \times \frac{N}{2}$ from the original image. For $0 \leq i \leq M - 1$ and $0 \leq j \leq \frac{N}{2} - 1$,

$$D(i, j) = I(i, 2j + 1) - I(i, 2j) \quad (1)$$

where $I(i, 2j + 1)$ and $I(i, 2j)$ are the odd-line field and the even-line field, respectively. For watermark embedding, we empty the histogram bins of -2 and 2 by shifting some pixel values in the difference image. If the difference value is greater than or equal to 2, we add one to the odd-line pixel. If the difference value is less than or equal to -2, we subtract one from the the odd-line pixel. Then, the modified difference

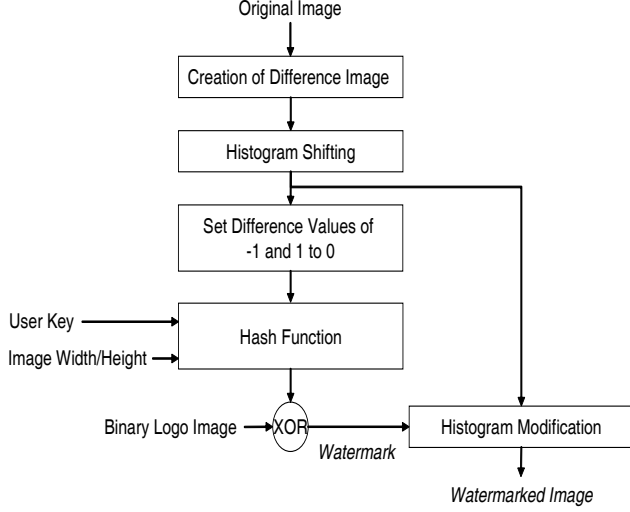


Fig. 1. Watermark embedding

image $\tilde{D}(i, j)$ can be represented as

$$\tilde{D}(i, j) = \tilde{I}(i, 2j + 1) - I(i, 2j) \quad (2)$$

where

$$\tilde{I}(i, 2j + 1) = \begin{cases} I(i, 2j + 1) + 1 & \text{if } D(i, j) \geq 2 \\ I(i, 2j + 1) - 1 & \text{if } D(i, j) \leq -2 \\ I(i, 2j + 1) & \text{otherwise} \end{cases} \quad (3)$$

In this paper, we use MD5 as a hash function. The MD5 algorithm takes as input a message of arbitrary length and produces as output a bit array of 128. It is conjectured that it is computationally infeasible to produce two messages having the same message digest, or to produce any message having a given prespecified target message digest. Let

$$\begin{aligned} \mathbf{D}_s &= \{D_s(i, j) : 0 \leq i \leq M - 1; 0 \leq j \leq \frac{N}{2} - 1\} \\ \tilde{\mathbf{D}}_s &= \{\tilde{D}_s(i, j) : 0 \leq i \leq M - 1; 0 \leq j \leq \frac{N}{2} - 1\} \end{aligned} \quad (4)$$

Then, we compute the hash $A(l)$ of length 128 as follows:

$$H(K, M, N, \tilde{\mathbf{D}}_s) = \{A(l) : 0 \leq l \leq 127\} \quad (5)$$

where K is a user key consisting of a sting of bits and $\tilde{\mathbf{D}}_s$ equals the corresponding element in \mathbf{D}_s except difference values of -1 and 1 are set to zero.

In order to generate a watermark $W(m, n)$ of the same size as embedding capacity C , we combine the hash $A(l)$ with the binary logo image $B(m, n)$ of size $P \times Q$ pixels using the bit-wise XOR operation. In the case that $C > 128$ and $C > P \times Q$, we form $\tilde{A}(l)$ and $\tilde{B}(m, n)$ by periodically replicating $A(l)$ and $B(m, n)$ to the length C , respectively. Then, the authentication watermark is computed as

$$W(m, n) = \tilde{A}(l) \oplus \tilde{B}(m, n) \quad (6)$$

where \oplus is the bit-wise XOR operation.

In the histogram modification process, the watermark is embedded into the modified difference image. The modified difference image is scanned. Once a pixel with the difference value of -1 or 1 is encountered, we check the watermark to be embedded. If the bit to be embedded is 1, we move the difference value of -1 to -2 by subtracting one from the odd-line pixel or 1 to 2 by adding one to the odd-line pixel. If the bit to be embedded is 0, we skip the pixel of the difference image until a pixel with the difference value -1 or 1 is encountered. In this case, there is no change in the histogram. Therefore, the watermarked old-line field $I_w(i, 2j + 1)$ is obtained as follows: If $W(m, n) = 1$ and $\tilde{D}(i, j) = 1$ or -1 ,

$$I_w(i, 2j + 1) = \begin{cases} \tilde{I}(i, 2j + 1) + 1 & \text{if } \tilde{D}(i, j) = 1 \\ \tilde{I}(i, 2j + 1) - 1 & \text{if } \tilde{D}(i, j) = -1 \end{cases} \quad (7)$$

and in all other cases,

$$I_w(i, 2j + 1) = \tilde{I}(i, 2j + 1) \quad (8)$$

and the watermarked even-line fields $I_w(i, 2j)$ is given by

$$I_w(i, 2j) = I(i, 2j) \quad (9)$$

The embedding capacity of this scheme equals to the number of pixels with the difference values of -1 and 1 in the difference image. A large number of pixel values of the difference image have a tendency to be distributed around 0. Using this property of the difference image, we can embed a large amount of data as compared to the original image itself.

Assume that there is no pixel with overflow and underflow in the original image. In the worst case, all pixels of the odd-line field will be added or subtracted by 1. The mean squared error of this case is 0.5. Hence, the PSNR of the watermarked image can be calculated as

$$\text{PSNR(dB)} = 10 \log_{10}(255^2 \cdot 2) \approx 51.14 \quad (10)$$

In short, the lower bound of the PSNR of the watermarked image is about 51.14 dB. This result is much higher than other lossless data hiding techniques.

2.2. Watermark Extraction and Recovery

Fig. 2 depicts the watermark extraction and recovery scheme. In this process, we verify that the image has not been tampered with and if the image is authentic, we reverse the watermarked image back to the original image.

For the received watermarked image $I_e(i, j)$ of size $M_e \times N_e$ pixels, we calculate the difference image $D_e(i, j)$ of size $M_e \times \frac{N_e}{2}$ pixels. The whole difference image is scanned. The authentication watermark $W_e(m, n)$ is extracted as follows:

$$W_e(m, n) = \begin{cases} 0 & \text{if } D_e(i, j) = -1 \text{ or } 1 \\ 1 & \text{if } D_e(i, j) = -2 \text{ or } 2 \end{cases} \quad (11)$$

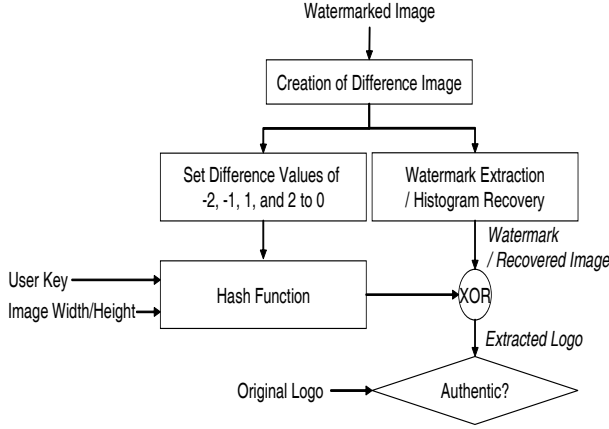


Fig. 2. Watermark extraction and recovery

Simultaneously, we reverse the watermarked image back to the original image by shifting some pixel values in the difference image. The whole difference image is scanned once again. The recovered odd-line field $I_r(i, 2j + 1)$ can be expressed as

$$I_r(i, 2j + 1) = \begin{cases} I_e(i, 2j + 1) - 1 & \text{if } D_e(i, j) \geq 2 \\ I_e(i, 2j + 1) + 1 & \text{if } D_e(i, j) \leq -2 \\ I_e(i, 2j + 1) & \text{otherwise} \end{cases} \quad (12)$$

Since we manipulate pixel values of only the odd-line field in the watermark embedding process, the recovered even-line field $I_r(i, 2j)$ is obtained by

$$I_r(i, 2j) = I_e(i, 2j) \quad (13)$$

Hence, the corresponding difference image $D_r(i, j)$ is calculated as

$$D_r(i, j) = I_r(i, 2j + 1) - I_r(i, 2j) \quad (14)$$

Let

$$\begin{aligned} \mathbf{D}_r &= \{D_r(i, j) : 0 \leq i \leq M - 1; 0 \leq j \leq \frac{N}{2} - 1\} \\ \tilde{\mathbf{D}}_r &= \{\tilde{D}_r(i, j) : 0 \leq i \leq M - 1; 0 \leq j \leq \frac{N}{2} - 1\} \end{aligned} \quad (15)$$

Then, we compute the hash $\hat{A}(k)$ of length 128 as

$$H(K, M_e, N_e, \tilde{\mathbf{D}}_r) = \{\tilde{A}_e(k) : 0 \leq k \leq 127\} \quad (16)$$

where $\tilde{\mathbf{D}}_r$ equals the corresponding element in \mathbf{D}_r except difference values of -2, -1, 1 and 2 are set to zero.

In order to extract the embedded binary logo image of the image, we combine the hash $A_e(k)$ with the extracted watermark $W_e(m, n)$ of length C using the bit-wise XOR operation. In the case that $C > 128$, we form $\tilde{A}_e(k)$ by periodically replicating $A_e(k)$ to the length C . Then, the embedded binary logo image $\tilde{B}_e(m, n)$ of the length C is computed as

$$\tilde{B}_e(m, n) = \tilde{A}_e(k) \oplus W_e(m, n) \quad (17)$$

In the case that $C > P \times Q$, $\tilde{B}_e(m, n)$ is the replication version of $B_e(m, n)$.

Finally, we form $\hat{B}(m, n)$ by periodically replicating $B(m, n)$ to the length C and verify the integrity of the image by comparing $\hat{B}_e(m, n)$ with $\hat{B}(m, n)$.

3. EXPERIMENTAL RESULTS AND ANALYSIS

In order to evaluate the performance of the proposed authentication scheme, we perform computer simulations on several test images of size 512×512 pixels. Fig. 3 shows a watermark which is a binary logo image of size 128×26 pixels, equivalent to a binary sequence of 3,328 bits.



Fig. 3. Binary logo image of 128×26 pixels

Fig. 4 shows the watermarked test images with the authentication watermark added by the proposed scheme. It is seen from the figure that since we modify pixel values slightly to embed the authentication watermark, there is no visible degradation due to embedding in the watermarked image.

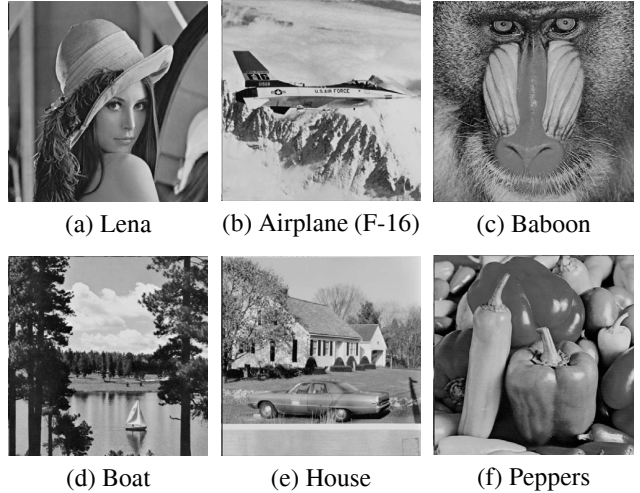


Fig. 4. Watermarked images

Table 1 summarizes the experimental results. This table shows that the PSNR values of all watermarked images are above 51.14 dB, as we theoretically proved in Section 2.1. The capacity ranges from 8,150 bits to 28,094 bits for $512 \times 512 \times 8$ test grayscale images. This result shows that the proposed scheme offers adequate capacity to address most applications. It is also seen from Table 1 that an image like Baboon, which contains significant texture, has considerably lower capacity than simple images such as Airplane(F-16).

Table 1. PSNR and capacity of the proposed scheme

Test images	PSNR (dB)	Capacity (bits)
Lena	52.21	26,990
Airplane (F-16)	52.35	28,094
Baboon	51.41	8,150
Boat	51.64	14,503
House	52.10	23,126
Peppers	51.43	17,143

In our experiments, we have considered several attacks of pixel value changes in the watermarked image. If the watermarked image is not changed, and if the correct key is used in the watermark extraction process, the appropriate binary logo image is extracted, as shown in Fig. 5(a). Note that the length of the extracted binary logo image is equal to the embedding capacity of the Lena image. In this case, we can reverse the watermarked image back to the original image. If the image is not watermarked or if the pixel values of the watermarked image are changed, the extracted binary logo image appears like random noise. If the inappropriate key is used in the extraction process, the embedded binary logo image is periodically deformed with period one line, as shown in Fig. 5(b). This is because of the properties of the MD5, that is, any string of data is hashed into a bit array of length 128.

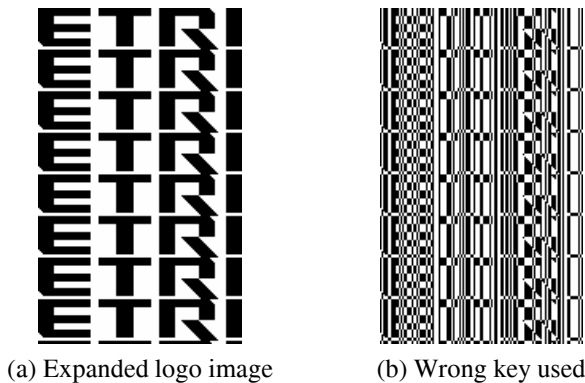


Fig. 5. Extracted logo images on Lena

Fig. 6 shows the PSNR versus capacity comparison on Lena among the proposed and other methods. Comparing the proposed method with compression-based methods [3, 5, 6], the gain of PSNR is at least 5 dB. This is because the methods rely on some form of compression to create space for embedding the payload. On the other hand, it can be seen that the payload size of the proposed method is relatively smaller than other methods because there is a trade-off between capacity and distortion.

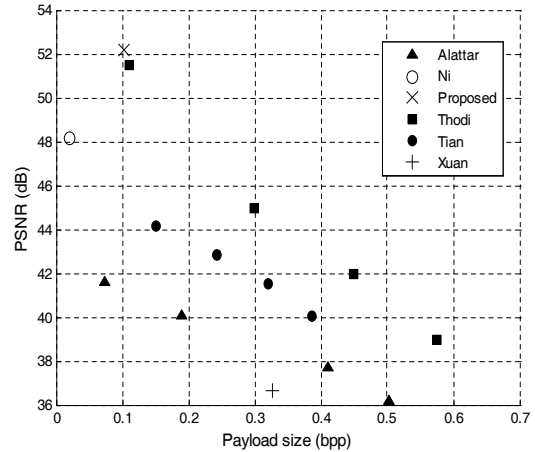


Fig. 6. PSNR vs. capacity comparison on Lena

4. CONCLUSIONS

In this paper, a new watermarking-based image authentication technique is presented. The technique provides low distortion, high embedding capacity, and lossless data hiding. Experimental results show that the proposed authentication scheme can detect if the key is incorrect, if the image is not watermarked, or if the image is changed in its pixel values. That is, we can verify the integrity of the image. Finally, we note that many authentication schemes based on watermarking have the ability to identify regions of the image that have been corrupted. We are currently studying this localization capability for our proposed scheme.

5. REFERENCES

- [1] I. Cox, M. Miller, and J. Bloom, *Digital Watermarking*, Morgan Kaufmann Publishers, 2001.
- [2] Z. Ni, Y. Shi, N. Ansari, and W. Su, "Reversible data hiding," *Proc. ISCAS 2003*, vol. 2, pp. 912–915, 2003.
- [3] A. Alattar, "Reversible watermark using difference expansion of triplets," *Proc. ICIP*, pp. 501–504, 2003.
- [4] D. Thodi and J. Rodríguez, "Reversible watermarking by prediction-error expansion," *6th IEEE Southwest Symposium on Image Analysis and Interpretation*, pp. 21–25, 2004.
- [5] J. Tian, "High capacity reversible data embedding and content authentication," *Proc. ICASSP*, pp. III-517–520, 2003.
- [6] G. Xuan, J. Zhu, J. Chen, Y. Shi, Z. Ni, and W. Su, "Distortionless data hiding based on interger wavelet transform," *IEE Electronics Letters*, pp. 1646–1648, 2002.