ELSEVIER

# Invertibility attack against watermarking based on forged algorithm and a countermeasure

Xinpeng Zhang *, Shuozhong Wang

*School of Communication and Information Engineering, Shanghai University, 149 Yanchang Road, Shanghai 200072, China*

## Abstract

It is shown in this paper that, even with a non-invertible watermarking algorithm or an asymmetric watermarking protocol, it is still possible to effect an invertibility attack, which relies on a forged watermarking algorithm, a counterfeit mark, and a fake key. Two examples are given to show the vulnerability of the unfortified non-invertible algorithm/asymmetric protocol. As a solution, a secure watermarking protocol is proposed, which establishes correlation between the watermarking algorithm and the embedded mark.
© 2004 Elsevier B.V. All rights reserved.

*Keywords:* Digital watermarking; Protocol; Invertibility attack

## 1. Introduction

As an effective means for the protection of intellectual property rights, watermarks can be embedded imperceptibly into multimedia data, typically containing information about the owner, origin, status, and destination of the product (Petitcolas et al., 1999; Hartung and Kutter, 1999). In the meantime, hostile attacks try to invalidate watermarking systems for illegal profits. The attacks can generally be classified into four categories: removal attack, geometric distortion attack, cryptographic attack, and protocol attack (Voloshynovskiy et al., 2001).

Protocol attacks do not aim at destroying the embedded information or disabling the synchronization. The goal is to demolish the very concept of the watermarking application itself. For example, copy attack (Kutter et al., 2000) is based on estimation of the embedded watermark in the spatial domain through a filtering process without knowing the watermarking algorithm and the key. The estimate is then used and inserted into the target data. In this way, the watermark is copied from a watermarked product into another.

Craver et al. (1997, 1998) devised another protocol attack, invertibility attack. It creates counterfeit watermarks in the same digital product to cause confusion therefore prevent rightful assertion of ownership. In this regard, it has been suggested

---
* Corresponding author. Tel.: +86-2156331435; fax: +86-2156331964.
E-mail addresses: zhangxinpeng@263.net (X. Zhang), shuo-wang@yc.shu.edu.cn (S. Wang).

(Swanson et al., 1998; Hwang et al., 1999) that non-invertible watermarking algorithms be used as they cannot be abused as tools in invertibility attacks. An alternative countermeasure is to design a special watermarking protocol in conjunction with an invertible algorithm to defeat the attack (Katzenbeisser and Veith, 2002).

In this paper, however, we show that loopholes still exist even though a non-invertible algorithm or protocol is used. Having obtained any particular watermarked multimedia data, an attacker can always tailor-make a non-invertible inserter and a corresponding detector, as well as a fake original and a workable key. Since the watermarking algorithm constructed as such conforms to the non-invertibility requirements, the attacker can claim ownership of the data. A novel watermarking protocol is then introduced to provide a remedy.

## 2. Invertibility attack and the forged-algorithm attack

An invertibility attack is described as follows (Craver et al., 1997, 1998). Suppose that the copyright owner A embeds his watermark $W_A$ into the host medium $I$ with an inserter $E_A$. A watermarked medium $I_A$ is thus produced that is available to both the legitimate user and a hacker named B. Assume that B possesses a watermarking tool consisting of an inserter $E_B$ and a detector $D_B$. If B is able to construct a watermark $W_B$ as well as a fake data $I'$ from $I_A$ such that (1) $E_B(I', W_B) = I_A$, (2) $D_B(I_A, W_B) = 1$, (here 1 means "watermarked", and 0 means "not watermarked"), and (3) $I'$ is similar to $I_A$, then the watermarking system $(E_B, D_B)$ is said to be invertible. According to the above operation, both the authentic watermark $W_A$ and the counterfeit mark $W_B$ can be detected from $I_A$ using $D_A$ and $D_B$, respectively. Therefore, the rightful owner of $I_A$ cannot be identified based on the extracted watermark.

In addition, a quasi-invertible watermarking tool can also be designed. With this tool the attacker can find $W_B$, $I'$, and $I_B$ so that (1) $E_B(I', W_B) = I_B$, (2) $D_B(I_B, W_B) = 1$, $D_B(I_A, W_B) = 1$, and (3) both $I'$ and $I_B$ are similar to $I_A$. Quasi-invertibility is less stringent than full-invertibility,

therefore more harmful to legitimate watermarking system users. In the following discussion, quasi-invertibility and full-invertibility are not distinguished and will simply be referred to as "invertibility".

To defeat invertibility attacks, Craver et al. suggested that non-invertible watermarking techniques be used for copyright protection. In other words, all legitimate watermarking algorithms must be designed as non-invertible. If a non-invertible watermarking system $(E_B, D_B)$ and a marked copy $I_A$ are given, the attacker cannot find a fake watermark $W_B$ and a fake original copy $I'$ satisfying all the three conditions listed in either of the previous two paragraphs.

In fact, however, after obtaining the marked copy $I_A$, an attacker can always construct a non-invertible watermarking system $(E_B, D_B)$ as well as $W_B$ and $I'$ satisfying the above-mentioned conditions. In this way, the attacker can publish the mechanism of his illicit watermarking system $(E_B, D_B)$ to show its non-invertibility, therefore declare the legitimacy of the fake mark $W_B$ in $I_A$. Note that the illicit watermarking system $(E_B, D_B)$ is tailor-made for the particular object $I_A$ and cannot be used to perform attack against any other media because of its non-invertibility. Nonetheless the copyright of $I_A$ is breached, and a second forged system can be found against another watermarked object. Such an object-specific attack is termed forged-algorithm attack in this paper.

Katzenbeisser and Veith proposed an alternative watermarking system or protocol for resisting invertibility attack. The method can be used in combination with any watermarking algorithm regardless of being non-invertible or not. It will be shown in the next section, however, that Katzenbeisser's method is also vulnerable to the forged-algorithm attack.

## 3. Scenarios of possible attack against existing non-invertible schemes

### 3.1. Attack against non-invertible algorithms

In this subsection, we show how a forged algorithm can be obtained to perform a successful

attack against a non-invertible watermarking system.

Without loss of generality, assume that the marked copy $I_A$ is a still image and the number of pixel is $N$. Randomly generate a binary sequence $R_B$ of length $N$ which consists of +1 and −1, and then subtract $R_B$ from $I_A$ to produce a fake original $I'$:

$$I' = I_A - R_B \tag{1}$$

where "−" is a subtraction operation between two vectors.

Define a one-way hash function $h$, which generates a binary sequence of length $N$. Arbitrarily select a key $k_B$ and compute

$$W_B = h(I', k_B) \tag{2}$$

In both $R_B$ and $W_B$, about $N/2$ elements are +1, and the rest are −1. So, one can find a permutation $\Pi$ for re-organizing elements of $W_B$ to make the permuted result very similar to $R_B$. For example, assuming that $W_B = [+1, +1, -1, -1, +1, -1, +1, -1, -1, +1, +1, -1, -1, -1, +1, +1, \ldots]$, and $R_B = [-1, +1, +1, -1, +1, -1, -1, +1, -1, -1, +1, +1, +1, -1, -1, +1, \ldots]$, the attacker may define a permutation $\Pi$: $[1 \rightarrow 3, 2 \rightarrow 2, 3 \rightarrow 4, 4 \rightarrow 1, 5 \rightarrow 8, 6 \rightarrow 7, 7 \rightarrow 5, 8 \rightarrow 6, 9 \rightarrow 10, 10 \rightarrow 11, 11 \rightarrow 12, 12 \rightarrow 9, 13 \rightarrow 15, 14 \rightarrow 14, 15 \rightarrow 16, 16 \rightarrow 12, \ldots]$ to make $\Pi(W_B) \approx R_B$, where $\alpha \rightarrow \beta$ means to place the $\alpha$th element on the $\beta$th position. Thus an illegal non-invertible watermarking system can be counterfeited using $\Pi$ and $h$ as follows.

The embedding procedure $E_B$:

1. Assuming that the original image is $I$, a secret key $k$ can be selected to compute a watermark by using the one-way hash function $h$:

$$W = h(I, k) \tag{3}$$

2. Permute $W$ with $\Pi$ to yield $W_1$, and add it into the original $I$ to produce a marked image.

Detecting procedure $D_B$:

1. If anyone announces that the rightful owner is himself, he must provide his original image and the key to compute the permuted mark $W_1$. The difference between marked and original images, $R$, can also bef obtained.
2. Computing the correlation coefficient $\rho$ between $W_1$ and $R$:

$$\rho = \frac{1}{N} \sum_{i=1}^{N} [R(i) \cdot W_1(i)] \tag{4}$$

When $\rho$ is greater than 1/2 by a specified margin, a "watermarked" decision is made, otherwise "not watermarked".

In this way, the attacker can provide his fake original $I'$ and a fake key $k_B$ to claim copyright. Because $W_1$, the permuted version of $W_B$, is very similar to $R_B$, the value of $\rho$ must be close to 1.

Fig. 1 sketches the forged embedding algorithm. It should be noted that the illegal watermarking system ($E_B$, $D_B$) is not invertible because a one-way hash function $h$ is involved. When a watermarking system, including the one-way hash function $h$ and the permutation $\Pi$, is given, it is virtually impossible to fabricate a fake original and a usable key. In the above procedure, the attacker not only "creates" the original image and the key, but also a permutation $\Pi$. In other words, he has successfully forged a workable watermarking algorithm. Although the forged watermarking system ($E_B$, $D_B$) cannot be used to attack any other watermarked image, the copyright of $I_A$ is breached. To
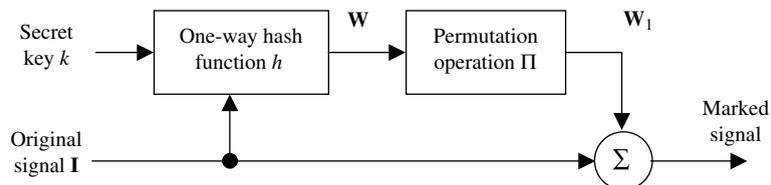


Fig. 1. Sketch of forged embedding algorithm.

attack another medium, the attacker must find a new permutation and construct another watermarking algorithm.

### 3.2. Forged-algorithm attack against the asymmetric watermarking protocol

In (Katzenbeisser and Veith, 2002), Katzenbeisser and Veith constructed an asymmetric watermarking system using a cryptographic signature on top of a traditional watermarking algorithm $E$, as sketched in Fig. 2. The purpose was to defeat the invertibility attack even if $E$ is invertible, since the asymmetry is not from $E$, but from the non-invertible signature mechanism that includes a public key $P$ and a private key $S$. The copyright owner A uses the original object $I$, and a watermark $W_A$, to derive $W'_A = W_A \| S_S(I, P) \| S_S(W_A, P)$, where the operator $\|$ concatenates two strings and $S_S$ is a signature generator. $W'_A$ is then embedded into $I$ using an embedding algorithm $I_A = E(I, W'_A, k_A)$, where $k_A$ is a key. In detection, given a claimed watermarked copy $I_A$, the corresponding original $I$, a watermark $W$, and a key $k$, $W'_A$ is computed and split into three concatenated parts: $W' = W'_1 \| W'_2 \| W'_3$. Check whether $D(I_A, I, W', k) = \text{TRUE}$, $V_S(I, W'_2, P) = \text{TRUE}$, and $V_S(W'_1, W'_3, P) = \text{TRUE}$, where $D$ is the detecting algorithm corresponding to $E$, and $V_S$ the signature verification process. Only if all the three are true, the watermark is accepted as genuine. If the attack's watermarking algorithm $(E_B, D_B)$ is given, it is virtually impossible to find a fake original $I'$, a fake mark $W_B$ and a fake key $k_B$ satis-fying all the three tests, even though the watermarking algorithm may be invertible.

However, the system in Fig. 2 is also susceptible to the forged-algorithm attack, which fabricates an illegal watermarking algorithm $(E_B, D_B)$ together with a fake original, a derived watermark and a key. This is shown in the following where $(E_B, D_B)$ is based on a dither modulation technique (Chen and Wornell, 1999).

Assuming that the length of $W'$ is always $M$, convert $I_A$ into $Y_A$ using an orthogonal transform such as DCT or DFT, and select a quantization step $\Delta$. The purpose of orthogonal transform is to spread the mark information $W'$ over the entire space of host object. Then, select $M$ transform coefficients in $Y_A$, denoted $y_1, y_2, \ldots, y_M$. The dither coefficients are derived:

$$y'_k = y_k + e_k \quad k = 1, 2, \ldots, M \tag{5}$$

where $e_k$ is randomly and independently picked from $[-\Delta/2, \Delta/2]$. The coefficients $y_1, y_2, \ldots,$ and $y_M$ are replaced with $y'_1, y'_2, \ldots,$ and $y'_M$ and inversely transformed to produce an object $I'$ as a fake original copy. The fake key $k_B$ is made up of $\Delta$ and the way of selecting $M$ coefficients from $Y_A$.

Select a fake $W_B$ to compute $W'_B$ according to the watermarking protocol. A vector of length $M$ is obtained

$$d_k = \text{round}\left(\frac{y_k}{\Delta}\right) \cdot \Delta - y_k \quad k = 1, 2, \ldots, M \tag{6}$$

where round $(\cdot)$ takes the nearest integer of its argument. The values of $d_k$ also lie within $[-\Delta/2, \Delta/2]$. A new dither vector $[D_1, D_2, \ldots, D_M]^T$ is used as the parameter of watermarking algorithm, where
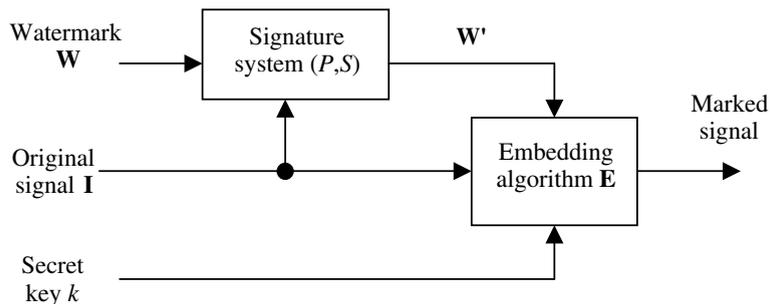


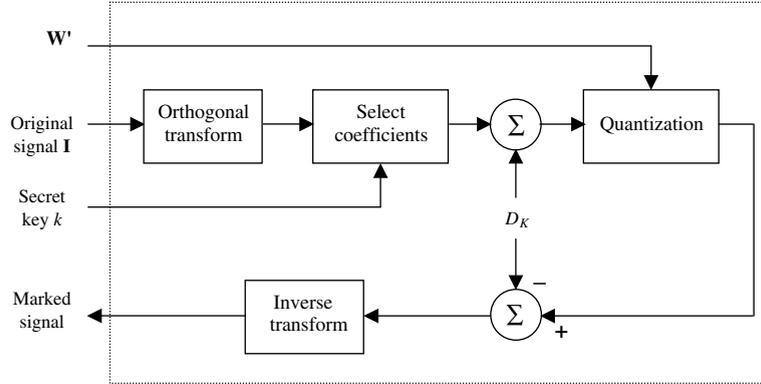Fig. 2. Sketch of asymmetric watermarking protocol.

Fig. 3. Forged embedding algorithm $E_\mathrm{B}$ in the asymmetric watermarking protocol.

$$D_k = \begin{cases} d_k, & \text{if} \quad W'_{\mathrm{B},k} = 0 \\ d_k - \Delta/2, & \text{if} \quad W'_{\mathrm{B},k} = 1, d_k \geqslant 0 \\ \Delta/2 + d_k, & \text{if} \quad W'_{\mathrm{B},k} = 1, d_k < 0 \end{cases}$$
$$k = 1, 2, \ldots, M \tag{7}$$

Thus, the attacker is able to construct its illegal watermarking algorithm as described in the following. Note that the original is not needed in detection.

The embedding procedure $E_\mathrm{B}$ (see Fig. 3):

1. Convert the original data using an orthogonal transform and select $M$ coefficients $y'_1, y'_2, \ldots, y'_M$ according to the key.
2. Embed $W' = W \| S_\mathrm{S}(I, P) \| S_\mathrm{S}(W, P)$ using the dither modulation method

$$y_k = \begin{cases} Q_0(y'_k + D_k) - D_k, & \text{if} \quad W'_k = 0 \\ Q_1(y'_k + D_k) - D_k, & \text{if} \quad W'_k = 1 \end{cases}$$
$$k = 1, 2, \ldots, M \tag{8}$$

where

$$Q_0(x) = \text{round}\left(\frac{x}{\Delta}\right) \cdot \Delta \tag{9}$$

$$Q_1(x) = \text{round}\left(\frac{x + \Delta/2}{\Delta}\right) \cdot \Delta - \frac{\Delta}{2} \tag{10}$$

3. Perform the inverse transform to produce a marked copy.

The detecting procedure $D_\mathrm{B}$:

1. Transform the marked copy and pick $M$ coefficients $y_1, y_2, \ldots, y_M$ according to the key.
2. Compute elements of the vector $V$:

$$V_k = 1 - \frac{1}{2}\left[(-1)^{\text{round}(\frac{y_k + D_k}{\Delta/2})} + 1\right]$$
$$k = 1, 2, \ldots, M \tag{11}$$

3. Compare $V$ with $W'$. If they are identical, the output is TRUE, otherwise FALSE.

When the attacker publishes his fake watermarking algorithm and announces his mark $W_\mathrm{B}$ and the key $k_\mathrm{B}$, $V$ and $W'_\mathrm{B}$ should be identical according to (7) and (11). Eq. (5) indicates that $e_k$ is the difference between the fake "original" $y'_k$ and the marked coefficient $y_k$, which obeys a uniform distribution within $[-\Delta/2, \Delta/2]$, similar to the distribution of quantization error due to embedding a legal watermark using $E_\mathrm{B}$. This way, the asymmetric watermarking system is broken under the forged-algorithm attack. As in the previous example, the attack is aimed at a specific object. Another fake watermarking algorithm must be constructed in an attack against a new object.

## 4. A secure watermarking protocol capable of resisting forged-algorithm attacks

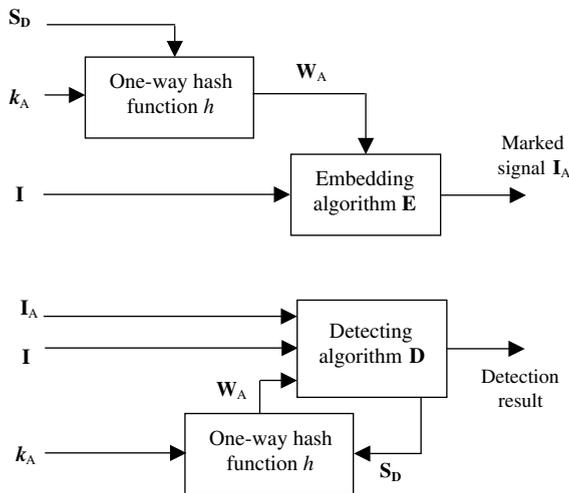As mentioned in Section 3, an attacker can arbitrarily construct a non-invertible watermarking

Fig. 4. The proposed secure watermarking protocol.

algorithm, or an invertible algorithm under an asymmetric protocol, to find a fake watermark in any media to falsely claim his "ownership". In order to defeat this kind of attack, correlation between the embedded signal and the watermarking algorithm should be established.

As the structure of the watermarking algorithm must be published according to Kerckhoff's principle, and the detecting software $D$ can be considered and treated as a binary sequence $S_D$, we propose a secure watermarking protocol as follows (see Fig. 4):

1. The protocol employs a one-way hash function $h$ and publicizes it.
2. The copyright owner A selects his key $k_A$ and compute $W_A = h(S_D, k_A)$ to be embedded into the original copy, $I_A = E(I, W_A) = E[I, h(S_D, k_A)]$.
3. If one claims rightful ownership of a multimedia object, he must provide his original copy $I$, the key $k$ and the detecting software $D$. A "watermarked" decision is made when $D[I_A, I, h(S_D, k_A)] = \text{TRUE}$, otherwise "not watermarked".

In the proposed protocol, the embedded data is a code derived from the detecting software $D$ and a secret key, but not the detecting software itself.

When an attacker constructs his counterfeit watermarking algorithm, he cannot obtain a suitable mark $W$ resulting from $h$, therefore is unable to make $D[I_A, I, h(S_D, k)] = \text{TRUE}$. Thus the forged-algorithm attack can no longer succeed.

## 5. Conclusion

In this paper, we first describe a type of invertibility attack, which can be performed when a non-invertible watermarking algorithm or an asymmetric protocol is used. If a watermarking algorithm is properly constructed, any watermark signal needed for ownership authentication can be derived from a marked copy so that the copyright is confused. In order to defeat this attack, a secure watermarking protocol is proposed, in which correlation between the detecting algorithm and the embedded mark is imperative. Because the attacker cannot generate the algorithm and the required watermark simultaneously, the forged-algorithm attack is defeated.

## Acknowledgement

## References

Chen, B., Wornell, G.W., Dither modulation: a new approach to digital watermarking and information embedding. In: Security and Watermarking of Multimedia Contents. Proc. SPIE, vol. 3657, San Jose, CA, January 1999, pp. 342–353.

Craver, S., Menon, N., Yeo, B.L., Yeung, M., 1997. On the invertibility of invisible watermarking techniques. Proc. IEEE Internat. Conf. on Images Process., Chicago, Illinois, USA, October 1997, pp. 540–543.

Craver, S., Menon, N., Yeo, B.L., Yeung, M., 1998. Resolving rightful ownerships with invisible watermarking techniques: limitations, attacks, and implications. IEEE J. Selected Areas Comm. 16 (4), 573–586.

Hartung, F., Kutter, M., 1999. Multimedia watermarking techniques. Proc. IEEE 87, 1079–1107.

Hwang, M.S., Chang, C.C., Hwang, K.F., 1999. A watermarking technique based on one-way hash functions. IEEE Trans. Consumer Electron. 45, 286–294.

Katzenbeisser, S., Veith, H., 2002. Securing symmetric watermarking schemes against protocol attacks. In: Security and

Watermarking of Multimedia Contents IV, Proc. SPIE, vol. 4675, San Jose, CA, January 2002, pp. 260–268.

Kutter, M., Voloshynovskiy, S., Herrigel, A., 2000. The watermark copy attack, in security and watermarking of multimedia contents II. Proc. SPIE 3971, 371–379.

Petitcolas, F.A.P., Anderson, R.J., Kuhn, M.G., 1999. Information hiding—A survey. Proc. IEEE 87, 1062–1078.

Swanson, M.D., Zhu, B., Tewfik, A.H., 1998. Multiresolution scene-based video watermarking using perceptual models. IEEE J. Selected Areas Comm. 16, 540–550.

Voloshynovskiy, S., Pereira, S., Pun, T., Eggers, J., Su, J., 2001. Attacks on digital watermarking: classification, estimation-based attacks, and benchmarks. IEEE Comm. Mag. (August), 118–126.