PRACTICAL KEY LENGTH OF WATERMARKING SYSTEMS

Patrick Bas[†], Teddy Furon^{*}, and François Cayre⁺

[†] CNRS-LAGIS, Lille, France
* INRIA research centre Rennes Bretagne Atlantique, Rennes, France
⁺ GipsaLab, Grenoble, France

ABSTRACT

The paper proposes a new approach for evaluating the security levels of digital watermarking schemes, which is more in line with the formulation proposed in cryptography. We first exhibit the class of equivalent decoding keys. These are the keys allowing a reliable decoding of contents watermarked with the secret key. Then, we evaluate the probability that the adversary picks an equivalent key. The smaller this probability, the higher the key length. This concept is illustrated on two main families of watermarking schemes: DC-QIM (Distortion Compensation Quantization Index Modulation) and SS (Spread Spectrum). The trade-off robustness-security is again verified and gives some counter-intuitive results: For instance, the security of SS is a decreasing function of the length of the secret vector at a fixed Document to Watermark power ratio. Additionally, under the Known Message Attack, the practical key length of the watermarking scheme rapidly decreases to 0 bits per symbol.

Index Terms— Watermarking , Security , Key Length , Quantization , Spread-spectrum

1. INTRODUCTION

The concept of security in watermarking has not a so long history compared to cryptography. The early works base watermarking security assessment on the information about the secret key that leaks from the observations [1, 2]. This is the translation of the definition of security for cryptosystems by C.E. Shannon. However, this conception rarely applies in modern cryptography which almost always relies on computational security¹.

Another difference comes from the fact that the secret keys may not be unique in watermarking. In cryptography, the exact knowledge of the secret key is needed to decrypt cipher texts or to encrypt messages. In watermarking, the secret key is usually a real vector: it is for example the dither of DC-QIM (Distortion Compensated Quantization Index Modulation) or the carrier of SS (Spread Spectrum) schemes. The embedder and the decoder generates this secret vector using the same seed of a pseudo-random generator but the adversary needs not to find back this seed. Indeed it is usually easier to find an estimation of the secret vector good enough to grant the decoding of watermarked contents or the embedding of hidden messages in content. Therefore, we consider that the secret key, denoted by k (or k for vectors) in the sequel, is the secret vector, not the seed. This idea has been first sketched in [3] and this paper investigates it thoroughly.

In cryptography the security of a primitive strongly relies on the size of the secret key. If the implementation of the scheme is not flawed by some security pitfalls, the adversary has no other option than performing an exhaustive search. This attack scans all the possible keys in order to find the one used for encryption. The size N of the cryptographic key in bits is consequently directly linked with the number of trials necessary to perform this exhaustive search attack (a.k.a. brute-force attack). If the key is a uniformly drawn binary word of length N, the probability to pick the right key is $p = 2^{-N}$, the size of the key ensemble is $p^{-1} = 2^N$ or in logarithmic scale $-\log_2(p) = N$ bits.

In this paper, we would like to measure watermarking security in the same manner. To do so, we will ask the following question: What is the probability p that the adversary picks up a convenient key? Then, the security level of the scheme will be called the *key length* measured by $-\log_2(p)$ in bits.

2. INFORMATION THEORETIC SECURITY

Let us denote K the random variable associated to the secret key, \mathcal{K} the space of the secret keys. Before producing any watermarked content, the designer draws the secret key kaccording to a given distribution $p_K(K = k)$. The adversary knows \mathcal{K} and $p_K(K = k)$ but he doesn't know the instantiation k. This lack of knowledge can be measured in bits by the entropy of the key $H(K) \triangleq -\oint_{\mathcal{K}} p_K(k) \log_2 p_K(k)$ (i.e., an integral if K is a continuous r.v. or a sum if Kis a discrete r.v.). Now, suppose the adversary sees N_o observations $O^{N_o} = (O_1, \ldots, O_{N_o})$. The nature of these observations defines the attack. In this paper, we restrict

This work was partly founded by the French National Research Agency program referenced ANR-10-CORD-019 under the Estampille project.

¹It is also true that we witness a recent trend to a return to information theoretical security in cryptography.



Fig. 1. The ensemble of keys



Fig. 2. The embedding domain

our attention to the Known Message Attack (KMA - an observation is a pair of a watermarked content and the embedded message). Thanks to these observations, the adversary can refine his knowledge about the key by constructing a posteriori distribution $p_K(k|O^{N_o})$. The information leakage is given in bits by the mutual information $I(K;O^{N_o})$, and the equivocation $h_e(N_o) \triangleq H(K|O^{N_o})$ measures how this leakage decreases the initial lack of information: $h_e(N_o) = H(K) - I(K;O^{N_o})$. The equivocation is a non increasing function. For most of the watermarking schemes, the information leakage is not null, and as the adversary keeps on observing, the equivocation decreases down to 0 (discrete r.v.) or $-\infty$ (continuous r.v.). This means that the adversary as collected enough observations so that he can uniquely identify the secret key k.

The subsections below show a geometrical interpretation of the concept for two watermarking schemes. Thanks to the observations (e.g. watermarked vectors $\mathbf{y}_1, \ldots, \mathbf{y}_{N_o}$ in Fig. 2), the adversary succeeds to restrict the set of possible keys to a smaller set $\mathcal{K}(O^{N_o}) \subset \mathcal{K}$ depicted in Fig. 1.

2.1. DC-QIM

Let us model the host signal by a vector $\mathbf{x} \in \mathbb{R}^{N_v}$. Consider a lattice $\Lambda \subset \mathbb{R}^{N_v}$. For each message $m \in \{1, 2, \dots, M\}$, a coset leader \mathbf{d}_m is defined such that $\cup_{m=1}^M \Lambda + \mathbf{d}_m$ is a finer lattice. Hiding message m in \mathbf{x} yields watermarked vector \mathbf{y} :

$$\mathbf{y} = e(\mathbf{x}, m, \mathbf{k}) = \mathbf{x} + \alpha (Q_{\Lambda}(\mathbf{x} - \mathbf{d}_m - \mathbf{k}) - \mathbf{x} + \mathbf{d}_m + \mathbf{k}), (1)$$

with $Q_{\Lambda}(.)$ the Euclidean quantizer on Λ . The key k is a N_v dimension vector applying a secret shift of the quantizer. Due

to the Λ -periodicity, the key ensemble \mathcal{K} is the Voronoi cell $\mathcal{V}(\Lambda) \triangleq \{\mathbf{v} \in \mathbb{R}^{N_v} | Q_{\Lambda}(\mathbf{v}) = \mathbf{0}\}.$

Paper [4] shows that the security is maximized if **k** has been uniformly drawn over $\mathcal{K} = \mathcal{V}(\Lambda)$, and that, under the KMA attack, the adversary succeeds to narrow the key ensemble down to $\mathcal{K}(O^{N_o}) = \bigcap_{i=1}^{N_o} \mathcal{D}_i$ with $\mathcal{D}_i \triangleq \tilde{\mathbf{y}}_i - \mathbf{d}_{m_i} - (1 - \alpha)\mathcal{V}(\Lambda)$ and $\tilde{\mathbf{y}} \triangleq \mathbf{y} - Q_{\Lambda}(\mathbf{y})$. The equivocation is then the expectation of the log-volume of this set: $h_e(N_o) = E[\log(\operatorname{vol}(\mathcal{K}_{N_o}))]$. $\mathcal{K}(O^{N_o})$ is in general hard to compute and [4] gives fast approximations.

2.2. Spread Spectrum

Consider a spread spectrum one-bit watermarking s.t. $\mathbf{y} = e(\mathbf{x}, m, \mathbf{k}) = \mathbf{x} + (-1)^m \mathbf{k}$, with $m \in \{-1, 1\}$. The host is modeled by a white Gaussian vector of power σ_X^2 and N_v samples. The secret key is usually drawn also as $\mathbf{K} \sim \mathcal{N}(\mathbf{0}, \sigma_K^2 \mathbf{I}_{N_v})$. This time, $\mathcal{K} = \mathbb{R}^{N_v}$ is not bounded. Yet, thanks to the AEP, for N_v sufficiently large, the key indeed lies in a bounded volume so-called the typical set with high probability: $\mathbb{P}[\mathbf{k} \in \mathcal{K}_{\epsilon}] > 1 - \epsilon$ with $\operatorname{vol}(\mathcal{K}_{\epsilon}) \leq 2^{H(\mathbf{K})+n\epsilon}$. In this very simple case, $\mathcal{K}_{\epsilon} = \{\mathbf{v} \in \mathbb{R}^{N_v} : |\frac{\|\mathbf{v}\|^2}{N_v \sigma_K^2} - 1| < 2\epsilon\}$ and $H(\mathbf{K}) = N_v/2 \cdot \log_2(2\pi e \sigma_K^2)$. Under the KMA, $h_e(N_o) = N_v/2 \cdot \log_2(2\pi e \frac{\sigma_X^2 \sigma_K^2}{\sigma_X^2 + N_o \sigma_K^2})$, and the estimation $\hat{\mathbf{K}}$ is given in [5, Eq.(3) (4)]. Again, thanks to the AEP, $h_e(N_o)$ can be seen as the log volume of the typical set of the estimation $\hat{\mathbf{K}}$.

3. PROBABILISTIC WATERMARKING SECURITY

Our new definition of the security is not centered on the information leakage or estimation of a key. It is based on the fact that the secret key may not be unique in digital watermarking because there exist equivalent keys.

From key k, a watermarking scheme derives an encoder $\mathbf{y} = e(\mathbf{x}, m, k)$ and a decoder $\hat{m} = d(\mathbf{y}, k)$ which can be thought as regions in the embedding domain. The decoding region is defined as $\mathcal{D}_m(k) \triangleq \{\mathbf{y} \in \mathbb{R}^{N_v} : d(\mathbf{y}, k) = m\}$. To hide message m, the encoder pushes the host vector \mathbf{x} deep inside $\mathcal{D}_m(k)$, and this creates an embedding region $\mathcal{E}_m(k)$. To provide robustness, $\mathcal{E}_m(k) \subset \mathcal{D}_m(k)$ s.t. if the vector extracted from an attacked content $\mathbf{z} = \mathbf{y} + \mathbf{n}$ goes out of $\mathcal{E}_m(k)$, \mathbf{z} might be still in $\mathcal{D}_m(k)$ and the correct message is decoded.

We introduce the set of equivalent decoding keys $\mathcal{K}_{eq}^{(d)}(k,\epsilon)$ as the set of keys that allow a decoding of the hidden messages embedded with k with probability $1 - \epsilon$:

$$\mathcal{K}_{eq}^{(d)}(k,\epsilon) = \{k' \in \mathcal{K} : \mathbb{P}[d(e(\mathbf{x},m,k),k') \neq m] \le \epsilon\}$$
(2)

In the same way, $\mathcal{K}_{eq}^{(e)}(k,\epsilon)$ is the set of keys that allow to embed messages which will be reliably decoded using key k. In Fig. 2, $k' \in \mathcal{K}_{eq}^{(d)}(k,0)$ because $\mathcal{E}_m(k) \subset \mathcal{D}_m(k')$, whereas $k' \notin \mathcal{K}_{eq}^{(e)}(k,0)$ because $\mathcal{E}_m(k') \notin \mathcal{D}_m(k)$. This paper only focuses on the equivalent decoding keys. The goal of the adversary is now to draw a key according to the set of observations $\mathcal{K}_{eq}(O^{N_o})$ which also belongs to \mathcal{K}_{eq} (see Fig.1).

The issue is then the probability $P^{(d)}(\epsilon)$ (or $P^{(e)}(\epsilon)$) that an adversary picks up an equivalent key. For instance, if the keys are uniformly distributed over a bounded set, $P^{(d)}(\epsilon) = \mathbb{E}[\operatorname{vol}(\mathcal{K}_{eq}^{(d)}(k,\epsilon))]/\operatorname{vol}(\mathcal{K})$ (see Fig. 1). Like in Sect. 2, we also would like to investigate how the observations O^{N_o} increases this probability. Again, if the estimator of k is uniformly distributed over $\mathcal{K}(O^{N_o})$, then $P^{(d)}(\epsilon, N_o) = \mathbb{E}[\operatorname{vol}(\mathcal{K}_{eq}^{(d)}(k,\epsilon) \cap \mathcal{K}(O^{N_o}))/\operatorname{vol}(\mathcal{K}(O^{N_o})))]$. Finally, the security level is then expressed in bits as the key length:

$$L(\epsilon, N_o) \triangleq -\log_2(P^{(d)}(\epsilon, N_o)) \quad \text{bits}, \tag{3}$$

to obtain an analogy with cryptography.

The next sections of this paper compute the key length of two classical watermarking schemes: DC-QIM with a cubic lattice (a.k.a. SCS) and Spread Spectrum.

4. INVESTIGATIONS ON DC-QIM

From now on, we suppose that Λ is a cubic lattice. We can then proceed component-wise and we introduce $\pi^{(d)}(\epsilon, N_o)$ and $\ell(\epsilon, N_o)$ the probability and key length per symbol, a symbol in $\{1, 2, \ldots, M\}$ being embedded per component. For a given component, the secret key is a scalar k, and regions $\mathcal{E}_m(k)$ and $\mathcal{D}_m(k)$ are two intervals of respective lengths $(1-\alpha)\Delta$ and Δ/M both centered on $k+m\Delta/M$. We assume that the message is embedded without error ($\alpha > (M-1)/M$) and that the adversary wants to decode without error ($\epsilon = 0$).

The secret key is uniformly drawn over the interval $\mathcal{K} = [-\Delta/2, \Delta/2]$ s.t. the information theoretical approach evaluates the security to $H(K) = \log_2 \Delta$ bits. Depending on the value of Δ , this quantity can be negative whose interpretability is difficult. On the contrary, the computation of the key length is straightforward. As illustrated in Fig. 3, we have $\mathcal{K}_{eq}^{(d)}(k,0) = [k'_{\min},k'_{\max}]$ and $\operatorname{vol}(\mathcal{K}_{eq}^{(d)}(k,0)) = \Delta(\alpha - 1 + 1/M)$. This yields the key length per symbol:

$$\ell(0,0) = -\log_2(\alpha + 1/M - 1) \quad \text{bits.} \tag{4}$$

Note that the probabilistic approach yields a key length independent of Δ contrary to the information theoretical approach, and $\ell(0,0) \rightarrow \infty$ for $\alpha \rightarrow (M-1)/M$. This means that only k' = k allows a decoding without errors because $\mathcal{E}_m(k) = \mathcal{D}_m(k), \forall m$. On the other hand, if $\alpha = 1$ (no distortion compensation), $\pi^{(d)}(\epsilon, N_o) = M^{-1}$. All samples with membedded inside are decoded as m' and there is one chance out of M that m' = m.

For $N_o = 1$ in the KMA setup, the information theoretical approach evaluates the security to $h_e(1) = \log_2((1 - \alpha)\Delta)$ bits [5, Eq.(16)]. $\mathcal{K}(O^{N_o})$ is defined by the feasible region \mathcal{D}_1 (see Sub. 2.1), which is in this case the interval of length $(1 - \alpha)\Delta$ centered on $y_1 + m\Delta/M$. Depending of the value



Fig. 3. Computation of vol $(\mathcal{K}_{eq}^{(d)}(k, 0))$ for DC-QIM.

of y_1 , we can compute the probability that a key belonging to the feasible set is included in the equivalent set, and its expectation enables to compute $\pi^{(d)}(0,1)$. Finally the key length is:

$$\ell(0,1) = \begin{cases} -\log_2 \frac{(\alpha + (1-M)/M)(5-5\alpha - 1/M)}{4(1-\alpha)^2} \text{ bits } & \text{if } \alpha \le \alpha' \\ 0 \text{ bit } & \text{if } \alpha > \alpha', \end{cases}$$

where α' is the root of equation $(\alpha + (1 - M)/M)(5 - 5\alpha - 1/M) = 4(1 - \alpha)^2$, $\alpha' \in [0, 1]$. The feasible set is always included in the equivalent region for $\alpha > \alpha'$.

For $N_o > 1$, we must use Monte-Carlo simulations. For a given run, we draw a key k, we generate N_o observations and compute $\mathcal{K}(O^{N_o})$, which is also an interval in this case, and its intersection with $\mathcal{K}_{eq}^{(d)}(k,0)$. This gives us the probability $\pi^{(d)}(0, N_o)$ (see Sect. 3). We finally take the log of the average of $\pi^{(d)}(0, N_o)$ over N_r runs.



Fig. 4. key length per symbol for DC-QIM and M = 2 vs. the distortion compensation parameter α .

Fig. 4 gives the key length per symbol. Parameter α is usually increased to gain some robustness, whereas the key length is a decreasing function. This illustrates the trade-off robustness-security. The probability to disclose the secret key over N_s symbol is $\pi^{(d)}(\epsilon, N_o)^{N_s}$. Therefore, the total key length is $L(\epsilon, N_o) = N_s \ell(\epsilon, N_o)$. Note that the key length in bits might be bigger than the vector dimension N_s because $\ell(\epsilon, N_o)$ can be bigger than 1 bit. This is due to the fact than the key is not a binary word but a vector in \mathbb{R}^{N_v} . Also, $L(\epsilon, N_o)$ should be clipped to the length of the binary seed of the PRNG. We conjecture that, in KMA, for a given $\alpha \in ((M-1)/M, 1]$, there exists a N_o^* s.t. $N_o \ge N_o^*$ sets in expectation $\ell(0, N_o) = 0$, i.e. the scheme is totally broken whatever the length N_s .

5. INVESTIGATIONS ON SPREAD SPECTRUM

We consider the simple case of a one-bit embedding with $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \sigma_X^2 \mathbf{I})$ and $\mathbf{y} = \mathbf{x} + (-1)^m \mathbf{k}$, where $m \in \{0, 1\}$. The decoder is correlation based: $d(\mathbf{y}, \mathbf{k}) = 0$ if $\mathbf{y}^\top \mathbf{k} > 0$, 1 else. When a noise is added giving $\mathbf{z} = \mathbf{y} + \mathbf{n}$ with $\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma_N^2 \mathbf{I})$, the Bit Error Rate equals $\Phi\left(-\|\mathbf{k}\|/\sqrt{\sigma_N^2 + \sigma_X^2}\right) (\Phi(.)$ the cumulative distribution function of a normal r.v.). Even without noise, the BER is not null and we define $\eta \triangleq \Phi(-\|\mathbf{k}\|/\sigma_X)$.

Now, the adversary draws a key \mathbf{k}' and decoding a content watermarked with \mathbf{k} yields $d(\mathbf{y}, \mathbf{k}') = 0$ if $\mathbf{x}^{\mathsf{T}}\mathbf{k}' + (-1)^m \mathbf{k}^{\mathsf{T}}\mathbf{k}' > 0$. Therefore, the BER is $\Phi(-\mathbf{k}^{\mathsf{T}}\mathbf{k}'/\sigma_X ||\mathbf{k}'||)$, which is lower than ϵ if

$$\frac{\mathbf{k}^{\mathsf{T}}\mathbf{k}'}{\|\mathbf{k}'\|\|\mathbf{k}\|} \ge \tau(\epsilon, \mathbf{k}) \triangleq -\frac{\sigma_X}{\|\mathbf{k}\|} \Phi^{-1}(\epsilon) = \frac{\Phi^{-1}(\epsilon)}{\Phi^{-1}(\eta)}.$$
 (5)

The LHS is the cosine of the angle between k and k' which is always lower than 1. Thus, there exist equivalent keys iff $\epsilon > \eta$. $\mathcal{K}_{eq}(\epsilon, \mathbf{k})$ is then a single hypercone of axis k and angle $\arccos(\Phi^{-1}(\epsilon)/\Phi^{-1}(\eta))$.

For $N_o = 0$, the probability of drawing a key $\mathbf{k}' \sim \mathcal{N}(\mathbf{0}, \sigma_K^2 \mathbf{I}_{N_v})$ inside $\mathcal{K}_{eq}(\epsilon, \mathbf{k})$ is the ratio of the solid angle of this hypercone and the full space. This equals $\beta(\epsilon, \mathbf{k}) \triangleq (1 - I_{\tau(\epsilon, \mathbf{k})^2}(1/2, (N_v - 1)/2))/2$ where I(.) is the regularized incomplete beta function. Finally, $\pi^{(d)}(\epsilon, 0) = \mathbb{E}[\beta(\epsilon, \mathbf{k})]$, where the expectation is over $p_{\mathbf{K}}(\mathbf{K} = \mathbf{k})$.

For $N_o > 0$, we suppose without loss of generality that the embedded messages were all set to 0. Then one estimator $\hat{\mathbf{k}}$ is the average of the watermarked contents, e.g. $\hat{\mathbf{k}} =$ $N_o^{-1} \sum_{i=1}^{N_o} \mathbf{x}_i + \mathbf{k}$ and $p(\hat{\mathbf{k}}|O^{N_o})$ is $\mathcal{N}(\mathbf{k}, N_o^{-1}\sigma_X^2 \mathbf{I}_{N_v})$. The probability of drawing an estimation inside $\mathcal{K}_{eq}(\epsilon, \mathbf{k})$ is upper bounded by the cumulative distribution function of a noncentral F-distribution variable of degrees of freedom $\nu_1 = 1$, $\nu_2 = N_v - 1$ and non centrality parameter $\lambda(\mathbf{k}) = N_o \frac{\|\mathbf{k}\|^2}{\sigma_X^2}$, weighted by the probability $\mathbb{P}[\mathbf{k}'^{\mathsf{T}}\mathbf{k} > 0]$:

$$\gamma(\epsilon, N_o, \mathbf{k}) \triangleq \left[1 - F\left(\frac{(N_v - 1)\tau(\epsilon, \mathbf{k})^2}{1 - \tau(\epsilon, \mathbf{k})^2}; \nu_1, \nu_2, \lambda(\mathbf{k})\right) \right] \\ * \Phi\left(\sqrt{\lambda(\mathbf{k})}\right)$$
(6)

Finally, $\pi^{(d)}(\epsilon, N_o) \leq \mathbb{E}[\gamma(\epsilon, N_o, \mathbf{k})]$, where the expectation is over $p_{\mathbf{K}}(\mathbf{K} = \mathbf{k})$.

Fig. 5 plots the key lengths for three different setups ($N_o = 0$, $N_o = 1$ and $N_o = 10$) and DWR = 10dB. The most important fact is that the key length is a decreasing function w.r.t. N_v , the length of k. This seems counter-intuitive and contradicts a claim of [3] (a key length proportional to N_v). This is



Fig. 5. key length per symbol for SS: DWR = 10dB, $\epsilon = 10^{-2}$, \star represents Monte-Carlo simulations using 10^3 random keys and observing 10^2 times N_o contents.

indeed normal since it stems from the trade-off robustness vs. security: a bigger N_v improves the robustness but decreases the key length. We also note the devastating effect of KMA: for $N_v = 200$ the key length decreases from approximately 50 bits to 10 bits for $N_o = 1$ and 10^{-3} bits for $N_o = 10$. Some key lengths were estimated by Monte-Carlo simulations when possible (ie. not too big) and they confirm the theoretical values. Again, these are key lengths per symbol and shall be multiplied by N_s for a multi-bit SS scheme.

6. CONCLUSION AND PERSPECTIVES

This novel definition of the security in watermarking enables to compute the security of a scheme regarding an exhaustive search strategy. It also enables to compute and analyse a watermarking scheme the same way as one would quantify the security of a cryptographic system. Yet this first analysis shows that the key length, which can be very important when the adversary doesn't have access to any observations; can also decrease dramaticaly whenever the adversary can uses observations.

7. REFERENCES

- M. Barni, F. Bartolini, and T. Furon, "A general framework for robust watermarking security," *Signal Processing*, vol. 83, no. 10, pp. 2069– 2084, Oct. 2003.
- [2] F. Cayre, C. Fontaine, and T. Furon, "Watermarking security: Theory and practice," *IEEE Trans. Signal Processing*, vol. 53, no. 10, pp. 3976 – 3987, Oct. 2005.
- [3] I. Cox, G. Doerr, and T. Furon, "Watermarking is not cryptography," in Proc. Int. Work. on Digital Watermarking, Jeju island, Korea, Nov. 2006, vol. 4283 of LNCS, Springer-Verlag.
- [4] L. Pérez-Freire, F. Pérez-González, T. Furon, and P. Comesańa, "Security of lattice-based data hiding against the known message attack," *IEEE Trans. on Information Forensics and Security*, vol. 1, no. 4, pp. 421–439, Dec. 2006.
- [5] L. Pérez-Freire and F. Pérez-González, "Spread-spectrum watermarking security," *IEEE Trans. on Information Forensics and Security*, vol. 4, no. 1, pp. 2–24, 2009.