

PAPER

On the Security of BioEncoding Based Cancelable Biometrics

Osama OUDA^{†,††a)}, Nonmember, Norimichi TSUMURA[†], and Toshiya NAKAGUCHI[†], Members

SUMMARY Proving the security of cancelable biometrics and other template protection techniques is a key prerequisite for the widespread deployment of biometric technologies. BioEncoding is a cancelable biometrics scheme that has been proposed recently to protect biometric templates represented as binary strings like iris codes. Unlike other template protection schemes, BioEncoding does not require user-specific keys or tokens. Moreover, it satisfies the requirements of untraceable biometrics without sacrificing the matching accuracy. However, the security of BioEncoding against smart attacks, such as correlation and optimization-based attacks, has to be proved before recommending it for practical deployment. In this paper, the security of BioEncoding, in terms of both non-invertibility and privacy protection, is analyzed. First, resistance of protected templates generated using BioEncoding against brute-force search attacks is revisited rigorously. Then, vulnerabilities of BioEncoding with respect to correlation attacks and optimization based attacks are identified and explained. Furthermore, an important modification to the BioEncoding algorithm is proposed to enhance its security against correlation attacks. The effect of integrating this modification into BioEncoding is validated and its impact on the matching accuracy is investigated empirically using CASIA-IrisV3-Interval dataset. Experimental results confirm the efficacy of the proposed modification and show that it has no negative impact on the matching accuracy.

key words: template protection, cancelable biometrics, BioEncoding, correlation attacks, optimization based attacks

1. Introduction

Although biometrics-based authentication systems exhibit many usability advantages over traditional authentication systems, they suffer from several security and privacy concerns [1]. As a result, many template protection techniques have been proposed in the last few years to deal with these issues [2]. Generally, template protection techniques may be classified into two main categories; namely, biometric encryption (BE) and cancelable biometrics (CB). In BE techniques, such as fuzzy commitment [3], fuzzy extractors [4], and fuzzy vaults [5], biometric templates are linked with a user-specific key to produce a biometrically encrypted pseudo-identity for the user so that the key can be released only if the true biometric template is present on verification. On the other hand, CB methods, such as distorting transforms [6], BioHashing [7], and BioEncoding [8], generate revocable protected templates from true biometric tem-

plates through applying different non-invertible transforms to true templates in different applications. Matching is done in the transform domain after applying the same transform (applied in enrollment) to a fresh template during authentication. Any template protection scheme should satisfy the following requirements [9]:

Accuracy A template protection scheme should not introduce significant degradation in the recognition performance of the unprotected biometric system.

Revocability It should be easy to revoke (cancel) a protected template if it is stolen or compromised.

Irreversibility Retrieving original templates from protected ones should be computationally infeasible.

Diversity It should be possible to generate large number of protected templates (to be used in different applications) for the same biometric.

Unlinkability It should not be possible for an adversary to determine whether different protected templates belong to the same user.

Although much attention has been given to proving the accuracy and revocability requirements in almost all template protection methods proposed in the literature so far, less attention has been paid to analyzing the security of such systems rigorously. Just recently, a few researchers have studied the security flaws of some template protection schemes. Scheirer and Boulton [10] showed *theoretically* that both the fingerprint biometric encryption algorithm of Soutar et al. [11] and the fuzzy vault scheme [4] are vulnerable to three different classes of attacks and they concluded that both techniques are not suitable for preserving privacy or enhancing security of biometric templates. Adler [12] discussed the vulnerability of the method in [11] to the hill climbing attack and the experimental results showed that an estimate of the enrolled biometric image could be regenerated and hence the stored secret could be released assuming that the system leaks information on the matching score. Kholmatov and Yanikoglu [13] realized the correlation attack against the fuzzy vault scheme using fingerprints, and the obtained results proved that the fuzzy vault scheme is indeed vulnerable to correlation attacks. Nagar et al. [14] analyzed the security of two well-known cancelable biometrics techniques, namely, distorting transforms for fingerprints [6] and BioHashing [7], and their analysis showed that both techniques are susceptible to invertibility attacks. Zhou and Kalker [15] argued that essential information about the original biometric features could be leaked if an attacker gained access to more BioHashes of the same user. Zhou et al. [16] described

Manuscript received 0, 0000.

Manuscript revised 0, 0000.

Final manuscript received 0, 0000.

[†]Graduate School of Advanced Integration Science, Chiba University, 1-33, Yayoi-cho, Inage-ku, Chiba, 263-8522, Japan.

^{††}Faculty of Computers and Information Sciences, Mansoura University, Mansoura 35516, Egypt.

a) E-mail: oouda@graduate.chiba-u.jp

DOI: 10.1587/transinf.E0.D.1

how the correlation of biometric features could be exploited to attack fuzzy extractors-based techniques. Simoens et al. [17] demonstrated how protected templates generated by code-offset [3] and bit-permutation sketches [4] can be linked and reversed. BioEncoding [8] is a recently proposed CB scheme for protecting standard iris codes* [18] that offers several advantages over other CB techniques. It has been shown in [8] that protected templates, referred to as *BioCodes*, generated from true templates can be used efficiently to verify the user identity without deteriorating the recognition accuracy achieved using original (unprotected) recognition systems. Besides, unlike other template protection systems, BioEncoding can be used as a one-factor (tokenless) method. Moreover, it is easy to implement and can be integrated simply with current iris recognition systems. Therefore, we believe that if the security of BioEncoding against different possible attacks is proved, it would be one of the most suitable template protection candidates for widespread deployment. In this paper, the security aspects of BioEncoding are discussed rigorously. First, we investigate the vulnerabilities of BioEncoding with respect to different categories of threats and attacks. Then, a simple yet effective modification to BioEncoding is proposed to enhance the irreversibility and privacy of BioCodes. The impact of integrating the proposed modification into BioEncoding on both security and matching accuracy is discussed. The rest of this paper is organized as follows. In Sect. 2, the main procedure of BioEncoding is reviewed. In Sect. 3, different security aspects of BioEncoding are discussed. In Sect. 4, we describe some methods for securing BioEncoding against reversibility and privacy threats. Section 5 presents a set of experiments for validating our theoretical security analysis. Finally, Section 6 concludes the paper.

2. BioEncoding Overview

The basic idea behind BioEncoding lies in employing random addressing of a set of randomly generated binary digits in order to achieve the irreversibility property of CB. The main advantage of BioEncoding is that no user-specific key/password needs to be associated with each user. Rather, the random sequence employed in the cancelable transformation process of BioEncoding can be set common to all users. Hence, this random sequence can be stored centrally in the application database. The transformation process of BioEncoding, illustrated in Fig. 1, can be summarized in the following steps [8]:

1. Group bits of the true binary template T into n/m words of fixed length m , where n is the number of bits in T and $'/'$ denotes integer division with truncation of the result toward zero.
2. Generate a (pseudo-) random sequence S of length $l = 2^m$ using a random seed that can be stored in a centralized storage (application database).

*In fact, BioEncoding can be applied to the binary representation of any biometrics modality

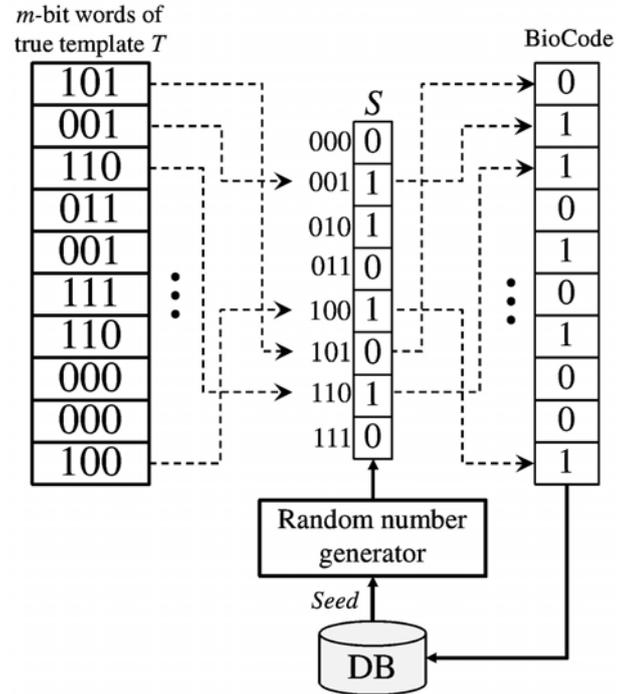


Fig. 1 Illustration of BioEncoding transformation process where $m = 3$.

3. Map each word in T to a single bit value in S whose location is addressed by the value in that word.
4. Constitute the protected *BioCode* from the set of n/m addressed bits.
5. Store the *BioCode* in the centralized storage and discard the original (unprotected) template.

An example of the transformation process of BioEncoding where $n = 30$, $m = 3$, and $l = 8$ is shown in Fig. 1.

Fortunately, the discrimination between the resulting BioCodes is due to the variation existing between the true templates and not resulted from the random sequence. Hence, the same random sequence could be used with all users and that is why BioEncoding can be employed without the need for physical user-specific tokens since there no user-specific information that need to be stored.

On the other hand, it is straightforward to show that BioEncoding meets the requirements of the cancelable biometrics construct. For the revocability requirement, in case of compromising the stored BioCodes, changing the random sequence and re-enrolling the users would generate a new set of protected BioCodes. Regarding the diversity requirement, it is possible to generate a large number of different BioCodes especially for large address words since 2^{2^m} different random sequences can be generated using address words of length m . For the non-invertibility requirement, the many-to-one nature of the transformation process can guarantee its irreversibility. In [8], it has been shown that BioEncoding is robust against the brute-force search attacks.

It is worth noting that the enrolment and verification templates might not be aligned sufficiently due to rotational inconsistencies caused by head tilt during the acquisition

of iris images. Therefore, the verification template should be shifted several times in both directions, as suggested by Daugman [18], and then the BioEncoding transformation process should be repeated after each shift. Only the BioCode that gives the smallest Hamming distance with the enrolment BioCode is considered as explained in [8].

3. Security Vulnerabilities of BioEncoding

The main motivation behind developing cancelable biometrics and other template protection constructs is to address two specific issues that are inherent to the use of biometrics in identity authentication; these are, protecting biometric features from unauthorized disclosure and preserving users' privacy. The importance of preserving biometric features undisclosed to adversaries is due to the fact that each person has a limited number of permanently associated biometric traits and hence if a biometric trait is compromised, revoking it would not be as easy as revoking compromised passwords or tokens. Moreover, sensitive personal information such as kinship, gender, or diseases that a person may be suffering from, could be disclosed if the true biometric features are revealed [17]. That is why protected templates are required to be *noninvertible*. However, even if it is ensured that biometric features cannot be disclosed to adversaries, users' privacy could still be exposed to attackers who may try to track users across applications, via using cross-matching between different databases, to, for example, determine whether they are registered in a particular application or not. That is why template protection techniques are required to generate a large number of *diverse* (un-linkable) protected templates from the same biometric data.

Thus, based on these two underlying objectives, the security of any template protection scheme should be assessed in accordance to two main criteria [19]: 1) irreversibility, and ii) diversity. In this section, we discuss the most important security issues and threats that are related to template protection systems in general and BioEncoded templates in particular according to the mentioned criteria.

3.1 Reversibility Attacks

With respect to reversibility attacks, three different categories of attacks are investigated; namely, brute-force search attacks, correlation attacks (a.k.a record multiplicity attacks) and optimization-based attacks. For all attacks, we assume that the attacker is familiar with the encoding algorithm and the random sequence S is known. Moreover, we assume that the attacker can gain access to the matching score in the optimization-based attacks.

3.1.1 Brute Force Search Attacks

In a brute-force attack, an attacker tries every possible solution in the solution space until he finds the one he searches for. A template protection system is secure enough if this attack is computationally infeasible and if no other attack is

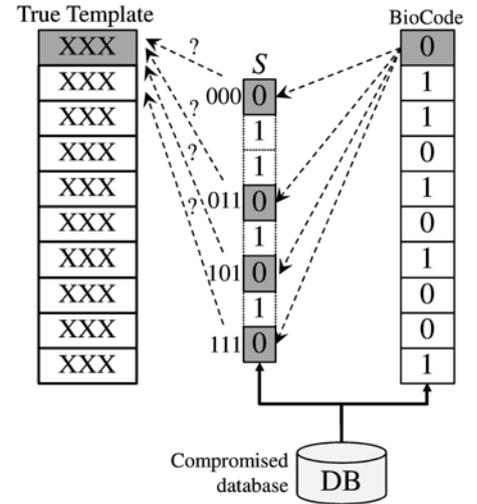


Fig. 2 Example of BioCode inversion via brute-force search ($m = 3$).

computationally less expensive than it.

In BioEncoding, the number of zeros in S (the random sequence used in the mapping process) is expected to be equal to the number of ones $\approx 2^{m-1}$, where m is the length of any address word in the true (binary) template. Therefore, every bit in a given BioCode could be originated from 2^{m-1} address words in the true template, as illustrated in Fig. 2 (for $m = 3$). If the number of bits in a given template is n , the number of bits in BioCodes generated from that template would be n/m bits and hence an attacker needs to make, at most, $(2^{m-1})^{n/m}$ ($\approx 2^n$, when m is large) trails in order to recover all the bits in the original template. That is, recovering the original template from a compromised BioCode (together with its corresponding random sequence S) is almost as difficult as guessing the values of all the bits in the true template, which is computationally infeasible.

However, it might be sufficient to recover specific percentage of the original bits rather than recovering the original template entirely. For example, different iris codes for the same eye could differ in 10 – 20% of the bits [4] and hence iris codes that are similar in 80% of their bits can be considered as belonging to the same eye. Therefore, it is important to investigate not only the probability of recovering the (exact) true template using a single guess but also to measure the similarity between the true template and the reversed pattern using a single random guess.

BioEncoding would be robust against brute-force search attacks if the normalized Hamming distance between the true template and a randomly guessed pattern ≈ 0.5 . Both high and low Hamming distances would leak some information about the true template since high distances would indicate larger similarity between the true template and the *inverse* of the guessed pattern.

If an adversary gained access to a protected BioCode along with its corresponding random sequence S , he would try to reverse the true template by reversing each zero or one in the BioCode through selecting (randomly) one address

Table 1 Hamming distances between 3-bit words

3-bit words				Hamming distance			
w_i	b_2	b_1	b_0	$d_H = 0$	$d_H = 1$	$d_H = 2$	$d_H = 3$
w_0	0	0	0	w_0	w_1, w_2, w_4	w_3, w_5, w_6	w_7
w_1	0	0	1	w_1	w_0, w_3, w_5	w_2, w_4, w_7	w_6
w_2	0	1	0	w_2	w_0, w_3, w_6	w_1, w_4, w_7	w_5
w_3	0	1	1	w_3	w_1, w_2, w_7	w_0, w_5, w_6	w_4
w_4	1	0	0	w_4	w_0, w_5, w_6	w_1, w_2, w_7	w_3
w_5	1	0	1	w_5	w_1, w_4, w_7	w_0, w_3, w_6	w_2
w_6	1	1	0	w_6	w_2, w_4, w_7	w_0, w_3, w_5	w_1
w_7	1	1	1	w_7	w_3, w_5, w_6	w_1, w_2, w_4	w_0

word from all words that address the value of ‘0’ or ‘1’, respectively. In order to formally analyze the effectiveness of this reversing strategy, we need to recall the following two facts: 1) the Hamming distance, d_H , between any two different m -bit words, w_i and w_j , could be a value that ranges from 1 to m , and 2) for any m -bit word w_i , there are $\binom{m}{k}$ words, w_j , that have the same length m and differ with w_i in exactly k bits, that is $d_H(w_i, w_j) = k$. For example, for any 3-bit word in Table 1, there are $\binom{3}{1}$ words differ in 1 bit, $\binom{3}{2}$ words differ in 2 bits, one word, $\binom{3}{3}$, differs in the three bits, and one word that does not differ in any bit (itself).

Denoting the true word and the reversed word by w_{true} and w_{rev} , respectively, and letting w_k to be any binary word that differs from w_{true} in exactly k bit positions ($k \neq 0$), the following events could be defined:

- $A : d_H(w_{rev}, w_{true}) = 0$.
- $A' : d_H(w_{rev}, w_{true}) \neq 0$.
- $B : w_k$ and w_{true} address the same bit value in S .
- $C : d_H(w_{rev}, w_{true}) = k, \quad k = 1$ to m (word size).

The probabilities of the above events would be:

$$P(A) = \frac{1}{2^{m-1}} = 2^{1-m} \quad (1)$$

$$P(A') = 1 - 2^{1-m} \quad (2)$$

$$P(B) = \binom{m}{k} \left(\frac{1}{2^m - 1} \right) \quad (3)$$

$$\begin{aligned} P(C) &= P(A'B) \\ &= P(A')P(B) \\ &= \binom{m}{k} \left(\frac{1 - 2^{1-m}}{2^m - 1} \right) \end{aligned} \quad (4)$$

Recall that the above probabilities are evaluated under the assumption that the number of zeros in S is equal to the number of ones = 2^{m-1} . Note also that the probability of C can be evaluated in terms of A' and B since the probability that the Hamming distance between w_{true} and w_{rev} is k ($1 \leq k \leq m$) can be expressed as the probability that the two words are not identical ($k \neq 0$) joined by the probability that the reversed word (which is at distance k from the

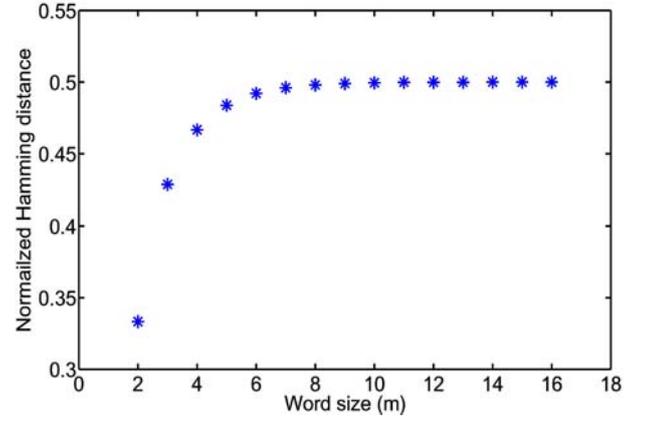


Fig. 3 The expected fractional Hamming distances between true templates and reversed patterns using one trial of a brute-force attack for different word sizes.

true word) addresses the same bit value addressed by the true word. Moreover, it is worth noting that both A' and B are independent since the probability of A' depends only on the word size m . Now, the expected bit error rate (BER), that indicates the expected normalized Hamming distance between the reversed template and the true template, could be formulated as follows:

$$BER = \frac{1}{m} \sum_{k=1}^m k P(d_H(w_{rev}, w_{true}) = k) \quad (5)$$

Figure 3 shows the expected normalized Hamming distances between the reversed patterns and the true templates for $m = 2$ to 16 calculated according to Eq. (5). It is clear from the figure that for small m values (< 4), large percentage of true templates could be revealed (approximately sixty percent of the bits in the true template could be disclosed when $m = 3$ and more than sixty five percent when $m = 2$) using one trial of this brute-force search attack. However, for $m \geq 4$, it is shown that BioEncoding is robust against the brute-force attacks since it is assured that the randomly guessed pattern and the true template would be uncorrelated (approximately 50% of the bits are different). Accordingly, we recommend to use address words of size $m \geq 4$ since the security of BioEncoding against brute-force attacks might be questionable for smaller lengths.

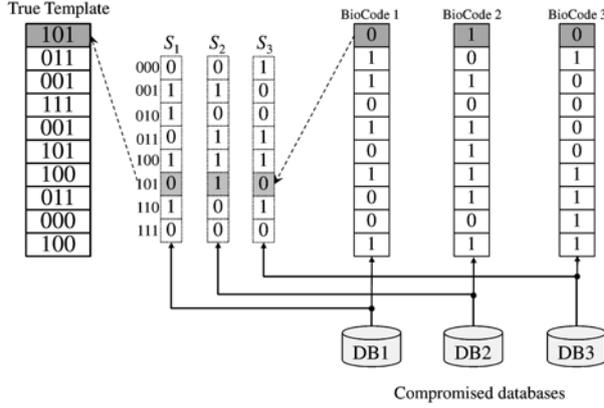


Fig. 4 Example of BioCode inversion via correlation attack with 3 compromised databases ($m = 3$).

3.1.2 Correlation Attacks

For many template protection schemes, even though it might be infeasible to obtain the original template from a single protected template, it might be possible for an attacker to get an exact (or at least approximate) version of the original template by correlating different protected templates created from the same biometric data if he could gain access to several protected templates (along with their associated keys used in the transformation/binding processes) that are used in different applications. This is sometimes called a record multiplicity attack [10].

One of the most vulnerable systems to the correlation attack is the fuzzy vault scheme [5]. The security of template protection schemes that are based on the fuzzy vault construct relies on adding randomly generated points (called chaff points) to a polynomial projection of the true biometric data (for example, fingerprint minutiae) to hide their very existence. As described in [10], given two or more fuzzy vaults generated using the same biometric data and different chaff points, it is possible to recover the true biometric data by correlating the compromised vaults since the only thing in common in these vaults will be the genuine data. Using only two vaults created from the same fingerprint minutiae for 200 different fingers, the correlation attack is realized against fuzzy fingerprint vault in [13] and the experimental results showed that 59% of vaults could be reconstructed successfully. Biometric encryption might be also susceptible to correlation based attacks. It is argued in [10] that since the same phase information are used in all user records, it would be simple to estimate the user's phase data, and hence an estimation of the full signal, if an adversary could acquire multiple samples of the same user in different applications. With regard to BioHashing, it has been shown recently that combining multiple BioHashes that are generated from the same biometric data can reveal essential information about the original features [15].

With respect to BioEncoding, compromising more than one BioEncoded template can not only let the attacker to re-

veal some information about the original biometric data, but also may enable him to recover the original features entirely even with a limited number of compromised templates. As explained in the previous Section and illustrated in Fig. 2, due to the many-to-one transformation adopted in BioEncoding, every bit in a (single) compromised BioCode could be originated from 2^{m-1} address words in the true template. Assuming that the random sequences that are employed in different applications are independent and identically distributed (for example, S_1 , S_2 and S_3 in Fig. 4), it would be expected that the fractional Hamming distance between any pair of these sequences to be 0.5 (i.e. 50% of the bits in both sequences would be similar).

Under this assumption, if the number of the compromised databases is 2, then the many-to-one transformation of BioEncoding would reduce from 2^{m-1} -to-1 to 2^{m-2} -to-1. That is, every corresponding two bits in the (two) compromised BioCodes could be originated from 2^{m-2} address words in the original template rather than 2^{m-1} words as in the case of compromising a single BioCode. For example, in Fig. 4, each '0' in 'BioCode1' could be reversed to 1 out of 4 (since $m = 3$, then $2^{m-1} = 2^2 = 4$) address words in S_1 ('000', '011', '101', and '111'), but if both 'BioCode1' and 'BioCode2' are compromised, there would be four patterns ('00', '01', '10', and '11'), composed from corresponding bits in both BioCodes, rather than only two patterns ('0' and '1') as in the case of compromising a single BioCode. Now the attacker would have only two candidates to choose from to reverse any pattern from these four patterns. For instance, the pattern '01' could be reversed to one of the two address words '011' and '101'. That is, in case of compromising 'BioCode1' only, the attacker needs to invert a 4-to-1 (2^{m-1} -to-1) transformation, however, in case of compromising both 'BioCode1' and 'BioCode2', the attacker needs to invert a 2-to-1 (2^{m-2} -to-1) transformation, as stated previously.

Generally, if the number of compromised BioCodes is k , the many-to-one transformation of BioEncoding would reduce from 2^{m-1} -to-1 to 2^{m-k} -to-1. If the number of compromised BioCodes reaches m (that is, $k = m$), the many-to-one transformation would reduce to a one-to-one transformation (i.e. invertible transformation). This implies that for address words of size m , it is possible to recover the entire features of the true biometrics data (or at least a close approximation of the original template) if an attacker could gain access to at least m databases. Figure 4 shows an example of recovering the entire true template from three compromised BioCodes for $m = 3$. In Sect. 4, we propose several approaches for securing BioEncoding against correlation attacks.

3.1.3 Optimization-based Attacks

If the matching score between the *protected* enrollment sample and the *protected* verification sample could be accessed by an adversary, it would be possible to break the protection system via optimization strategies, such as hill-

climbing, by iteratively making small modifications to an initial (random) verification sample retaining only modifications that enhance the matching score, between the resulting transformed sample and the stored protected template, until a sufficient match is achieved.

Adler [12] realized a quantized hill-climbing attack against the biometric encryption technique with facial images. Experimental results showed that an approximate match between the original image and a randomly chosen initial image can be obtained after a number of iterations. Theoretically, protection schemes that are based on a similarity score as well as biometric encryption schemes that employ short error correcting codes are vulnerable to optimization-based attacks [9].

Fortunately, although BioEncoding relies on a score-based matching, it is not susceptible to this kind of attacks due to the strong many-to-one nature of its transformation process. As described in Section 3.1.1, the complete pre-image of a given protected BioCode contains $2^{n(m-1)/m}$ templates. Therefore, if an attacker started with an initial random binary pattern of the same size as the original template and after a number of iterations he ended up with a pattern that could generate an exactly similar BioCode, the probability that this pattern is identical to the original template would be $1/2^{n(m-1)/m}$.

3.2 Privacy Attacks

Even if it is computationally infeasible to recover the original biometrics data from protected ones, determining whether two protected templates are driven from the same biometric might pose a potential threat to users' privacy. One of the main goals of template protection systems is to prevent cross-matching across different applications through using un-linkable pseudo-identities for the same user across different databases. In fact, satisfying this objective is mainly related to the diversity requirement of cancelable biometrics.

In BioEncoding, the lack of diversity could be originated from two main reasons. The first reason is due to the high correlation existing among local biometric features [16] and the second reason is due to similarity that might be found among random sequences employed in the transformation process. The latter problem could be addressed simply through testing random sequences before use to ensure that sequences employed in different applications are uncorrelated. Regarding the first problem, different approaches to diminish correlation among biometric features are described in the next section.

4. Securing BioEncoded Templates

What makes BioEncoding vulnerable to correlation attacks is that the values of address words do not change across the different applications of the transformation process and hence all the corresponding bits (bits at same location i) in different compromised BioCodes are always addressed by

the same value; that is, value of the address word at location i in the original template (assuming that this value is always the same regardless of the intra-user variations). For example, the first bit in all the three BioCodes shown in Fig. 4 are addressed by the same address word ('101') in the true template. Therefore, modifying the values of address words before every application of BioEncoding would make its security against correlation attacks as robust as its security against brute-force search attacks.

We suggest three different approaches to hinder the record multiplicity attack on BioEncoding. The essence of the three approaches is based on changing the values of address words, in the true template, before every application of the mapping process.

- The first approach is to use different lengths of address words in different applications. If an attacker gained access to two BioCodes belonging to the same user, he would not be able to link the i th bit in the two BioCodes to the same word in the true template since the value of the i th address word used in the first application would be different from the value of the i th address word in the second application due to the change of word size. Although this approach could make BioEncoding more robust against record multiplicity attacks, it restricts the renewability capacity of BioEncoding since one cannot increase the size of address words at will. Address words of large lengths would result in small BioCodes that might be more vulnerable to simple brute-force attacks.
- A more efficient approach is to permute the original template, before applying the BioEncoding transformation, using different (secret) permutations in different applications. Figure 5 shows a block diagram of the enrollment and verification modules of the proposed modified BioEncoding using this approach. Unlike the first approach, the same size of address words could be used in different applications since the values of corresponding address words are ensured to be different as a result of the permutation process. Obviously, the secret permutation should be stored in the centralized database to be used during verification. Storing this random permutation key would not affect the security of the system since what is encoded is the permuted template and hence this key would be useless to the attacker unless he/she could reverse the encoded sequence which is practically infeasible for $m \geq 4$ as explained in subsection 3.1.1.
- Alternatively, the original template could be XORed with a different random sequence of the same length before applying BioEncoding in each new application. Similar to the permutation approach, changing the values of address words is guaranteed after this XORing process. Also, the random sequence which is XORed with the true template needs to be stored to be used during verification. Only if the attacker could reverse

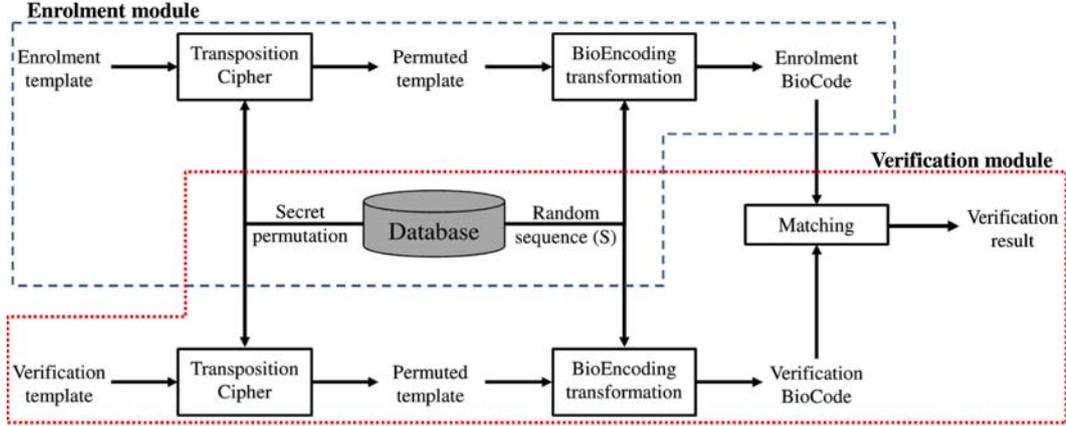


Fig. 5 Block diagram of the enrollment and verification modules of the proposed modified BioEncoding (permutation approach).

the encoded (and ciphered) template, which is computationally infeasible for $m \geq 4$, this random sequence could be XORed with the reversed template to recover the original template.

Employing the above approaches would make inverting the protected BioCodes using record multiplicity attack requires an attacker to perform as many number of trials as required in simple brute-force search attacks at a cost of storing an extra random sequence (in case of XORing approach) or a random permutation key (in case of the permutation approach). It is worth noting that if an adversary gained access to the permutation key (or the random string in the XORing approach), he would be able to obtain the original template from the permuted (or XORed) template since both permutation and XORing operations are reversible. However, as illustrated in Fig. 6, obtaining the permuted (or XORed) template itself from a compromised BioCode is infeasible since the BioEncoding transformation process is non-invertible. That is, disclosing the permutation key would not pose any security threat on the proposed modified BioEncoding method. As a result, similar to base BioEncoding, the proposed modified scheme can still be used without employing user-specific tokens or passwords.

Moreover, permuting the true template randomly (or XORing it with randomly-generated binary string) would diminish the strong correlation between the original features before applying BioEncoding. This would enhance the diversity requirement significantly since applying BioEncoding to the same biometric data after randomization would be as applying BioEncoding to unrelated biometrics templates (i.e. belonging to different users) and hence it would be more difficult for an attacker to use the linkage attack to determine whether two BioCodes are belonging to the same original template.

5. Experiments and Discussion

We have conducted several experiments to validate our security analysis of BioEncoding, testify the effectiveness of

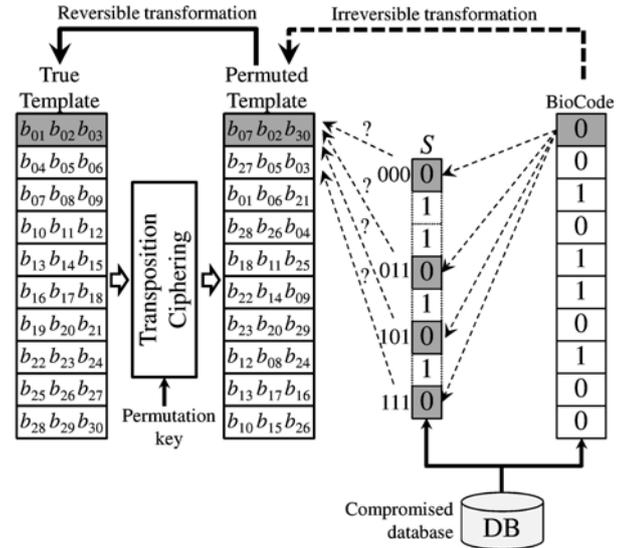


Fig. 6 Obtaining the true template from a compromised BioCode is not feasible even if the permutation key is known to attackers.

the proposed approaches, and investigate the impact of integrating these approaches into BioEncoding on the matching accuracy. The publicly available iris image database collected by the Chinese Academy of Science-Institute of Automation, CASIA-IrisV3-Interval [20] was used in the experiments. This database contains 2639 images captured from 396 different classes (eyes). All the images are 8-bit grayscale images with a resolution of 320×280 pixels.

The objective of the first experiment was to validate our theoretical analysis concerning the resistance of BioEncoding against brute-force attacks. In this experiment, iris codes were generated from all images in the database using the open source MATLAB implementation provided by Masek and Kovese [21] and BioCodes were then derived from all iris codes for different word sizes ($m = 2$ to 16). For each BioCode, zeros and ones are reversed randomly to one of the address words that address '0' and '1' in the corresponding

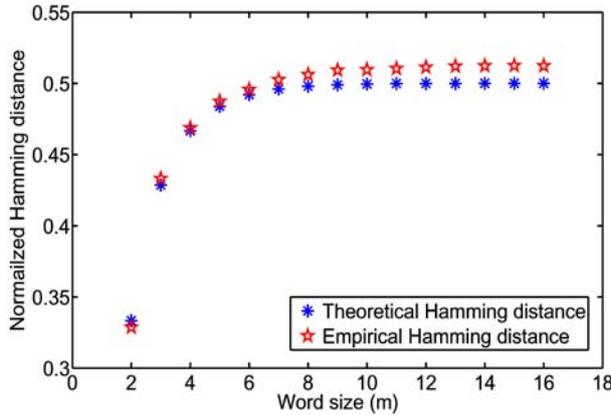


Fig. 7 Hamming distances between true iris codes and code patterns recovered via single brute-force guess for different word lengths.

random sequence, respectively. The Hamming distance between the reversed templates and the true iris codes were then obtained and the average distance is calculated for every m value. This experiment was repeated 100 times using different random sequences and the results were averaged to reduce the statistical fluctuations. Figure 7 shows the empirical (average) normalized Hamming distances for different word sizes (from 2 to 16) along with the expected values computed theoretically using Eq. (5). As shown in the figure, the empirical results are almost identical to the results calculated based on the theoretical interpretation presented in subsection 3.1.1. This conformity between theoretical and empirical results supports our recommendation of avoiding using address words of length $m < 4$.

To measure the robustness of BioEncoding against correlation attacks, we searched for patterns that consist of the corresponding bits in different BioCodes derived from the same iris code (for example, the pattern '010' in Fig. 4) in all the patterns which are composed of corresponding bits in the associated random sequences and the address word of the first matched pattern is selected as the reversed word (for example, the word '101' in Fig. 4) for that pattern. To evaluate the success of this attack, the similarity between the reversed template and the true iris code is measured using the normalized Hamming distance. We tested different word sizes ($m = 3, 4, 5$ and 6) assuming different number of compromised BioCodes (1 to 6) and the whole process is repeated 100 times using different random sequences for each configuration. The results obtained from this experiment, shown in Fig. 8, indicate that the percentage of the recovered bits in true iris codes increases when the number of compromised BioCodes increases. More importantly, the results imply that for BioCodes generated using address words of length m , more than 75% of the true template can be recovered if the number of the compromised BioCodes is m .

The above experiment, using the same setup, was repeated to test the effectiveness of randomly permuting the true template (or XORing it with a random string of same

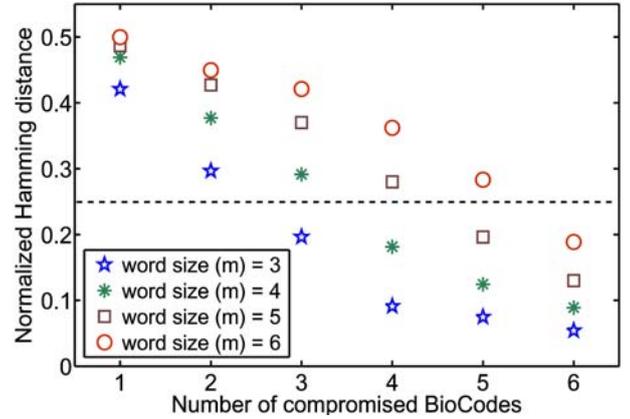


Fig. 8 Hamming distances between true iris codes and code patterns recovered using correlating different number of BioCodes (base BioEncoding).

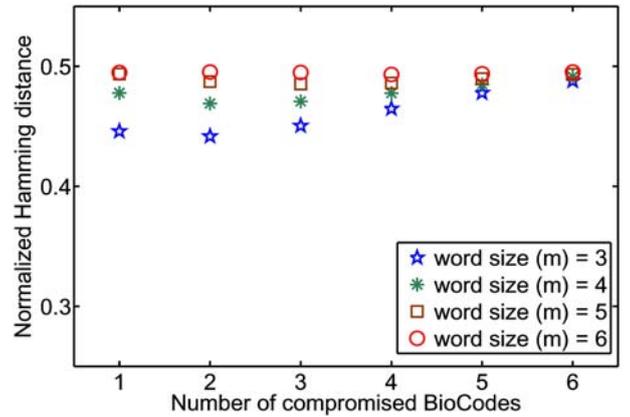


Fig. 9 Hamming distances between true iris codes and code patterns recovered using correlating different number of BioCodes (modified BioEncoding using XORing).

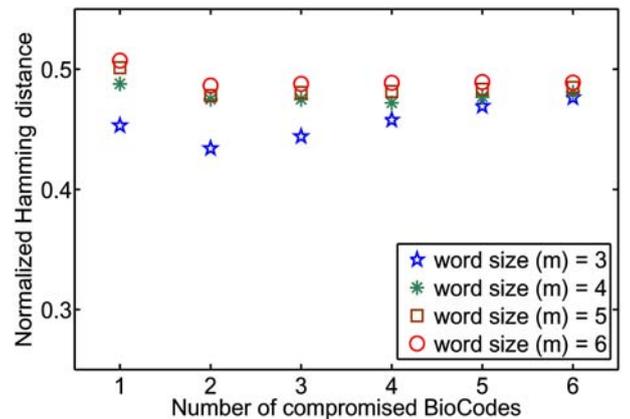


Fig. 10 Hamming distances between true iris codes and code patterns recovered using correlating different number of BioCodes (modified BioEncoding using permutation).

length) before applying the random addressing process of BioEncoding. Figures 9 and 10 show the obtained results for the cases of XORing and permuting, respectively. Slight

differences between the normalized Hamming distances that correspond to one and two compromised BioCodes, respectively, can be seen in these figures. This is most likely due to the statistical fluctuations that result from using random permutation and XORing keys as well as the random sequences employed in the BioEncoding cancelable transformation process. However, in general, both figures show similar results that illustrate the efficacy of integrating the proposed ciphering step in the BioEncoding technique. Using this step, the resistance of BioEncoding against correlation attacks would become very close to its resistance against brute-force search attacks since in every new application of BioEncoding, different patterns are used for the same iris code. Furthermore, correlating more BioCodes derived from different permutations of the same iris code would produce more random templates as the number of BioCodes increases since the probability of the average mismatch between the reversed words and the true words increases with increasing the number of BioCodes employed in the attack until the mismatch reaches 50% (totally random) as shown in Fig. 9 and Fig. 10.

In order to validate the impact of the ciphering step on the diversity of BioCodes derived from the same biometric data, 100 different BioCodes (with word size = 6, and employing different random sequences) were derived from each iris code in the database. We measured the Hamming distances between different pairs of BioCodes generated from the same iris code using base BioEncoding and improved BioEncoding (using permutation). The Hamming distance distributions for both cases are shown in Fig. 11. It is clear from the figure that the standard deviation of the Hamming distances distribution of the modified technique is less than that of the base technique as a result of the permutation step which implies that the dissimilarity between BioCodes derived from same iris codes has increased and hence the diversity property of BioEncoding has been improved.

Finally, to investigate the impact of integrating the ciphering step into BioEncoding on the matching accuracy, a subset that contains 740 iris images (74 classes, 10 images/class) were selected from the adopted dataset. Then, the genuine and imposter matching scores were calculated for unprotected iris codes, BioCodes generated via base BioEncoding and BioCodes derived via the modified BioEncoding (using both approaches: permutation and XORing). All BioCodes were generated using address words of length 5. For genuine comparisons, the first template of each class was matched against the other remaining nine templates of the same eye, and for imposter comparisons, the first template of each class was matched against all templates of all the other classes. The equal error rates (EER) for the four test scenarios are listed in Table 2 and their ROC curves are shown in Fig. 12. The obtained results show that integrating the ciphering step into BioEncoding has no negative impact on the matching accuracy of the protection system.

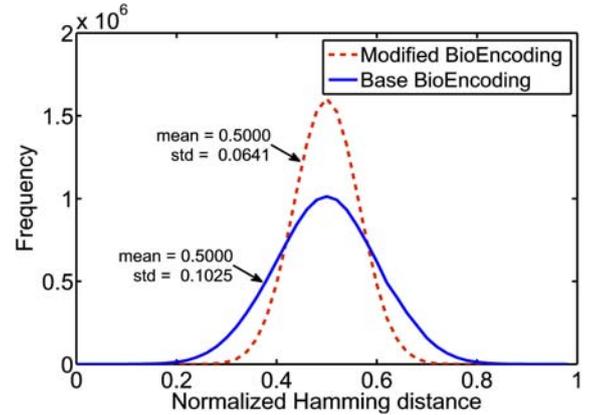


Fig. 11 Distributions of Hamming distances between different BioCodes derived from the same iris codes using base and modified BioEncoding.

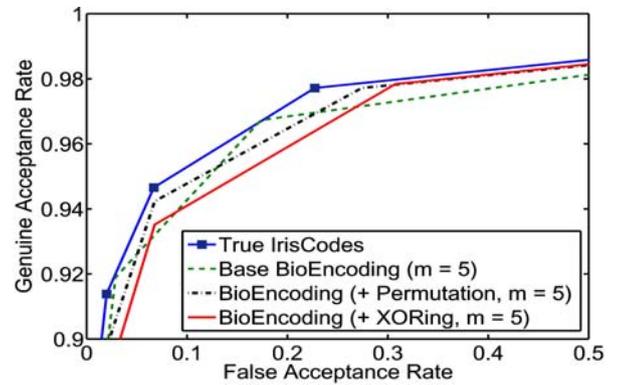


Fig. 12 ROC curves iris codes and BioCodes derived using base and modified BioEncoding.

Table 2 EER (Equal Error Rate) values of iris codes and BioCodes derived using base and modified BioEncoding

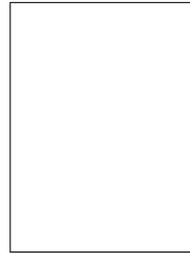
Method	EER
No protection	6.02
Base BioEncoding ($m = 5$)	6.34
Mod. BioEncoding (using permutation, $m = 5$)	6.27
Mod. BioEncoding (using XORing, $m = 5$)	6.63

6. Conclusion

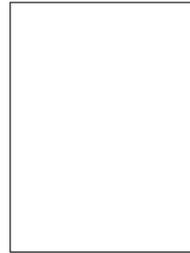
In this paper, the security aspects of a recently proposed cancelable biometrics scheme, BioEncoding, were analyzed with respect to irreversibility and diversity. Three different categories of attacks were investigated: brute-force search attacks, correlation attacks and optimization-based attacks. It has been shown that although BioEncoding is secure against brute-force attacks and optimization-based attacks, it is vulnerable to correlation attacks. We proposed three different approaches to enhance the security of BioEncoding against correlation attacks. Experimental results using CASIA-V3-Interval dataset validated our analysis and confirmed the effectiveness of the proposed modifications to BioEncoding in terms of security and accuracy.

References

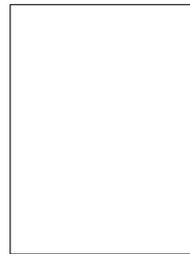
- [1] S. Prabhakar, S. Pankanti, and A.K. Jain, "Biometric recognition: Security and privacy concerns," *IEEE Secur. Privacy Mag.*, vol.1, no.2, pp.33-42, Mar. 2003.
- [2] A.K. Jain, K. Nandakumar, and A. Nagar, "Biometric template security," *EURASIP J. Adv. Signal Process.*, p.579416, 2008.
- [3] A. Juels and M.A. Wattenberg, "A fuzzy commitment scheme," *Proc. 6th. ACM Conf. Computer & Communications Security*, pp.28-36, 1999.
- [4] F. Hao, R. Anderson, and J. Daugman, "Combining crypto with biometrics effectively," *IEEE Trans. Comput.*, vol.55, no.9, pp.1081-1088, Sep. 2006.
- [5] A. Jules and M. Sudan, "A fuzzy vault scheme," *Proc. IEEE Int. Symp. Info. Theory*, p.408, 2002.
- [6] N.K. Ratha, S. Chikkerur, J.H. Connell, and R. Bolle, "Generating cancelable fingerprint templates," *IEEE Trans. on Pattern Anal. Mach. Intell.*, vol.29, no.4, pp.561-572, Apr.2007.
- [7] A.B.J. Teoh, D.C.L. Ngo, and A. Goh, "BioHashing: Two factor authentication featuring fingerprint data and tokenised random number," *Pattern Recogn.*, vol.37, no.11, pp.2245-2255, Nov. 2004.
- [8] O. Ouda, N. Tsumura, and T. Nakaguchi, "BioEncoding: A reliable tokenless cancelable biometrics scheme for protecting IrisCodes," *IEICE Trans. Inf & Syst.*, vol.E93-D, no.7, pp.1878-1888, July 2010.
- [9] A. Cavoukian and A. Stoianov, "Biometric encryption: The new breed of untraceable biometrics." In: N. V. Boulgouris, K. N. Plataniotis, and Micheli-Tzanakou, E. (eds.), *Biometrics: fundamentals, theory, and systems*. Wiley-IEEE Press, 2009
- [10] W.J. Scheirer and T.E. Boulton, "Cracking fuzzy vaults and biometric encryption," *Proc. IEEE Biometrics Symp.*, Baltimore, Md, USA, Sep. 2007.
- [11] C. Soutar, D. Roberge, A. Stoianov, R. Gilroy, and B.V.K.V. Kumar, "Biometric encryption," In: Nichols, R.K. (ed.), *ICSA Guide to Cryptography*, McGraw-Hill New York, 1999
- [12] A. Adler, "Vulnerabilities in biometric encryption systems," *Proc. 5th. Int. Conf. Audio- and Video-Based Biometric Person Authentication (AVBPA)*, LNCS 3546, pp.1100-1109, 2005.
- [13] A. Kholmatov and B. Yanikoglu, "Realization of correlation attack against fuzzy vault scheme," *Proc. SPIE*, vol.6819, pp.1-7, 2008.
- [14] A. Nagar, K. Nandakumar, and A.K. Jain, "Biometric template transformation: a security analysis," *Proc. SPIE*, vol.7541, 2010.
- [15] X. Zhou and T. Kalker, "On the security of BioHashing," *Proc. SPIE*, vol.7541, 2010.
- [16] X. Zhou, S. Wolthusen, C. Busch, and A. Kuijper, "Feature correlation attack on biometric privacy protection schemes," *Proc. 5th. Int. Conf. on Intelligent Information Hiding and Multimedia Signal Processing (IHMSP)*, pp.1061-1065, Kyoto, Japan, Sep. 2009.
- [17] K. Simoons, P. Tuyls, and B. Preneel, "Privacy weaknesses in biometric sketches," *Proc. IEEE Symp. Security and Privacy*, pp.188-203, 2009.
- [18] J. Daugman, "How iris recognition works," *IEEE Trans. Circ. Syst. Video Technol.*, vol.14, no.1, pp.21-30, 2004.
- [19] A. Nagar and A.K. Jain, "On the security of non-invertible fingerprint template transforms," *Proc. IEEE Worksh. Inf. Forens. Secur. (WIFS)*, pp.81-85, Dec. 2009.
- [20] CASIA iris image database, Available: <http://www.cbsr.ia.ac.cn/Databases.htm>.
- [21] L. Masek, P. Kovesi. MATLAB source code for a biometric identification system based on iris patterns. The School of Computer Science and Software Engineering, The University of Western Australia. 2003



Osama Ouda was born in Kuwait, on August 1979. He received the B.Sc. from the faculty of Computers and Information Sciences, Mansoura University, Egypt in 2000 and the M.Sc. from the faculty of Computers and Information Sciences, Ain Shams University, Egypt in 2007. He is currently pursuing the doctoral degree at Chiba University, Chiba, Japan. His scientific interests include information security, biometric encryption and pattern recognition.



Norimichi Tsumura was born in Wakayama, Japan, on April 1967. He received the B.E., M.E. and D.E. in applied physics from Osaka University in 1990, 1992 and 1995, respectively. He moved to the Department of Information and Image Sciences, Chiba University in April 1995, as assistant professor. He is currently associate professor since 2002. He was visiting scientist in University of Rochester from March 1999 to January 2000. He also was researcher at PREST, Japan Science and Technology Corporation (JST) from 2001 to 2003. He got the Optics Prize for Young Scientists (The Optical Society of Japan) in 1995, and Applied Optics Prize for the excellent research and presentation (The Japan Society of Applied Optics) in 2000. He received the Charles E. Ives award in 2002 from the IS&T. He is interested in the color image processing, computer vision, computer graphics and biomedical optics.



Toshiya Nakaguchi was born in Kobe, Japan, on April, 1975. He received the B.E., M.E., and Ph.D. degrees from Sophia University, Tokyo, Japan in 1998, 2000, and 2003, respectively. He was a research fellow supported by Japan Society for the Promotion of Science from April 2001 to March 2003. From 2006 to 2007, he was a research fellow in Center of Excellence in Visceral Biomechanics and Pain, in Aalborg Denmark, supported by CIRIUS, Danish Ministry of Education from 2006 to 2007.

Currently, he is an Assistant Professor of imaging science at the Graduate School of Advanced Integration Science, Chiba University, Chiba Japan. His current research interests include the computer assisted surgery and medical training, medical image analysis, real-time image processing, and image quality evaluation.