# PRACTICAL SECURE BIOMETRICS USING SET INTERSECTION AS A SIMILARITY MEASURE

Daniel Socek, Dubravko Ćulibrk

*CoreTex Systems LLC, 2851 S Ocean Blvd. 5L, Boca Raton, FL 33432, USA*


Vladimir Božović

*Dept. of Mathematical Sciences, Florida Atlantic University, Boca Raton, FL 33431, USA*

Keywords: Biometrics security, fuzzy one-way functions.

Abstract: A novel scheme for securing biometric templates of variable size and order is proposed. The proposed scheme is based on new similarity measure approach, namely the set intersection, which strongly resembles the methodology used in most current state-of-the-art biometrics matching systems. The applicability of the new scheme is compared with that of the existing principal schemes, and it is shown that the new scheme has clear advantages over the existing approaches.

## 1 INTRODUCTION

Identity theft represents the fastest growing type of fraud in the United States (Elbirt, 2005). While identity theft often occurs because of the victim's negligence, it can also occur as a result of direct tampering with the authentication system by a criminal.

Authentication systems based on user's biometric data have several advantages over other authentication methods. The main advantages of biometric-based authentication is the simplicity of use and a limited risk of losing, stealing, or forging users' biological identifiers. On the other hand, the major disadvantage of biometrics-based authentication is the *non-renewability* of biological identifiers. This is a particularly significant issue regarding the identity theft problem.

Biometric-based authentication with the same biometrics is likely to be used in multiple application systems. For example, a fingerprint-based authentication could be used to gain access to multiple systems or facilities. If a biometric template is stolen from a authentication system, criminals can abuse it in the present or future time in multiple venues. In addition, to respect valid privacy concerns by the users, such as corrupt employees at the trusted institutions that have access to a database of biometric templates, the templates should not be stored as plaintext (in its clear form). One solution to the problem is to make use of tamper resistant systems; however, the use of such systems could be infeasible in a given system setup.

Biometric templates often contain condensed discriminatory information about the biometric uniqueness of the user. For instance, in case of fingerprints, the system often stores the discriminatory set of minutiae points. With this information, an adversary can bypass the access control system or extract certain system-specific keys provided that tampering with the system at that level is feasible. In addition, this information could potentially also be used to perform attacks even from the topmost sensor level by creating fake biometric identifiers with the same discriminatory biometric features. For instance, given fingerprint minutiae, an attacker can construct a fake fingerprint that has the same discriminatory information. Methods for creating fake fingerprints such as SFINGE by Cappelli, Miao and Maltoni (Cappelli et al., 2002) or synthetic generation technique by Araque et al. (Araque et al., 2002) can be used for exactly that purpose. Uludag and Jain (Uludag et al., 2004) described many attacks on fingerprint-based identification systems using a fake fingerprint such as rubber or silicon finger, and alike. Similar arguments are also applicable to the other types of biometrics.

Clearly, standard cryptographic one-way primitives are unsuitable for this purpose since the biometric identifiers are fuzzy (not exactly reproducible)

as a result of imperfections in both acquisition and feature extraction methodologies. As a result, several schemes for storing biometric templates securely were proposed recently. In Section 2 of this paper, we present a brief summary of principal work in this area and point out a number of limitations of several state-of-the-art methods for securing biometric templates. In Section 3, we propose a novel approach to securing biometric templates that has several clear advantages over other principal approaches. Finally, conclusions and a number of topics for further research are given in Section 4.

## 2 RELATED WORK

Before describing and analyzing properties of the principal schemes that have been proposed up to date, and also to set the stage for later discussion, several preliminary definitions and concepts are presented next.

### 2.1 Basic Definitions

The design of a scheme for securing biometric templates is constrained with a type of biometric feature vector that is extracted from the sensory information. Properties of feature vectors representing biometric templates heavily depend on the type of biometric data involved, capability of a sensor, and the corresponding feature extraction algorithm. These properties include the types of errors introduced during data acquisition process, as well as the expected range of values and similarity thresholds.

Typically, two types of biometrics templates (feature vectors) often appear in practice: (1) templates with points that have constant size and order, here denoted by *type I* templates, and (2) templates with points having variable size and order, denoted by *type II* templates. For example, type I biometric templates often appear in face recognition systems where feature vectors are singular value decomposition of a face image, or in iris recognition systems such as IrisCode (Hao et al., 2006). Fingerprint and palm print minutiae-based recognition systems, which constitute what are the most common biometric systems (Maltoni et al., 2003) work with type II templates. Schemes for securing biometric templates are in general designed for a particular template type.

In terms of application requirements, there are several types of schemes for securing biometric templates. In work by Dodis et al. (Dodis et al., 2004; Dodis et al., 2006), two types of schemes are defined:

1. *Secure sketch* – This scheme essentially allows for the precise reconstruction of a noisy input. Given an input $x$, the scheme produces a public value $f(x)$, called secure sketch, from which no information about $x$ can be deduced (i.e. $f$ is a one-way function). The scheme can recover the original value of $x$ solely from $f(x)$ and $y$ if and only if $y$ is similar to $x$ according to some similarity measure, denoted with $y \sim x$.

2. *Fuzzy extractor* – For a given input $x$ this scheme produces a public value $f(x)$ and a secret value $k$. Function $f$ is a one-way map so that no information about $x$ can be deduced from $f(x)$. The scheme is able to recover $k$ solely from $y$ and $f(x)$ if and only if $y \sim x$. In practice, $k$ is often used as a secret key for further cryptographic processing.

In (Dodis et al., 2006), it was also shown that it is always possible to construct fuzzy extractors from secure sketches. Intuitively this means that secure sketches comply with a stronger condition (or requirement) than fuzzy extractors do. However, in a number of biometrics-based security applications, even fuzzy extractors comply to a stronger requirement than what suffices in practice.

When concerned with pure verification or identification applications, ability to determine whether a new template matches the stored one is a sufficient requirement. In general, a match is declared when two templates are similar, or, in other words, with similarity measure greater than some threshold $t$ (also referred to as the *similarity bound*). Note that the similarity function is not necessarily a metric. We define a *threshold-based similarity measuring scheme S* to be a scheme that for given one-way transformed value $f(x)$ and a template $y$ determines whether the original template $x$ and $y$ are similar or not:

$$S(f(x),y) = \begin{cases} similar, & \text{if } s(x,y) > t; \\ not\ similar, & \text{if } s(x,y) \leq t, \end{cases}$$

where $s(x,y)$ denotes a similarity measure of $x$ and $y$. Strictly speaking, this kind of scheme is slightly more limited than a scheme that can compute the actual value of $s(x,y)$ from $f(x)$ and $y$; however, almost all biometrics security systems are based on a threshold similarity measure approach.

It is not too difficult to observe that both secure sketches and fuzzy extractors are also threshold-based similarity measuring schemes. It may be of interest to have schemes which are threshold-based similarity measuring schemes that are strictly not secure sketches.

## 2.2 Previously Proposed Schemes

To secure biometric templates of type I, Juels and Wattenberg proposed a scheme called *fuzzy commitment*. This conceptually simple scheme is based on error correcting codes. Let $\mathcal{F}$ be a field, and $\mathcal{C}$ the set of vectors of some $t$-error correcting code. Let $x \in \mathcal{F}^n$ denote a biometric feature vector. Assuming that all codewords lie in $\mathcal{F}^n$, a codeword $c$ is selected uniformly at random from $\mathcal{C}$ and difference $\varepsilon = c - x$ is computed. Next, a suitable one-way function $h$ is selected, and the pair $(\varepsilon, h(c))$ is published, representing the output of fuzzy commitment scheme.

To reconstruct the original feature vector $x$, a similar vector $y$ is required, where the measure of similarity is given by a certain metric. If the usual Hamming distance between $c' = \varepsilon + y$ and $c$ is less than $t$, the error correcting capability of the code $\mathcal{C}$, then it is possible to reconstruct $c$ and consequently $x$. Since the feature vectors are required to be from $\mathcal{F}^n$, the scheme can be applied only to type I feature vectors, where constant size and order is assumed. Fuzzy commitment is a secure sketch scheme. A scheme based on fuzzy vault principle was constructed and successfully applied for securing a particular type of iris templates, called IrisCode, as described in (Hao et al., 2006).

Juels and Sudan in (Juels and Sudan, 2002; Juels and Sudan, 2006) proposed a scheme, called *fuzzy vault*, that slightly extends the applicability of a scheme from (Juels and Wattenberg, 1999) by allowing for the order invariance of feature vector coordinates. This scheme substantially relies on Reed-Solomon error correcting codes, where the codewords are polynomials over a finite field $\mathcal{F}$. Given a feature vector (set) $x \subset \mathcal{F}$ and a secret value $k$, a polynomial $p \in \mathcal{F}[X]$ is selected so that it encodes $k$ in some way (e.g., has an embedding of $k$ in its coefficients). Then an evaluation of the elements of $x$ against $p$ is computed and, along with these points, a number of random *chaff* points that do not lie on $p$ is added to a public collection $R$.

To recover $k$, a set $y$ similar to $x$ must be presented. If $y \sim x$, then $y$ contains many points that lie on $p$. Using error correction procedure, it is possible to reconstruct $p$ exactly, and thereby $k$. If $y$ is not similar to $x$, it does not overlap substantially with $x$ and thus it is not possible to reconstruct $p$ using the error correction mechanism of Reed-Solomon code. By observing the public value $R$, it is infeasible to learn $k$ due to the presence of many chaff points. This is also a secure sketch scheme. While fuzzy vault does allow for a variable order, it does require feature vector sizes to be of the fixed length, thus still not fully

supporting biometrics feature vectors of type II. Several schemes based on fuzzy vault principle were reported for fingerprint data in (Clancy et al., 2003) and (Uludag et al., 2004).

One of the most serious attacks considered for fuzzy vault-based schemes is the *multiple-use attack* that the original authors did not consider in their security model. Under the multiple-use attack, the adversary has public information obtained from multiple authentication systems regarding user $U$. The multiple-use attack is successful if it is possible to compromise the secret information about $U$ (in whole or in part) from analyzing the public information about $U$ from multiple systems. Schemes based on fuzzy vault and generally any schemes that are based on the principle of *chaffing and winnowing* (Rivest, 1998) are weak against multiple-use attack.

Suppose the same user is enrolled in $k > 1$ authentication systems which are all based on the same kind of biometric (e.g. fingerprint) and which all use the fuzzy vault scheme for securing biometric feature vectors. For simplification, let us assume that the user's biometric feature vector in all systems was $x = \{x_1, \ldots, x_t\}$, since almost the same arguments apply when these vectors are *similar*. Recall that the public information that is stored in system $i$ is a collection $R^{(i)}$ that contains $t$ points $(x_1, p^{(i)}(x_1)), \ldots, (x_t, p^{(i)}(t))$ and $m^{(i)}$ chaff points $(r_1^{(i)}, s_1^{(i)}), \ldots, (r_{m^{(i)}}^{(i)}, s_{m^{(i)}}^{(i)})$. According to the fuzzy vault specification chaff points are selected uniformly at random from $\mathcal{U} - x$, where $\mathcal{U}$ denotes the universe of feature vector coordinates. If $R_x^{(i)}$ denotes the restriction of $R^{(i)}$ to the x-axis, then

$$\lim_{k \to \infty}(R_x^{(1)} \cap R_x^{(2)} \cap \ldots \cap R_x^{(k)}) = x$$

unless chaff points always entirely cover the remaining universe $\mathcal{U} - x$ or some fixed parts of it. Moreover, if we take a simple case when $r = |R_x^{(i)}| - t \ll |\mathcal{U}|$ for $i = 1, 2$, then

$$Prob(R_x^{(1)} \cap R_x^{(2)} = x) = \frac{\binom{|\mathcal{U}|-t-r}{r}}{\binom{|\mathcal{U}|-t}{r}+1} \approx 1,$$

where $|\mathcal{U}|$ denotes the cardinality of set $\mathcal{U}$. In other words, if the number of randomly selected chaff points is much smaller than the size of the universe $\mathcal{U}$, the intersection of chaff points of the same person taken from two authentication systems will almost certainly be empty.

In (Juels and Sudan, 2002; Juels and Sudan, 2006) it is shown that the number of different polynomials that agree on $t$ is small if the size of collection $R$ is small. Thus, in order to ensure security from that

point of view, the authors recommend taking a large number of chaff points. Yet, the authors do not *require* to always cover the entire remaining universe $\mathcal{U} - x$ with chaff. Indeed, this is probably infeasible when dealing with larger universes. However, to avoid the multiple-use attack as described here, the entire remaining universe or fixed part of it must be covered by chaff. That is, $R_x^{(i)} = \mathcal{U}'$ for all $i$ where $\mathcal{U}'$ is a subset of $\mathcal{U}$ (likely $\mathcal{U}' = \mathcal{U}$) that provides a large number of polynomials that agree on $t$ points and also a computationally infeasible search space.

In (Boyen, 2004), Boyen considered the issues of multiple uses of the same fuzzy secret in a general fuzzy extractor scheme. Boyen pointed out that in the security model of fuzzy extractors such issue must be addressed and related security risks accounted for.

Dodis et al. in (Dodis et al., 2004; Dodis et al., 2006) proposed a scheme that allows for securing biometric feature vectors of type II. This scheme, called *PinSketch*, relies on *t*-error correcting (BCH) code *C*. In order to simplify description, let us assume $H$ to be a parity check matrix of the code $C$ over some finite field $\mathcal{F}$. For a given feature vector $x$ which belongs to $\mathcal{F}^n$, the scheme computes output $syn(x) = Hx$, which is referred to as the *syndrome* of vector $x$.

In the reconstruction phase, $syn(y)$ is computed for a given vector $y$. Let $\delta = syn(x) - syn(y)$. It is easy to see that there exists at most one vector $v$ such that $syn(v) = \delta$ and $weight(v) \leq t$. One of the nice features of binary BCH codes is possibility of computing $supp(v)$ given $syn(v)$ and vice versa, where $supp(v)$ represents the listing of positions where $v$ has nonzero coordinate. Computing of $supp(v)$ for a given $syn(v)$ is the key step in the reconstruction phase. If a distance metric $d(x,y) \leq t$ then $supp(v) = x \triangle y$, and in that case the original set could be reconstructed by $x = y \triangle supp(x)$. PinSketch is a secure sketch scheme that supports biometrics feature vectors of type II.

## 2.3 Applicability Critique of Error Correcting-based Schemes for Type II Templates

From the mathematical point view, the most suitable method for measuring *similarity* between two sets is by their symmetric set difference. However, this quite reasonable mathematical choice is often a limitation for practical use. Let us try to illustrate this problem in the case where it is needed to measure closeness between two sets $A$ and $B$ that represent biometric (fingerprint) personal data, of not necessarily different persons. This is an inevitable step in the process of verification or identification. Reconstruction of $A$, using *similar* set $B$ will be successful if and only if

$|A \triangle B| \leq t$, where $t$ is a given parameter that controls the closeness between sets. It seems that error correcting codes are a suitable choice for reconstructing $A$ from a noisy input $B$. Here, $t$ is the error correcting bound of the chosen code.

We argue that the use of error correcting codes and consequently the Hamming distance as a measure of similarity between type II feature vectors is not an adequate choice. For instance, in the PinSketch scheme (Dodis et al., 2006), templates are represented as characteristic vectors with respect to universe $\mathcal{U}$. Therefore, the symmetric difference is simply related to the Hamming distance between characteristic vectors. In a typical application of PinSketch, such as fingerprint identification, the scheme has a substantial applicability issue. The number of minutiae, according to many statistical analyses of fingerprints lies with high probability in the interval between 20 and 80 (Amengual et al., 1997). Thus, choice of the error correcting bound $t$ that is used in this scheme seems to be its main shortcoming.

Considering that size of the universe is not large, $t$ must be chosen in a way not to compromise security. For instance, if a template set is of size 15, then setting $t > 12$ would not be an adequate choice, since an adversary could test all elements or 2-subsets of the universe (which is feasible for a universe of fingerprint minutiae) and use error correction to obtain the template set. On the other hand $t$ must be set to provide proper authentication. Due to imperfections in the template extraction it is common to have spurious minutiae and some real minutiae that are not recognized. Thus, symmetric difference between newly presented and stored template could became relatively large, yet the intersection could still be large enough for authentication of $B$ as $A$ with high confidence. For example, suppose $|A| = 20$ and $|\mathcal{U}| \approx 10^6$ with possibly nonuniform distribution. Therefore, $t$ could be at most 17. If we accept *twelve point matching rule* as valid, and if $|B| = 22$ and $|A \cap B| = 12$ then $B$ will not be authenticated as $A$ although intersection is large enough to confirm the identity. Even if do not accept *twelve point matching rule*, it is possible to construct many examples where symmetric difference does not appear as an adequate choice for *similarity* measure. In most minutia-based authentication systems similarity is measured using the number of points that agree in the best possible alignment of two sets of minutiae using translation, rotation and potentially scaling. Therefore, the set intersection is a more appropriate similarity measure in practice.

The authors of fuzzy vault (Juels and Sudan, 2002; Juels and Sudan, 2006) indicated that the scheme is applicable to feature vectors with fixed size and vari-

able ordering which limits the practical use of the scheme to type I vectors. Even if it is possible to extend the fuzzy vault scheme to work with the type II feature vectors, the scheme would face the similar applicability issues since it is based on error correction approach. As an artifact of fuzzy vault where the entire universe is covered by chaff due to multiple use attack and the requirement about the minimal number of different polynomials that agree on $t$ points, the similarity measure is not achieved with symmetric set difference but with ordinary set difference $B - A$. This slightly better scenario is still inappropriate since it is possible to have cases where both $A \cap B$ and $B - A$ are relatively large, in which case the fuzzy vault scheme would give a false rejection.

In this work we design a scalable secure scheme applicable to type II biometric templates, such as fingerprint minutiae which are currently the most common biometric templates (Maltoni et al., 2003).

# 3 OUR APPROACH

Let $\mathcal{F}$ be a finite field. Given an encoding of biometric templates into the field $\mathcal{F}$, it is common to denote $\mathcal{F}$ as the *universe* $\mathcal{U}$. In this setting, biometric templates correspond to the subsets of $\mathcal{U}$. The key observation is that the size of the universe is typically much larger than the size of a subset representing a biometric template, but still in a range that allows feasible exhaustive search. For instance, the size of the universe representing fingerprint minutiae is approximately in the range of $10^5$-$10^7$, depending on technical characteristics of the sensor, yet the size of a biometric template is between 20 and 80 with high probability. In further analysis, we will assume $|\mathcal{U}| \gg |A|$, where $A$ represents a template set.

Accuracy of the extraction of biometric data depends on several factors, but mostly on the sensory technology for data acquisition and image processing algorithms for biometric template extraction. Due to these imperfections, it cannot be expected that newly submitted templates perfectly match the stored ones. It is not uncommon to have, under certain scenarios, just part of the fingerprint that needs to be identified. Therefore, a scheme for secure authentication needs to have a necessary level of tolerance with respect to possible incompleteness and inaccuracy of submitted templates. The tolerance threshold for our scheme can be easily customized regarding the particular application.

## 3.1 Scheme Description

Let $\mathcal{G}$ be a finite field where $|\mathcal{G}| = p^k$, assuming that $p^k$ provides a large keyspace, e.g. $p^k > 2^{100}$. Let $m_1$ and $m_2$ be integers such that $m_1 \leq |A| \leq m_2$ for all subsets $A$ representing biometric templates. Suppose that $\ell$ is an integer chosen in so that

$$\binom{m_2}{\ell} \leq 2^{k_1} \ll 2^{k_2} \leq \binom{|\mathcal{U}|}{\ell}.$$

In general, it is required for $k_1$ to be small enough to allow for a feasible search through the set of $\ell$-subsets of any given template $A$. On the other hand, it is required for $k_2$ to be large enough, making it infeasible to search through all $\ell$-subsets of the universe $\mathcal{U}$. As an illustration, under the assumption that the distribution of points of $A$ is uniform over $\mathcal{U}$, if $|\mathcal{U}| \approx 10^6$ and $m_2 = 100$, even with a choice of $\ell = 3$ the size of $\binom{|\mathcal{U}|}{\ell}$ is approximately $2^{60}$ which is a larger search space than that of DES. For the same parameters, the size of $\binom{m_2}{\ell}$ is just 161700. The generation of public one-way transformation of the given template in the proposed scheme is as follows:

1. Let $A = \{a_1, a_2, \ldots, a_n\}$ the input biometric template. Randomly choose $s \in \mathcal{G}$ and using an $\ell$-out-of-$n$ perfect secret sharing scheme, create $n$ shares of $s$ denoted by $s_1, \ldots, s_n$.

2. Choose a secure cryptographic hash function $h$ and obtain set $\{h(sa_1), h(sa_2), \ldots, h(sa_n)\}$, where $sa_i$ means concatenation of $s$ and $a_i$. It is required that the chosen hash function is both preimage resistant and collision-resistant.

3. Define a discrete function $f_A : \mathcal{U} \rightarrow \mathcal{G}$ in the following way

$$f_A(x) = \begin{cases} s_i, & \text{if } x = a_i; \\ y_x, & \text{if } x \notin A, \end{cases}$$

where the values $y_x$ are chosen uniformly at random.

4. Store $f_A(x)$, $H_A = \{h(sa_1), h(sa_2), \ldots, h(sa_n)\}$ and $h(s)$ as a one-way public transformation of $A$.

The recovery process in our scheme is performed in the following way:

1. For a given set $B = \{b_1, \ldots, b_m\}$, for all $\ell$-subsets of $B$, denoted by $B_1, \ldots, B_{\binom{m}{\ell}}$, do the following:

   (a) Evaluate $f_A(B_i)$.

   (b) Using the reconstruction method provided by the secret sharing scheme, obtain $s'$ from $f_A(B_i)$.

(c) Compute $h(s')$; if $h(s') = h(s)$, then assume $s = s'$, compute $H_B = \{h(s'b_1), \ldots, h(s'b_m)\}$, and then output $|H_A \cap H_B|\} = |A \cap B| \geq \ell$ and terminate.

2. If for all $\ell$-subsets of $B$ no termination was reached, output $|A \cap B| < \ell$ and terminate.

In our scheme, $s$ corresponds to the extracted key from the definition of fuzzy extractor. Moreover, with minor modifications the proposed scheme can also be turned into a secure sketch scheme where original set $A$ can be completely reproduced. The algorithm determines a threshold-based similarity of templates $A$ and $B$ using set intersection as a similarity measure, which reflects the same principle used in most minutia-based recognition methods. The algorithm outputs $|A \cap B|$ if $|A \cap B| \geq \ell$. Once $|A \cap B|$ has been obtained, it is to be decided if the authentication threshold has been achieved.

The authentication bound is not substantially involved in our scheme, which is not the case in the previous schemes. The only requirement related to the authentication bound is that it must be greater than or equal to the security bound $\ell$.

## 3.2 Security and Applicability Aspects

To address the security of our method, it is essential to discuss issues regarding the distribution of the source data. The attacker's goal is to learn information about the original template $A$ given only the public values $f_A(x)$, $H_A$ and $h(S)$. Note that the multiple use attack is not applicable to our scheme since the entire universe $\mathcal{U}$ is covered by uniformly random values according to $f_A$.

A reasonable question that arises from the analysis of the proposed scheme is how the assumption of strictly uniform distribution could be relaxed for some practical applications. We show how to set parameters of our scheme in the case of fingerprint authentication.

In our scheme, for the enrollment template $A$ and a probe $B$ that originates from the same subject as $A$, we assume that $|A \cap B| = \lceil t|A| \rceil$ for $t \in (0, 1)$.

Let $X$ be a random variable that describes the number of unsuccessful attempts before getting a qualified subset, i.e. a set from $A \cap B$. Clearly $X$ has a negative hypergeometric distribution. If $a^{(b)} = a(a-1) \cdots (a-b+1)$ than the distribution of $X$ is

$$Prob(X = r) = \frac{bw^{(r-1)}}{c^{(r)}},$$

where $b = \binom{\lceil t|A| \rceil}{\ell}$, $c = \binom{|B|}{\ell}$ and $w = c - b$.

Then, the mathematical expectation of $X$ is given by

$$EX = \frac{c+1}{b+1}.$$

Next, we show some concrete parameters that give a clear view of the computational complexity of the searching process for an $\ell$-subset in $A \cap B$. In Table 1 we fix parameter $t = 0.5$ and for simplicity, we fix the sizes of $A$ and $B$ to be equal although this is not required by our construction.

Table 1: The expected number of attempts needed to find an $\ell$-subset of $A \cap B$ for various sizes of $A$ and $B$ when $t = 0.5$.

| $|A| = |B| = 80$ | | | | |
| --- | --- | --- | --- | --- |
| $\ell$ | 8 | 10 | 12 | 15 |
| EX | 377 | 1943 | 10784 | 164968 |

| $|A| = |B| = 60$ | | | | |
| --- | --- | --- | --- | --- |
| $\ell$ | 8 | 10 | 12 | 15 |
| EX | 438 | 2510 | 16179 | 342928 |

| $|A| = |B| = 40$ | | | | |
| --- | --- | --- | --- | --- |
| $\ell$ | 8 | 10 | 12 | 15 |
| EX | 611 | 4588 | 44351 | 2594347 |

| $|A| = |B| = 30$ | | | | |
| --- | --- | --- | --- | --- |
| $\ell$ | 8 | 10 | 12 | 15 |
| EX | 910 | 10002 | 189679 | 77558761 |

If we set $t$ to be slightly higher, for example $t = 0.6$, then the expected values significantly change, as depicted in Table 2. For many authentication systems it is not unreasonable to expect that set $B$, which originates from the same subject as $A$, have at least 60% common points with $A$.

Table 2: The expected number of attempts needed to find an $\ell$-subset of $A \cap B$ for various sizes of $A$ and $B$ when $t = 0.6$.

| $|A| = |B| = 80$ | | | | |
| --- | --- | --- | --- | --- |
| $\ell$ | 8 | 10 | 12 | 15 |
| EX | 77 | 252 | 865 | 6070 |

| $|A| = |B| = 60$ | | | | |
| --- | --- | --- | --- | --- |
| $\ell$ | 8 | 10 | 12 | 15 |
| EX | 85 | 297 | 1118 | 9554 |

| $|A| = |B| = 40$ | | | | |
| --- | --- | --- | --- | --- |
| $\ell$ | 8 | 10 | 12 | 15 |
| EX | 105 | 433 | 2067 | 30765 |

| $|A| = |B| = 30$ | | | | |
| --- | --- | --- | --- | --- |
| $\ell$ | 8 | 10 | 12 | 15 |
| EX | 134 | 687 | 4659 | 189863 |

Although parameter $t$ is not included in the construction of the scheme, it is useful to have a presumption on the expectation for $t$. Taking into consideration the particular application and by doing a preliminary statistical analysis on the accuracy of the template extraction system, an estimation for $t$ can be achieved. When higher level of security is required, $t$ generally must be higher. Consequently, it is possible to choose larger $\ell$ and still have a high efficiency in the task of finding $\ell$-subset from $A \cap B$.

For instance, for certain high-security authentication, the threshold of common points between the new and stored template could be set to at least 80% of the stored template set. In that case, even setting $\ell \geq 20$ results in efficient performance of our scheme. Table 3 shows the case when $\ell = 20$.

Table 3: The expected number of attempts needed to find 20-subset of $A \cap B$ when $t = 0.8$ and $|A| = |B| = n$.

| $n$ | 30 | 40 | 60 | 80 |
|-----|------|-----|-----|-----|
| $EX$ | 2828 | 611 | 251 | 181 |

There have been a number of attempts to explain the minutiae distribution. Most recent papers tracking this subject come from the Michigan State University group which mainly dealt with the questions of individuality of fingerprints and how similar two randomly chosen fingerprint templates could be. This problem was partially inspired by a recent challenge to the generally accepted *twelve points matching rule* in some US courts.

The statistical model of distribution of minutiae points has not been established due to very complex nature of the problem. The distribution of minutiae that has been proposed in (Dass et al., 2005) is a so-called *mixed distribution*. This distribution appears to be more appropriate than the uniform distribution regarding the statistical data collection taken from three large publicly available databases of fingerprints (Dass et al., 2005). However, note that all results heavily depend on the quality of acquired fingerprint data and the extraction method used in the experiments.

The result which could be of particular importance for our security model is a result about the probability that two random fingerprint templates of 36 minutiae share more than 12 points. If $P(36, 36, 12)$ denotes this probability and assuming the mixed distribution, it can be shown that $P(36, 36, 12) \approx 6 \times 10^{-7}$. In our scheme, if $\ell = 12$ then an attacker could try to get stored set $A$ of 36 minutiae by choosing a random subset $B$ of 36 elements of the universe $\mathcal{U}$, hoping that $|A \cap B| \geq 12$. However, the only way the attacker can know if the chosen subset $B$ contains more than 12 elements of the stored template $A$ is by running through all 12-subsets of $B$. Thus, the probability of an attacker's success is $\approx 6 \times 10^{-7} \times \frac{1}{\binom{36}{12}} \approx 2^{-56}$. That makes this kind of attack inefficient especially if we set $\ell$ to be higher than 12.

We would like to stress that the previously mentioned results are dependent on the effectiveness of the automated minutiae extraction methods which are only of moderate reliability.

It must be understood that the nonuniformity of the universe of certain biometrics influences all proposed schemes regarding security issues. For the schemes based on error correction codes, nonuniformity affects the error correction bound. Consequently, it produces an increase of the False Rejection Rate (FRR). In our scheme, it induces an increase of the parameter $\ell$ that causes a higher computational cost.

# 4 CONCLUSIONS

We proposed a novel scheme for securing biometric templates of variable size and order. Unlike previously proposed schemes, our scheme uses set intersection as the similarity measure between the enrollment template and a probe. This principle reflects matching criteria used in most minutia-based authentication systems, and as such offers better applicability than the schemes based on error correcting approach. We showed that the scheme is scalable and has a relaxed dependency on the similarity bound. Finally we demonstrated how to set the parameters of the proposed scheme in order to achieve both high security and broad applicability even when the minutiae distribution is nonuniform.

# ACKNOWLEDGEMENTS

# REFERENCES

Amengual, J., Juan, A., Pérez, J., Prat, F., Sáez, S., and Vilar, J. (1997). Real-time minutiae extraction in fingerprint images. In *International Conference on Image Processing and Its Applications (IPA97)*, volume 2, pages 871–875.

Araque, J., Baena, M., Chalela, B., Navarro, D., and Vizcaya, P. (2002). Synthesis of fingerprint images. In *16th International Conference on Pattern Recognition*, volume 2, pages 422–425.

Boyen, X. (2004). Reusable cryptographic fuzzy extractors. In *ACM Conference on Computer and Communications Security (CCS 2004)*, pages 82–91. New-York: ACM Press.

Cappelli, R., Maio, D., and D.Maltoni (2002). Synthetic fingerprint-database generation. In *16th International Conference on Pattern Recognition*, volume 3, pages 744–747.

Clancy, T. C., Kiyavash, N., and Lin, D. J. (2003). Secure smartcard-based fingerprint authentication. In *ACM SIGMM Workshop on Biometrics Methods and Applications (WBMA '03)*, pages 45–52. ACM Press.

Dass, S. C., Zhu, Y., and Jain, A. K. (2005). Statistical models for assessing the individuality of fingerprints. In *IEEE Workshop on Automatic Identification Advanced Technologies (AUTOID '05)*, pages 3–9. IEEE Computer Society.

Dodis, Y., Reyzin, L., and Smith, A. (2004). Fuzzy extractors: How to generate strong keys from biometrics and other noisy data. In *EUROCRYPT 2004, Interlaken, Switzerland*, pages 523–540.

Dodis, Y., Ostrovsky, R., Reyzin, L., and Smith, A. (April 28, 2006). Fuzzy extractors: How to generate strong keys from biometrics and other noisy data. Retrieved April 4, 2007, from http://www.citebase.org/abstract?id=oai:arXiv.org:cs/0602007

Elbirt, A. J. (Summer 2005). Who are you? how to protect against identity theft. *IEEE Technology and Society Magazine*.

Hao, F., Anderson, R., and Daugman, J. (2006). Combining crypto with biometrics effectively. *IEEE Transactions on Computers*, 55(9):1081–1088.

Juels, A. and Sudan, M. (2002). A fuzzy vault scheme. In *IEEE International Symposium on Information Theory (ISIT 2002), Lausanne, Switzerland*.

Juels, A. and Sudan, M. (2006). A fuzzy vault scheme. *Designs, Codes and Cryptography*, 38(2):237–257.

Juels, A. and Wattenberg, M. (1999). A fuzzy commitment scheme. In *ACM Conference on Computer and Communications Security*, pages 28–36.

Maltoni, D., Maio, D., Jain, A. K., and Prabhakar, S. (2003). *Handbook of Fingerprint Recognition*. Springer-Verlag.

Rivest, R. L. (April 24, 1998). Chaffing and winnowing: Confidentiality without encryption. Retrieved April 4, 2007, from http://theory.lcs.mit.edu/~rivest/chaffing.txt

Uludag, U., Pankanti, S., Prabhakar, S., and Jain, A. (2004). Biometric cryptosystems: Issues and challenges. *IEEE Special Issue on Enabling Security Technologies for Digital Rights Management*, 92(6):948–960.