

1 *Type of the Paper (Article)*

## 2 **Security and Cryptographic Challenges for** 3 **Authentication based on Biometrics Data**

4 **Stefania Loredana Nita<sup>1</sup>, Marius Iulian Mihailescu<sup>2,\*</sup> and Valentin Corneliu Pau<sup>3</sup>**

5 <sup>1</sup> University of Bucharest, Department of Computer Science; [stefanialoredanani@gmail.com](mailto:stefanialoredanani@gmail.com)

6 <sup>2</sup> RCCL, Miami, Florida; [marius.mihailescu@hotmail.com](mailto:marius.mihailescu@hotmail.com)

7 <sup>3</sup> The Academy of Romanian Scientists, Bucharest, Romania; [v\\_pau@utm.ro](mailto:v_pau@utm.ro)

8 \* Correspondence: [marius.mihailescu@hotmail.com](mailto:marius.mihailescu@hotmail.com); Tel.: +40740030310

9

10 **Abstract:** Authentication systems based on biometrics characteristics and data represents one of the  
11 most important trend in the evolution of our world. In the near future, biometrics systems will be  
12 everywhere in the society, such as government, education, smart cities, banks etc. Due to its  
13 uniqueness characteristic, biometrics systems will become also vulnerable, privacy being one of the  
14 most important challenge. The classic cryptographic primitives are not sufficient to assure a strong  
15 level of secureness for privacy. The following work paper represents an effort to present the main  
16 cryptographic techniques and algorithms that can give us the possibility to raise a certain level of  
17 secureness for privacy. We will show their own challenges (strengths and weaknesses). We will  
18 demonstrate how we can use the most common and well-known techniques and algorithms in order  
19 to get a maximum efficiency and a high level in assuring the integrity of the biometrics data.

20 **Keywords:** classification, machine learning, chaos-based cryptography, Hadoop, data clustering,  
21 biometrics.

22

---

### 23 **1. Introduction**

24 Biometric represents a science through which a system can identify in a unique way the  
25 individual based on his physiological (face, fingerprint, iris, hand geometry, retina etc.) and  
26 behavioral (voice, gait, signature, keystroke etc.) traits. Due to the rich of applications and systems  
27 that occur at this moment on the market, the authentication technology is very widespread from e-  
28 commerce applications, door access, and smart city technologies to Internet-of-Things technologies  
29 and applications. Once these solutions were declared as providing better security in authentication,  
30 automatically issues with privacy preserving and data integrity has been found didn't stop to occur.  
31 Indeed, biometrics systems become more convenient to be used as compared with different classic  
32 authentication systems such as token based (e.g. ID cards) or knowledge based (e.g. passwords) [12].  
33 A biometric-based authentication system is using two different operating modes: identification and  
34 verification modes. These two different modes are very important as they represent security gaps,  
35 which are exploited by unwanted users (e.g. hackers/crackers).

36 *What is happening in identification mode?*

37 In identification mode, the system has a clear goal to carry out the one-to-many comparison  
38 meant to set up the individual identity. With other words, the user's identity and the templates that  
39 are stored in the database are compared and based on the result a decision will be made. The purpose  
40 of the identification process is to answer to the question "Who am I?" If we are talking about the

41 implementation process, we can say that they are time consuming for deploying and needs a huge  
42 amount of time for processing in order to find the proper match within the database [12].

43 *What is happening in verification mode?*

44 In verification mode, the system is carrying out the one-to-one comparison. This type o  
45 comparison is used to set up the individual identity. The user is claiming that identity and the system  
46 has the duty to verify if the claim is genuine or forged. The goal of the verification process is to answer  
47 to the question "Am I who I say I am?" [12].

48 An interesting question raised by different professionals is "What measurements for the  
49 biological characteristics is making them to be qualified as a biometric?" Most of the human  
50 physiological and/or behavioral characteristics used in the system for authentication as long as they  
51 are satisfying the following requirements:

- 52 - *Performance*: this requirement is quite important as the characteristic should be enough invariant.  
53 The respect has to be assure for the matching criterion over a period.
- 54 - *Distinctiveness*: by choosing two persons should be sufficient different in terms of the  
55 characteristic.
- 56 - *Universality*: the criteria consist in its unique characteristic that has to be for each person.
- 57 - *Collectability*: the requirement is a metric that is quantitatively measured.

58 When we are dealing with a real life biometric system, there is a number of issues that has to be  
59 taken into consideration in order to take a complete advantage of the full system and to combat those  
60 security issues that could occur on different section:

- 61 - *Performance*, the accuracy and speed, two main characteristics that refers to the achievable  
62 recognition, are required to achieve the desired recognition accuracy and speed. Also,  
63 operational and environmental factors are affecting the accuracy and speed;
- 64 - *Acceptability*, a factor that will indicates which people are willing to accept the use of a particular  
65 biometric identity in terms of characteristic using in a daily lives;
- 66 - *Circumvention* reflects how easy is to fool a system using different methods meant to steal data  
67 and to corrupt the integrity of the data.

68 As we will be able to see, the following content of the present work paper will cover the most  
69 important aspects of security flaws that are raised by each of the components of the biometric system.  
70 The work paper will show and demonstrate theoretical and practical two main benefits: (1) we will  
71 demonstrate how we can protect the integrity of the biometric data using Machine Learning  
72 Classification and what benefits we can obtain. Another benefit, and (2) applying chaos-based  
73 cryptography over encrypted biometric encrypted data. The two main method for secure and  
74 guarantee the integrity of the biometric encrypted data will be demonstrated in a professional system  
75 architecture based on Hadoop and Data clustering.

76 The new challenges from the last two years represents a very important alarm signal for both,  
77 academy and business environments. The security threats and gaps in assuring the privacy and  
78 integrity found, represents one of the most important occasion from which we have to take the  
79 maximum advantage in creating new theoretical and practical security frameworks.

80 Below we have underline the main new challenges that in our opinion will create new research  
81 directions for security and cryptography field, such as:

- 82 - Using Machine Learning Classification over the encrypted biometric data;
- 83 - Encryption of biometric data in a Data Clustering environment;

84 - Encryption of biometric data using Chaos-based cryptography;  
 85 In order to be able to follow the ideas and to understand how they will be applied in a real  
 86 environment, we need to understand main four modules of the biometric system and their  
 87 vulnerabilities.

88 The mentioned challenges raised above will be treated and the issues solved using Chaos-based  
 89 cryptography and machine-learning classification applied in cluster environment using  
 90 authentication based biometrics.

## 91 2. Machine Learning Classification Over Biometric Encrypted Data

92 In this section, we present two machine-learning techniques applied on biometric encrypted data:  
 93 hyperplane decision and Naïve Bayes.

### 94 2.1. Preliminaries

#### 95 2.1.1. Machine learning

96 “Machine learning is a field of computer science that gives the ability to learn without being explicitly  
 97 programmed” – Samuel Arthur [1].

98 The below techniques have been applied and implemented in a software solution in order to simulate the methods  
 99 on biometric data. The learning techniques implemented are described below. The solution software was  
 100 implemented in .NET Framework 4.5 using C# and Microsoft SQL Server 2016. Due to the status of the  
 101 software application, we cannot provide the source code in this work paper. For those who are interested about  
 102 the application they can visit the web page <https://www.researchgate.net/project/Biometrics-Analysis-Tool>.  
 103 [Tool](https://www.researchgate.net/project/Biometrics-Analysis-Tool).

104 There are four types of learning techniques implemented are:

- 105 - *Supervised learning* is a type of inductive learning based on training sets, in which, the agent  
 106 receives a set of inputs and their corresponding outputs. The task of the agent is to learn the links  
 107 between every input and its corresponding output and to generate a template function that will  
 108 be able to solve problems for new inputs.
- 109 - *Semi-supervised learning*. In this type of learning, the agent receives an incomplete training set.
- 110 - *Unsupervised learning* is not using training sets, but the agent needs to discover on its own  
 111 different patterns in dataset.
- 112 - *Reinforcement learning* is a type of learning in which the training data is given as feedback for the  
 113 agent, such that if its output is “good” it receives a reward, otherwise it receives a punishment.  
 114 The target of the agent is to maximize its reward, providing better and better outputs. The  
 115 meaning of “good” output is different depending on the environment in which the agent is used  
 116 [4].

117 *Classification* represents a machine learning technique (included in supervised learning) in which the  
 118 inputs are divided into two or more classes. The input is a *feature vector*  $v = (v_1, \dots, v_n) \in \mathbb{R}^n$  that  
 119 will be classified by applying a classification function  $f_m: \mathbb{R}^n \rightarrow \{x_1, \dots, x_c\}$  on  $v$ , and the output is  
 120  $x_{c^*} = f_m(v)$ , where  $c^* \in \{1, \dots, c\}$ ;  $x_{c^*}$  is the class in which  $v$  falls, based on model  $m$ .

121 In this case, the feature vector will represent all the biometric data over which the classification data  
 122 is applied.

$$123 \quad v = (v_1, \dots, v_n) \in \mathbb{R}^n \quad (1)$$

$$124 \quad \text{biometric vector} = (b_v) = (b_{v_1}, \dots, b_{v_n})$$

$$125 \quad f_m: \mathbb{R}^n \rightarrow \{x_1, \dots, x_c\} \text{ over } b_v \text{ and the output is } x_{c^*} = f_m(v) \text{ where } c^* \in \{1, \dots, c\}$$

126 The  $x_{c^*}$  represents the class in which  $b_v$  falls, being based on model  $m$ .  
 127 Two important classification algorithms are Naïve Bayes and hyperplane decision-based classifier.  
 128 *Naïve Bayes*. The model  $m$  of this classifier is based on probabilities: the probability that class  $x_i$   
 129 occurs is  $\{p(X = x_i)\}_{i=1}^c$  and the probability that the element  $v_j$  of  $v$  occurs in the particular class  
 130  $x_i$ . The classification function is:

$$131 \quad x_{c^*} = \max_{i \in \{1, \dots, c\}} p(X = x_i | V = v) \quad (2)$$

$$132 \quad = \max_{i \in \{1, \dots, c\}} p(X = x_i, V_1 = v_1, \dots, V_n = v_n)$$

133 In order to obtain the second equality Bayes rules was applied.  
 134 In biometrics, the Naïve Bayes will help us to classify the biometric data and to obtain a faster time  
 135 for processing. The classification is done based on the biometric characteristics.  
 136 *Hyperplane decision*. The model  $m$  contains  $c$  vectors in  $\mathbb{R}^n$  ( $m = \{m_i\}_{i=1}^c$ ) and the classification  
 137 function is [2]:

$$138 \quad x_{c^*} = \max_{i \in \{1, \dots, c\}} \langle m_i, v \rangle, \quad (3)$$

139 where  $\langle m_i, v \rangle$  represents the inner product between  $m_i$  and  $v$ , in a hypothesis space  $H$   
 140 having defined an inner product  $\langle \cdot, \cdot \rangle$ .

141 In biometrics, the hyperplane decision will help us to identify and to classify in a much better way  
 142 the biometrics 3D representations of face recognition.

143

#### 144 2.1.2. Cryptography

145 *Cryptosystem*. It has more components: plaintext space  $P$ , ciphertext space  $C$ , key space  $K$ , and  
 146 encryption functions  $E = \{E_k | k \in K\}$ , and decryption functions  $D = \{D_k | k \in K\}$ :

$$147 \quad \forall e \in K \exists d \in K \text{ such that } D_d(E_e(m)) = m, \forall m \in P \quad (4)$$

148 *Encryption scheme*. It represents a particularization of a cryptosystem. There are two types of  
 149 encryption schemes: *symmetric key* schemes (in which the same key is used for encryption and also  
 150 decryption) and *public key* encryption schemes (in which a public key is used to encrypt messages  
 151 and a private key is used to decrypt messages).

152 *Homomorphic encryption*. It is a special type of encryption which allows to apply functions over  
 153 encrypted message, resulting also an encrypted result, which, when decrypted it is the same as  
 154 applying the same function over unencrypted message. This is a powerful cryptographic technique,  
 155 because it increases the security of data, as the operations are applied on encrypted data, resulting  
 156 an encrypted output, which will be decrypted only the users that own the decryption key.  
 157 Unfortunately, at this moment a *fully homomorphic encryption* scheme (that allows to apply *any*  
 158 function on encrypted data) does not exist, because the existing computational capabilities are  
 159 overwhelm. An example of partial homomorphic encryption is RSA cryptosystem [3], where the  
 160 homomorphic operation is multiplication:

$$161 \quad E(m_1) \cdot E(m_2) = m_1^r m_2^r \bmod n = (m_1 m_2)^r \bmod n = E(m_1 \cdot m_2), \quad (5)$$

162 where the public key is modulus  $n$  and the exponent is  $r$ , and encryption function is

$$163 \quad E(m) = m^r \bmod n \quad (6)$$

164 Applying in biometrics, the function has to be changed in order to allow the operations on  
165 bits to be applied for each bit separately.

166

## 167 2.2. Techniques

168 In this section, we present the constructions of the above classification techniques proposed and  
169 improved by the authors of [4], such that they could be applied on encrypted data. The authors of [4]  
170 have shown that their methods are successfully applied on large real datasets, including in face  
171 detection, which became widely used in biometric authentication.

### 172 2.2.1. Auxiliary algorithms

173 In [4] the cryptosystems that have been used are Quadratic Residuosity [5] (where  $P = \mathbb{Z}_2$ ) and  
174 Paillier cryptosystem (where  $P = \mathbb{Z}_N$  and  $N$  is modulus of Paillier) [6]. Further, notation  $(b)_{QR}$   
175 means the bit  $b$  is encrypted with Quadratic Residuosity,  $(m)_{PA}$  means the integer  $m$  is encrypted  
176 with Paillier,  $SK_{QR}$  and  $PK_{QR}$  are secret and public key for Quadratic Residuosity and  $SK_{PA}$  and  
177  $PK_{PA}$  are secret and public key for Paillier.

178 The entities implied in this sections are two parties A and B for building blocks and C (client) and S  
179 (server) for classifiers.

180 Authors of [4] have defined some auxiliary operations: comparison with unencrypted inputs,  
181 comparison with encrypted inputs, reversed comparison over encrypted data, negative integers  
182 comparison and sign determination, which will be used in protocols defined below. The below  
183 algorithm will demonstrate how the encryption over biometric data can work bit by bit.

---

#### Algorithm 1 – max over encrypted data [4]

---

**Input A:**  $k$  integers encrypted using Paillier  $((a_1)_{PA}, \dots, (a_k)_{PA})$ , the length  $l$  of  $a_i$  (in bits),  $PK_{QR}$  and  $SK_{QR}$

**Input B:**  $SK_{PA}, PK_{PA}$ , the length  $l$  in bits

**Output A:**  $\max a_i$

A: generate random permutation  $\pi$  over  $\{1, \dots, k\}$

A:  $(\max)_{PA} := a_{\pi(i)}$

B:  $m := 1$

**for**  $i = 2$  **to**  $k$  **do**

$b_i = \max \leq a_{\pi(i)}$

A: randomly generate integers  $r_i, s_i := (0, 2^{\lambda+l}) \cap \mathbb{Z}$

A:  $(m'_i)_{PA} := (\max)_{PA} \cdot (r_i)_{PA} \quad \triangleright m'_i = \max + r_i$

A:  $(a'_i)_{PA} := (a'_{\pi(i)})_{PA} \cdot (s_i)_{PA} \quad \triangleright a'_i = a_{\pi(i)} + s_i$

A: send  $(m'_i)_{PA}$  and  $(a'_i)_{PA}$  to B

**if**  $b_i$  is true **then**

B:  $m := i$

B:  $(v_i)_{PA} := Refresh(a'_i)_{PA} \quad \triangleright v_i = a'_i$

**else**

---

---

B:  $(v_i)_{PA} := Refresh(m'_i)_{PA} \quad \triangleright v_i = m'_i$

**end if**

B: send  $(v_i)_{PA}$  to A

B: send  $(b_i)_{PA}$

A:  $(max)_{PA} := (v)_{PA} \cdot (g^{-1} \cdot (b_i)_{PA})^{r_i} \cdot ((b_i)_{PA})^{s_i} \quad \triangleright max = v_i + (b_i - 1) \cdot r_i -$   
 $b_i \cdot t_i$

**end for**

B: send  $m$  to A

A: output  $\pi^{-1}(m)$

---

184 Table 1 - max over encrypted data [4]

185 The next algorithm will show how to change an encryption scheme  $E_1$  into an encryption scheme  
 186  $E_2$ , both having the same plaintext size  $M$ . The authors of [4] supposed that  $E_1$  and  $E_2$  are  
 187 additively homomorphic, and semantically secure. In the above algorithm  $(c)_i$  means  $c$  is  
 188 encrypted using encryption scheme  $E_i, i \in \{1, 2\}$ .

---

**Algorithm 2** – Change encryption scheme [4]

---

**Input A:**  $(c)_1, PK_1$  and  $PK_2$

**Input B:**  $SK_1, SK_2$

**Output A:**  $(c)_2$

A: pick  $r \in M$

A: send  $(c')_1 := (c)_1 \cdot (r)_1$  to B

B: decrypt  $(c')_1$  and re-encrypt with  $E_2$

B: send  $(c')_2$  to A

A:  $(c)_2 = (c')_2 \cdot (r)_2^{-1}$

A: output  $(c)_2$

---

189 Table 2 – changing encryption scheme [4]

---

**Algorithm 3** – Private inner product [4]

---

**Input A:**  $a = (a_1, \dots, a_d) \in \mathbb{Z}^d, PK_{PA}$

**Input B:**  $b = (b_1, \dots, b_d) \in \mathbb{Z}^d, SK_{PA}$

**Output A:**  $(\langle a, b \rangle)_{PA}$

B: encrypt  $b$

B: send  $(b_i)_{PA}$  to A

A: compute  $(v)_{PA} = \prod_i (b_i)_{PA}^{x_i} \text{ mod } N^2 \quad \triangleright v = \sum b_i a_i$

A: re-randomize

A: output  $(v)_{PA}$

---

190 Table 3 – Private inner product [4]

191 The authors of [4] proved that these algorithms are secured in *honest but curious* model.

## 192 2.2.2. Naïve Bayes

193 In order to apply Naïve Bayes on encrypted data, there are needed more transformations. In [4], it  
 194 was working with the logarithm of the probability distribution:

$$195 \quad x_c^* = \max_{i \in \{1, \dots, c\}} \log p(X = x_i | V = v)$$

$$196 \quad = \max_{i \in \{1, \dots, c\}} \left\{ \log p(X = x_i) + \sum_{j=1}^n \log p(V_j = v_j | C = c_i) \right\} \quad (7)$$

197 The two types of auxiliary table were used:

198 - One table in which are stored values  $P(i) = \lceil K \log p(X = x_i) \rceil$

199 - One table for every feature  $j$  and class  $i$ :  $T_{i,j}(q) \approx \lceil K \log p(Y_j = q | V = v_i) \rceil, \forall q \in D_j$

200 where  $D_j$  represent the domain of possible values of  $v_j$ , and  $K \in \mathbb{N}$  is a constant.

201 Now to apply Naïve Bayes on encrypted data, the client needs to compute  $(p_i)_{PA} =$

202  $(P_i)_{PA} \prod_{j=1}^n (T_{i,j}(v_j))_{PA}$ , and then uses Algorithm 1 to obtain  $\max p_i$ .

### 203 2.2.3. Hyperplane decision

204 For hyperplane decision the things are quite simple, because the client will calculate

205  $\langle (m_i, v) \rangle_{PA}, i \in \{1, \dots, c\}$ , using Algorithm 3 for inner product, and then Algorithm 1 is used to

206 compute *max* on encrypted inner product.

## 207 3. Data Clustering in Cloud Computing

### 208 3.1. Preliminaries

#### 209 3.1.1. Fixed-width clustering algorithm

210 Clustering is a machine learning technique (included in unsupervised learning) in which a set of  
211 objects is partitioned into groups called clusters. Different from classification, in clustering there are  
212 no predefined clusters, thus the algorithm needs to find relationships or similarities between objects  
213 [8].

214 Fixed-width clustering (FWC) algorithm is based on a distance measure. The steps of FWC are the  
215 following:

216 1. From a given dataset  $D$  with an established cluster width  $w$ , generate a random set of  $m$   
217 clusters:  $C_i, i \in \{1, \dots, m\}$ .

218 2. Compute Euclidean distance between every point  $p_j, j = \{1, \dots, n\}$  and every cluster  $C_i$ , using  
219 the formula:

$$220 \quad d_{ij}(c_i, p_j) = \sqrt{(c_{ix} - p_{jx})^2 + (c_{iy} - p_{jy})^2} \quad (8)$$

221 3. If  $d_{ij}(c_i, p_j) \leq w$ , then  $p_j$  belongs to  $C_i$  cluster; adjust the centroid of  $C_i$  by computing the  
222 mean of the points that  $C_i$  contains at this moment, using the formula ( $n$  is the number of points  
223 in  $C_i$ ):

$$224 \quad \text{centroid}(C_i) = \left( \frac{p_{1x} + \dots + p_{nx}}{n}, \frac{p_{1y} + \dots + p_{ny}}{n} \right) \quad (9)$$

225 4. If  $d_{ij}(c_i, p_j) > w$ , then  $p_j$  is the new centroid of  $C_i$ .

226 5. Reiterate steps 2, 3, 4 until the end of  $D$ .

### 227 3.1.2. MapReduce

228 MapReduce [7] is a programming model for large datasets processing, which works exclusively on  
 229  $(key, value)$  pairs. It consists in three steps: *map*, *shuffle*, *reduce*, and the user needs to define *map*  
 230 and *reduce* functions. The basic idea is that *map* function takes as input a set of  $(key, value)$  pairs  
 231 and outputs an intermediary set of  $(key, value)$ . These outputs are processed by *shuffle* function,  
 232 which groups all intermediary values corresponding to an intermediary key, and sends them to  
 233 *reduce* function. The *reduce* function will try to join these values in order to decrease the number  
 234 of values. Usually just one value will result for the input key per *reduce* invocation.

235 A well-known software framework is Hadoop MapReduce [8] that allows processing in parallel large  
 236 datasets (multi-terabytes of data). We will not give more details, for a comprehensive view please  
 237 read [8].

### 238 3.2. FWC algorithm with MapReduce

239 In [9] the authors propose a distributed version of FWC algorithm, using MapReduce on a large  
 240 number of virtual machines (VMs), as follows:

- 241 - Inputs: dataset  $D$  and the set of clusters  $C_1, \dots, C_m$
- 242 - Partitioning: the  $N$  points of  $D$  are allocated to the  $M$  available VMs (if  $\frac{M}{N}$  is not integer, then  
 243 the remaining points are allocated to the last VM).
- 244 - Map function. The input is dataset  $D$  encrypted and kept into Hadoop Distributed File System  
 245 (HDFS) as  $(key, value)$  pairs, where *key* represents the position of *value* in a data file and  
 246 *value* represents the encryption of numerical of the data point. The data files are global and sent  
 247 to all mappers. The *map* function in proposed model computes the squared Euclidean distance  
 248 (in order to shun the square root):

$$249 \quad E(d_{ij}(c_i, p_j)) = (E(c_{ix}) - E(p_{jx}))^2 + (E(c_{iy}) - E(p_{jy}))^2 \quad (10)$$

250 The *output* of map is a set of  $(key, value)$  pairs, where *key* is the position of *value* into a  
 251 data file and *value* is the distance  $E(d_{ij}) \stackrel{\text{def}}{=} E(d_{ij}(c_i, p_j))$ . Note that  $E(v)$  means value  $v$  is  
 252 encrypted using function  $E$ .

- 253 - *Reduce* function. The output of a *map* function becomes the input for a *reduce* function. The  
 254 *reduce* function needs to find a minimum distance between every point  $E(p_j)$  and every  
 255 centroid  $E(c_i)$  and then to put data point  $E(p_j)$  into corresponding cluster (the one that  
 256 corresponds to the minimum distance) [8].

257 Next, we present the pseudo-code algorithms for *map* and *reduce* functions as the authors provide  
 258 them in [9].

259

260



---

**Algorithm 4** – Map function [9]

---

**Input:** encrypted dataset  $E(D)$ **Output:**  $\langle key, ctxt(value) \rangle \rightarrow \langle index, encrypted\ distance\ E(d_{ij}) \rangle$ Initialization: Choose a random set of clusters  $c_1, \dots, c_m$  from a given dataset  $E(D)$ 

index = 0

for ( $i=0$  to  $D.length$ ) do    for ( $j=0$  to  $c.length$ ) do         $E(d_{ij}) = computeDist(E(p_j), E(c_i))$         index =  $i + j$ 

end

end

End

Take index as key

Construct value as an encrypted numerical value  $E(d_{ij})$ Output  $\langle key, ctxt(value) \rangle$  pair

---

261

Table 4 – *map* function for distributed version for FWC

---

**Algorithm 5** – Reduce function [9]

---

**Input:**  $\langle index, encrypted\ distance\ E(d_{ij}) \rangle$ **Output:**  $\langle ctxt(key), ctxt(value) \rangle \rightarrow \langle E(c_i), E(p_j) \rangle$ Initialization:  $E(minDis)$ for ( $i=0$  to  $D.length$ ) do     $E(minDis) = \min(d_{i1}, \dots, d_{ij})$     if  $E(minDis) \leq w$  then        assign( $E(p_j), E(c_i)$ )        update( $E(c_i)$ )

else

        createNewCluster( $E(p_j)$ )

end

end

Take  $E(p_j)$  as keyConstruct value as a numerical value  $E(c_i)$ Output  $\langle ctxt(key), ctxt(value) \rangle$  pair

---

262

Table 4 – *reduce* function for distributed version for FWC

---

263 We can easily see the potential of this approach of FWC algorithm in biometry. For example, it  
 264 could be used on large datasets of human face datasets in order to find different groups among the  
 265 subjects by analyzing the face images. An advantage of the proposed model in [9] is that  
 266 computations are made on encrypted data, assuring data privacy.

267

268

#### 269 4. Our proposed solution

270 This section will describe the proposed solution, which has been implemented already practical, and  
 271 the results obtained where positive in encrypting data. The idea was developed and presented in  
 272 details in [10] and [11].

273 The below scheme can be adapted with success for all the mentioned algorithms mentioned above.  
 274 All the algorithms were been implemented with success in C, C++, Java, and C#. The source code will  
 275 be available at the link mentioned below.

276 The scheme has been created in a very flexible manner giving the possibility to be adapted  
 277 accordingly to any type of algorithms especially for those ones from machine learning field.

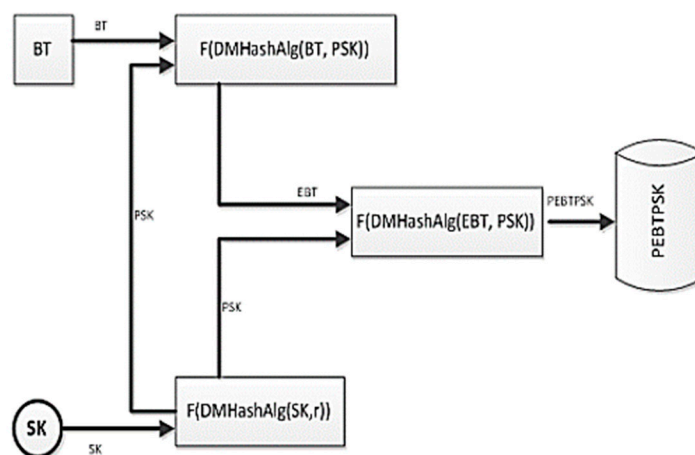
278 The algorithm has two components:

279 1 – The session key algorithm, and

280 2 – The scheme used for enrollment with data integrity checking and validating for the biometric  
 281 data.

282 In Figure 4 – 1, we can see that the permutation functions are applied on all the components that play  
 283 a role in the enrollment scheme. The permutation function will allow the original value to be re-  
 284 arranged in such a way that will be very difficult to understand something from permuted value. The  
 285 original values influence the encrypted system. The content of data could be protected using the  
 286 complex active action, which comes from the characteristic of the system mentioned above. The  
 287 behavior will result into a random series, which could be utilized to data encryption from secret  
 288 communication. Thus, an appropriate key controls encryption or decryption of data content. Machine  
 289 learning is also an appropriate choice for using a hash function due to one-way property. The model  
 290 described in Figure 4 - 2 stores the encrypted Biometric Template using Session Key. A session key is  
 291 the process that generates a random encoding and decoding key which ensures the privacy of a  
 292 session of communications. A similar example can be found in [5] and [6].

293



294

295 Fig. 4 – 1. Data integrity checking and validation for biometric data

296 As some notations for the scheme presented in Figure 4 – 1, the followings have to be considered:

- 297
- BT – biometric data

- 298 • SK – represents the session key, which is generated using a one of the algorithms, which were presented  
299 in Chapter 2 and 3 and combined, with elements of machine learning described in Section 2.1.1. The BT  
300 contains the biometric vector ( $b_v$ ) as we have discussed previous.
- 301 • PSK – represents the permuted session key, which is used to generate the extended version of permuted  
302 transformation of the session, key (SK).  $F(SK)$  represents a function used for permutation which can be  
303 used with any hash function generated based on the main ideas presented above.
- 304 • EBT – represents the biometric data, which are encrypted. In order to generate the biometric template  
305 encrypted, the hash function construction is applied  $F(DMHashAlg(BT, PSK))$ . The hash function is  
306 based on a simple XOR function and both functions  $F(DMHashAlg(SK, r))$  and  $F(DMHashAlg(EBt,$   
307  $PSL))$  functions are used together beside the hash functions with the permutation of the bits SK, EBT  
308 and ESK.
- 309 • PEBTPSK – the permuted biometric data and also the permuted session key (SK) will be used to  
310 generate the final step in order to concatenate the biometric template  $F(DMHashAlg(EBT, PSL))$ . In  
311 order to assure the decryption process, the biometric pattern is using the functions and session key,  
312 which has been used to encrypt the template.

313 In the end, the presented idea represents a new scheme for assuring the data integrity of the biometric  
314 data. The algorithm has been implemented with success and it was tested with positive results on a  
315 set of 700 unique biometric data of 523 subjects.

316 To access the source code of the application, please, visit the following web address:  
317 <https://www.researchgate.net/project/Biometrics-Analysis-Tool>.

## 318 5. Conclusions

319 Applying cryptographic mechanisms and machine learning over biometric data is not an easy task  
320 to accomplish. This fact can be due to the high complexity of how the biometric data are scanned and  
321 read from the user and transfer into the system. Every time that the evolution of technology is making  
322 important advances the complexity of assuring the security and integrity of the data is becoming a  
323 real pain for developers and designers of authentication systems on biometrics.

324 We have proven that applying machine learning techniques and cryptography mechanisms can be a  
325 task that can be accomplished. The most important aspect on which we have focused in this work  
326 paper is how the parameters can be represented and adapted within the algorithms and techniques  
327 used in cryptography and machine-learning. The complexity of the algorithms and the time  
328 processing can represented a problem but not so critical at this time for any of the system  
329 configuration based on the highest requirements possible.

330 The algorithms presented in Chapter 3 and Chapter 4 where implemented with success and we have  
331 obtained positive results.

332 All the results can be viewed at <https://www.researchgate.net/project/Biometrics-Analysis-Tool>. *The*  
333 *software used in analysis and simulation is presented at the mentioned web site. Due to copyrights, we didn't*  
334 *make the source code available at this moment.*

## 335 References

- 336 1. Samuel AL: Some studies in machine learning using the game of checkers. IBM Journal of research and  
337 development. 2000;44.1.2:206-226.
- 338 2. Nasrabadi NM: Pattern recognition and machine learning. Journal of electronic imaging. 2007; 16.4:  
339 049901.

- 340 3. Rivest RL, Shamir A, Adleman L: A method for obtaining digital signatures and public-key  
341 cryptosystems. *Communications of the ACM*. 1976;21(2):120-126.
- 342 4. Bost R, Popa RA, Tu S, Goldwasser S: Machine learning classification over encrypted data. In *NDSS*.  
343 2015;4324:4325.
- 344 5. Goldwasser S, Micali S: Probabilistic encryption & how to play mental poker keeping secret all partial  
345 information. In: *Proceedings of the fourteenth annual ACM symposium on Theory of computing*; May  
346 1982: ACM; 1982. p. 365-377
- 347 6. Paillier P: Public-key cryptosystems based on composite degree residuosity classes. In: *International*  
348 *Conference on the Theory and Applications of Cryptographic Techniques*, May 1999; Berlin,  
349 Heidelberg: Springer; 1999. p. 223-238
- 350 7. Dean J, Ghemawat S: MapReduce: simplified data processing on large clusters. *Communications of*  
351 *the ACM*. 2007; 51(1):107-113.
- 352 8. Taylor RC: An overview of the Hadoop/MapReduce/HBase framework and its current applications in  
353 bioinformatics. In *BMC bioinformatics*. 2010;11(12):S1.
- 354 9. Alabdulatif A, Khalil I, Reynolds M, Kumarage H, Yi X: Privacy-preserving data clustering in cloud  
355 computing based on fully homomorphic encryption. In *PACIS 2017: Societal Transformation Through*  
356 *IS/IT*. Association for Information Systems (AIS). 2017:1-13.
- 357 10. Mihailescu Marius Iulian, Nita Stefania Loredana, *Security of Biometrics Authentication Protocols. Theory*  
358 *and Practice Applications*, LAP LAMBERT Academic Publishing, ISBN-10: 3659777498, ISBN-13: 978-  
359 3659777493. Publishing Date: 11 September 2015, Language: English, 232 pages.
- 360 11. Mihailescu Marius Iulian, New Enrollment Scheme for Biometric Template using Hash Chaos-based  
361 cryptography, Elsevier - *Procedia Engineering*, Volume 69, 2014, pages 1459-1468, ISSN: 1877-7058.
- 362 12. Mihailescu Marius Iulian, Racuciu Ciprian, Grecu Dan Laurentiu, Nita Stefania Loredana, A Multi-  
363 Factor Authentication Scheme Including Biometric Characteristics as One Factor, 1<sup>st</sup> International  
364 Conference Sea-Conf, pp.:348-353, *Mircea cel Batran Naval Academy Scientific Bulletin*, Volume XVIII,  
365 Issue 1, ISSN: 2392-8956, ISSN-L: 1454-864X, CNCS Code: 884, 14-16 May 2015, Constanta, Romania.