



<sup>1,2</sup>. Mohamed SOLTANE

## FACE, VOICE AND SIGNATURE MULTI-MODAL BIOMETRIC VERIFICATION FUSION SYSTEMS

<sup>1</sup>Electrical Engineering & Computing Department, Faculty of Sciences & Technology,  
Doctor Yahia Fares University of Medea, 26000 Medea, ALGERIA

<sup>2</sup> Laboratoire des Systèmes Électroniques Avancées (LSEA)

**Abstract:** In this paper, the use of finite Gaussian mixture modal (GMM) based Expectation Maximization (EM) estimated algorithms for score level data fusion is proposed. Automated biometric systems for human identification measure a "signature" of the human body, compare the resulting characteristic to a database, and render an application dependent decision. These biometric systems for personal authentication and identification are based upon physiological or behavioral features which are typically distinctive, Multi-biometric systems, which consolidate information from multiple biometric sources, are gaining popularity because they are able to overcome limitations such as non-universality, noisy sensor data, large intra-user variations and susceptibility to spoof attacks that are commonly encountered in mono modal biometric systems. Simulation show that finite mixture modal (GMM) is quite effective in modelling the genuine and impostor score densities, fusion based the resulting density estimates achieves a significant performance on eNTERFACE 2005 multi-biometric database based on dynamic face, signature and speech modalities.

**Keywords:** Biometry, Multi-Modal, Authentication, Face Recognition, Speaker and Signature Verification, data Fusion, Adaptive Bayesian decision, GMM& EM

### 1. INTRODUCTION

BIOMETRIC is a Greek composite word stemming from the synthesis of bio and metric, meaning life measurement. In this context, the science of biometrics is concerned with the accurate measurement of unique biological characteristics of an individual in order to securely identify them to a computer or other electronic system. Biological characteristics measured usually include fingerprints, voice patterns, retinal and iris scans, face patterns, and even the chemical composition of an individual's DNA [1]. Biometrics authentication (BA) (*Am I whom I claim I am?*) involves confirming or denying a person's *claimed identity* based on his/her physiological or behavioral characteristics [3]. BA is becoming an important alternative to traditional authentication methods such as keys ("something one has", i.e., by possession) or PIN numbers ("something one knows", i.e., by knowledge) because it is essentially "who one is", i.e., by biometric information. Therefore, it is not susceptible to misplacement or forgetfulness [4]. These biometric systems for personal authentication and identification are based upon physiological or behavioral features which are typically distinctive, although time varying, such as fingerprints, hand geometry, face, voice, lip movement, gait, and iris patterns. Multi-biometric systems, which consolidate information from multiple biometric sources, are gaining popularity because they are able to overcome limitations such as non-universality, noisy sensor data, large intra-user variations and susceptibility to spoof attacks that are commonly encountered in mono-biometric systems. Some works based on multi-modal biometric identity verification systems has been reported in literature. Ben-Yacoub et al. [15] evaluated five binary classifiers on combinations of face and voice modalities (XM2VTS database). They found that (i) a support vector machine and bayesian classifier achieved almost the same performances; and (ii) both outperformed Fisher's linear discriminant, a C4.5 decision tree, and a multilayer perceptron. Korves et al. [16] compared various parametric techniques on the BSSR1 dataset. That study showed that the Best Linear technique performed consistently well, in sharp contrast to many alternative parametric techniques, including simple sum of z-scores, Fisher's linear discriminant analysis, and an implementation of sum of probabilities based on a normal (Gaussian) assumption. Jain et al. [20] propose a framework for optimal combination of match scores that is based on the likelihood ratio test. The distributions of genuine and impostor match scores are modeled as finite Gaussian mixture model. The proposed fusion approach is general in its ability to handle (i) discrete values in biometric match score distributions, (ii) arbitrary scales and distributions of match scores, (iii) correlation between the scores of multiple matchers and (iv) sample quality of multiple biometric sources. The performance of complete

likelihood ratio based fusion rule was evaluated on the three partitions of the NIST-BSSR1 database and the XM2VTS-Benchmark database. As expected, likelihood ratio based fusion leads to significant improvement in the performance compared to the best single modality on all the four databases. At a false accept rate (FAR) of 0.01%. Jain et al. [17] applied the sum of scores, max-score, and min-score fusion methods to normalized scores of face, fingerprint and hand geometry biometrics (database of 100 users, based on a fixed TAR). The normalized scores were obtained by using one of the following techniques: simple distance-to-similarity transformation with no change in scale (STrans), min-max, z-score, median-MAD, double sigmoid, tanh, and Parzen. They found that (a) the min-max, z-score, and tanh normalization schemes followed by a simple sum of scores outperformed other methods; (b) tanh is better than min-max and z-score when densities are unknown; and (c) optimizing the weighting of each biometric on a user-by-user basis outperforms generic weightings of biometrics. Kittler et al. [23] proposed a multi-modal person verification system, using three experts: frontal face, face profile, and voice. The best combination results are obtained for a simple sum rule. Snelick et al. [18] compared combinations of z-score, min-max, tanh and adaptive (two-quadratics, logistic and quadric-line-quadric) normalization methods and simple sum, min score, max score, matcher weighting, and user weighting fusion methods (database of about 1000 users, at a fixed FAR). They found that (a) fusing COTS fingerprint and face biometrics does outperform mono-modal COTS systems, but the high performance of mono-modal COTS systems limits the magnitude of the performance gain; (b) for open-population applications (e.g., airports) with unknown posterior densities, min-max normalization and simple-sum fusion are effective; (c) for closed-population applications (e.g. an office), where repeated user samples and their statistics can be accumulated, QLQ adaptive normalization and user weighting fusion methods are effective. Teoh et al. [19] Applied a modified  $k$ -NN and evidence theoretic  $k$ -NN classifier based on Dempster-safer theory, and it found that the best result is obtained using the evidence theoretic  $k$ -NN classifier as it introduces low FAR and FRR compared to both the ordinary and modified  $k$ -NN. Youssef Elmir et al. [31] present a hierarchical strategy fusion based on multimodal biometric system. The strategy presented relies on a combination of several biometric traits using a multi-level biometric fusion hierarchy. The multi-level biometric fusion includes a pre-classification fusion with optimal feature selection and a post-classification fusion that is based on the similarity of the maximum of matching scores. Mandeep Kaur et al. [32] discusses about Multimodal Biometric System such as signature and speech modalities which are used to overcome some of the problems of uni-modal systems like noise in sensed data, intra-class variations, distinctiveness, and spoof attacks.

A multi-modal biometric verification system based on facial, signature and vocal modalities is described in this paper. Both face images and speech biometrics are chosen due to their complementary characteristics, physiology, and behavior. In multimodal systems, complementary input modalities provide the system with non-redundant information whereas redundant input modalities allow increasing both the accuracy of the fused information by reducing overall uncertainty and the reliability of the system in case of noisy information from a single modality. Information in one modality may be used to disambiguate information in the other ones. The enhancement of precision and reliability is the potential result of integrating modalities and/or measurements sensed by multiple sensors [5].

## 2. AUTHENTICATION TRAITS

### 2.1. Face Extraction and Recognition

Face recognition, authentication and identification are often confused. Face recognition is a general topic that includes both face identification and face authentication (also called verification). On one hand, face authentication is concerned with validating a claimed identity based on the image of a face, and either accepting or rejecting the identity claim (one-to-one matching). On the other hand, the goal of face identification is to identify a person based on the image of a face. This face image has to be compared with all the registered persons (one-to-many matching). Thus, the key issue in face recognition is to extract the meaningful features that characterize a human face. Hence there are two major tasks for that: Face detection and face verification.

**2.1.1. Face detection:** Face detection is concerned with finding whether or not there are any faces in a given image (usually in gray scale) and, if present, return the image location and content of each face. This is the first step of any fully automatic system that analyzes the information contained in faces (e.g., identity, gender, expression, age, race and pose). While earlier work dealt mainly with upright frontal faces, several systems have been developed that are able to detect faces fairly accurately with in-plane or out-of-plane rotations in real time. For biometric systems that use faces as non-intrusive input modules, it is imperative to locate faces in a scene before any recognition algorithm can be applied. An intelligent vision based user interface should be able to tell the attention focus of the user (i.e., where the user is looking at) in order to respond accordingly. To detect facial features accurately for applications such as digital cosmetics, faces need to be located and registered first to facilitate further processing. It is evident that face detection plays an important and critical role for the success of any face processing systems.

On the results presented on this paper only size normalization of the extracted faces was used. All face images were resized to 130x150 pixels, applying a bi-cubic interpolation. After this stage, it is also developed a position correction algorithm based on

detecting the eyes into the face and applying a rotation and resize to align the eyes of all pictures in the same coordinates. The face detection and segmentation tasks presented in this paper was performed based on ‘Face analysis in Polar Frequency Domain’ proposed by Yossi Z. et al. [11]. First it extract the Fourier-Bessel (FB) coefficients from the images. Next, it compute the Cartesian distance between all the Fourier-Bessel transformation (FBT) representations and re-define each object by its distance to all other objects. Images were transformed by a FBT up to the 30<sup>th</sup> Bessel order and 6<sup>th</sup> root with angular resolution of 3°, thus obtaining to 372 coefficients. These coefficients correspond to a frequency range of up to 30 and 3 cycles/image of angular and radial frequency, respectively. Figure 1. Shows the face and eyes detections for different users from the database, and figure 2. Shows the face normalization for the same users.



Figure 1. Face & Eyes Detections for different users.



Figure 2. Face Normalization for the above users.

**Polar Frequency Analysis:** The FB series is useful to describe the radial and angular components in images [11]. FBT analysis starts by converting the coordinates of a region of interest from Cartesian  $(x, y)$  to polar  $(r, \theta)$ . The  $f(r, \theta)$  function is represented by the two-dimensional FB series, defined as:

$$f(r, \theta) = \sum_{i=1}^{\infty} \sum_{n=1}^{\infty} A_{n,i} J_n(\alpha_{n,i} r) \cos(n\theta) + \sum_{i=1}^{\infty} \sum_{n=1}^{\infty} B_{n,i} J_n(\alpha_{n,i} r) \sin(n\theta) \quad (1)$$

Where  $J_n$  is the Bessel function of order  $n$ ,  $f(R, \theta) = 0$  and  $0 \leq r \leq R$ .  $\alpha_{n,i}$  is the  $i^{th}$  root of the  $J_n$  function, i.e. the zero crossing value satisfying  $J_n(\alpha_{n,i}) = 0$  is the radial distance to the edge of the image. The orthogonal coefficients  $A_{n,i}$  and  $B_{n,i}$  are given by:

$$A_{0,i} = \frac{1}{\pi R^2 J_1^2(\alpha_{0,i})} \int_{\theta=0}^{\theta=2\pi} \int_{r=0}^{r=R} f(r, \theta) r J_0\left(\frac{\alpha_{0,i} r}{R}\right) dr d\theta \quad (2)$$

if  $B_{0,i} = 0$  and  $n = 0$ ;

$$\begin{bmatrix} A_{n,i} \\ B_{n,i} \end{bmatrix} = \frac{2}{\pi R^2 J_{n+1}^2(\alpha_{n,i})} \int_{\theta=0}^{\theta=2\pi} \int_{r=0}^{r=R} f(r, \theta) r J_n\left(\frac{\alpha_{n,i} r}{R}\right) \begin{bmatrix} \cos(n\theta) \\ \sin(n\theta) \end{bmatrix} dr d\theta \quad (3)$$

if  $n > 0$ .

An alternative method to polar frequency analysis is to represent images by polar Fourier transform descriptors. The polar Fourier transform is a well known mathematical operation where, after converting the image coordinates from Cartesian to polar, as described above; a conventional Fourier transformation is applied. These descriptors are directly related to radial and angular components, but are not identical to the coefficients extracted by the FBT.

### 2.1.2. Face Verification:

**Feature Extraction:** The so-called “eigenfaces” method [10] is one of the most popular methods for face recognition. It is based on the Principal Components Analysis (PCA) of the face images in a training set. The main idea is that since all human faces share certain common characteristics, pixels in a set of face images will be highly correlated. The K-L (Karhunen-Loeve) transform can be used to project face images to a different vector space that is of reduced dimensionality where features will be uncorrelated. In the new space

nearest neighbor classifiers can be used for classification. Euclidean distances  $d$  in the projection space are mapped into the  $[0,1]$  interval of the real line using the mapping function:  $f = d / (1+d)$ . It is easily seen that  $f$  is also a metric with distance values in  $[0,1]$ . Thus, the decomposition of a face image into an eigenface space provides a set of features. The maximum number of features is restricted to the number of images used to compute the KL transform, although usually only the more relevant features are selected, removing the ones associated with the smallest eigenvalues. Two different approaches, database training stage and the operational stage [10]. The concept verification system is illustrated in figure 4.

**The training stage:** Face spaces are eigenvectors of the covariance matrix corresponding to the original face images, and since they are face-like in appearance, they are so are called Eigenfaces.

Consider the training set of face images be  $i_1, i_2, \dots, i_m$ ; the average face of the set is defined as:

$$\bar{i} = \frac{1}{M} \sum_{j=1}^M i_j \tag{4}$$

where  $M$  is the total number of images.

Each face differs from the average by the vector  $\phi_n = i_n - \bar{i}$ . A covariance matrix is constructed where:

$$C = \sum_{j=1}^M \phi_j \phi_j^T = AA^T \tag{5}$$

where  $A = [\phi_1 \ \phi_2 \ \dots \ \phi_M]$ .

Then, the eigenvectors  $v_k$  and the eigenvalues  $\lambda_k$  with a symmetric matrix  $C$  are calculated.  $v_k$  Determines the linear combination of  $M$  difference images with  $\phi$  to form the Eigen faces:

$$u_l = \sum_{k=1}^M v_{lk} \phi_k \quad l = 1, \dots, M \tag{6}$$

From these Eigen faces,  $K (< M)$  Eigenfaces are selected corresponding to the  $K$  highest eigenvalues.

At the training stage, a set of normalized face images,  $\{i\}$ , that best describe the distribution of the raining facial images in a lower dimensional subspace (Eigen face) is computed by the following operation:

$$\omega_k = u_k (i_n - \bar{i}) \tag{7}$$

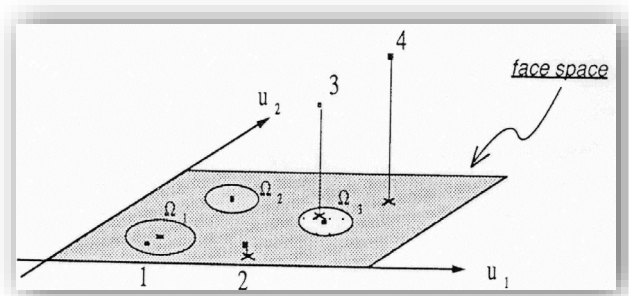
where  $n = 1, \dots, M$  and  $k = 1, \dots, K$ .

After that, the training facial images are projected onto the Eigen space,  $\Omega_i$ , to generate representations of the facial images in Eigen face:

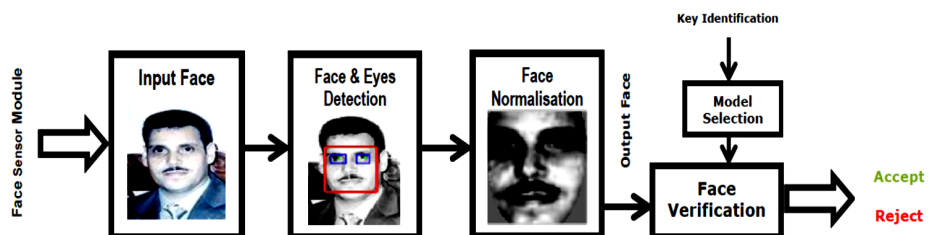
$$\Omega_i = (\omega_{n1}, \omega_{n2}, \dots, \omega_{nk}) \tag{8}$$

where  $n = 1, 2, \dots, M$ .

**The operational stage:** This approach is based on the same principles as standard PCA, explained in the training stage. The difference is that an eigenface space is extracted for each user. Thus, when a claimant wants to verify its identity, its vectorized face image is projected exclusively into the claimed user eigenface space and the corresponding likelihood is computed. The advantage of this approach is that it allows a more accurate model of the user's most relevant information, where the first eigenfaces are directly the most representative user's face information. Another interesting point of this method is its scalability in terms of the number of users. Adding a new user or new pictures of an already registered user only requires to compute or recompute the specific eigenface space, but not the whole dataset base as in the standard approach. For verification systems, the computation of the claimant's likelihood to be a specific user is independent on the number of users in the dataset. On the contrary, for identification systems, the number of operations increases in a proportional way with the number of users, because as many projections as different users are required. In the verification system described in this article, the independent user Eigen face approach has been chosen. Each user's Eigen face space was computed which 16 frames extracted from the database still faces.



**Figure 3.** Simplified version of the face space illustrating the four results of the projection of an image onto the face space. In this case there are two Eigen faces,  $u_1$  and  $u_2$  [10].



**Figure 4.** Face Verification Concept System



### 2.2. Voice Analysis and Feature Extraction

Gaussian Mixture Models (GMMs), is the main tool used in text-independent speaker verification, in which can be trained using the Expectation Maximization (EM) algorithm [8][12]. In this work the speech modality, is authenticated with a multi-lingual text-independent speaker verification system. The speech trait is comprised of two main components as shown in figure 5: speech feature extraction and a Gaussian Mixture Model (GMM) classifier. The speech signal is analyzed on a frame by frame basis, with a typical frame length of 20 ms and a frame advance of 10 ms [14]. For each frame, a dimensional feature vector is extracted, the discrete Fourier spectrum is obtained via a fast Fourier transform from which magnitude squared spectrum is computed and put it through a bank of filters. The critical band warping is done following an approximation to the Mel-frequency scale which is linear up to 1000 Hz and logarithmic above 1000 Hz. The Mel-scale cepstral coefficients are computed from the outputs of the filter bank [7]. The state

of the art speech feature extraction schemes (Mel frequency cepstral coefficients (MFCC) is based on auditory processing on the spectrum of speech signal and cepstral representation of the resulting features [2]. One of the powerful properties of cepstrum is the fact that any periodicities, or repeated patterns, in a spectrum will be mapped to one or two specific components in the cepstrum. If a spectrum contains several harmonic series, they will be separated in a way similar to the way the spectrum separates repetitive time patterns in the waveform. The description of the different steps to exhibit features characteristics of an audio sample with MFCC is showed in figure 6.

The distribution of feature vectors for each person is modeled by a GMM. The parameters of the Gaussian mixture probability density function are estimated with Expectation Maximization (EM) algorithm [8]. Given a claim for person  $C$ 's identity and a set of feature vectors  $X = \{\vec{x}_i\}_{i=1}^{N_V}$  supporting the claim, the average log likelihood of the claimant being the true claimant is calculated using:

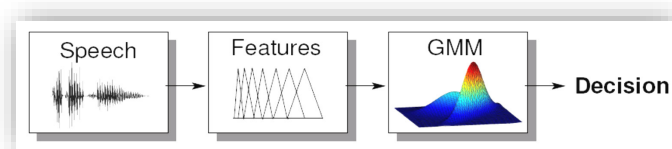


Figure 5. Acoustic Speech Analysis

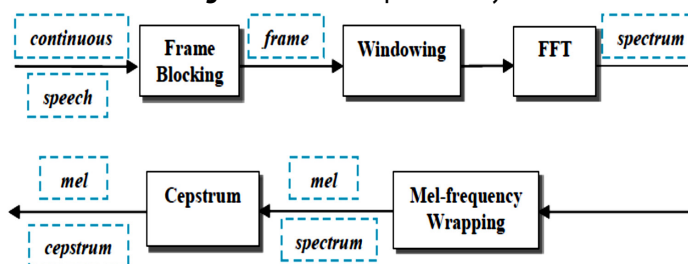


Figure 6. MFCC calculation Block diagram [7]

where

$$\mathcal{L}(X|\lambda_C) = \frac{1}{N_V} \sum_{i=1}^{N_V} \log p(\vec{x}_i|\lambda_C) \tag{9}$$

$$p(\vec{x}|\lambda) = \sum_{j=1}^{N_M} m_j \mathcal{N}(\vec{x}; \vec{\mu}_j; \Sigma_j) \tag{10}$$

and

$$\lambda = \{m_j, \vec{\mu}_j, \Sigma_j\}_{j=1}^{N_M} \tag{11}$$

Here  $\lambda_C$  is the model for person  $C$ .  $N_M$  is the number of mixtures,  $m_j$  is the weight for mixture  $j$  (with constraint  $\sum_{j=1}^{N_M} m_j = 1$ ), and  $\mathcal{N}(\vec{x}; \vec{\mu}, \Sigma)$  is a multi-variate Gaussian function with mean  $\vec{\mu}$  and diagonal covariance matrix  $\Sigma$ . Given a set  $\{\lambda_b\}_{b=1}^B$  of  $B$  background person models for person  $C$ , the average log likelihood of the claimant being an impostor is found using:

$$\mathcal{L}(X|\lambda_{\bar{C}}) = \log \left[ \frac{1}{B} \sum_{b=1}^B \exp \mathcal{L}(X|\lambda_b) \right] \tag{12}$$

The set of background person models is found using the method described in [9]. An opinion on the claim is found using:

$$o = \mathcal{L}(X|\lambda_C) - \mathcal{L}(X|\lambda_{\bar{C}}) \tag{13}$$

The opinion reflects the likelihood that a given claimant is the true claimant (i.e., a low opinion suggests that the claimant is an impostor, while a high opinion suggests that the claimant is the true claimant).

### 2.3. Signature Verification Systems

Handwritten signature is one of the first accepted civilian and forensic biometric identification technique in our society [28][29][30]. Human verification is normally very accurate in identifying genuine signatures. A signature verification system must be able to detect forgeries and at the same time reduce rejection of genuine signatures. The signature verification problem can be classified into categories: offline and online. Offline signature verification does not use dynamic information that is used extensively in online signature verification systems. This paper investigates the problem of offline signature verification. The problem of offline signature verification has been faced by taking into account three different types of forgeries: random forgeries, produced without knowing either the name of the signer or the shape of his signature; simple forgeries, produced knowing the name of the signer but without having an example of his signature; and skilled forgeries, produced by people who, looking at an original instance of the signature, attempt to imitate it as closely as possible.

**Feature Extraction:** The coordinate trajectories  $(x_n, y_n)$  and pressure signal  $p_n$  are the components of the unprocessed feature vectors  $u_n = [x_n, y_n, p_n]^T$  extracted from the signature signal [28][29][30], where  $n = 1, \dots, N_s$ , and  $N_s$  is the duration of the signature in time samples. Signature trajectories are then pre-processed by subtracting the centre of mass followed by rotation alignment based on the average path tangent angle. An extended set of discrete-time functions are derived from the pre-processed trajectories consisting of sample estimations of various dynamic properties. As a result, the parameterised signature  $O$  consists in the sequence of feature vectors  $o_n = [x_n, y_n, p_n, \theta_n, v_n, \dot{x}_n, \dot{y}_n]^T$ ,  $n = 1, \dots, N_s$ , where the upper dot notation represents an approximation to the first order time derivative and  $\theta$  and  $v$  stand respectively for path tangent angle, path velocity magnitude.



Figure 7. Wacom Graphire3 digitizing Tablet PC

$$v_i = \sqrt{\dot{x}_i^2 + \dot{y}_i^2} \text{ and } \theta_i = \arctan(\dot{y}_i, \dot{x}_i) \quad (14)$$

and

$$\dot{x}_i = x_i - x_{i-1} \text{ and } \dot{y}_i = y_i - y_{i-1} \quad (15)$$

A whitening linear transformation is finally applied to each discrete-time function so as to obtain zero mean and unit standard deviation function values. Seven dimensional feature vectors are used for GMM processing described in the following section. Figure 9 shows  $x$ -,  $y$ -,  $p$ - and velocity signals of an example signature.

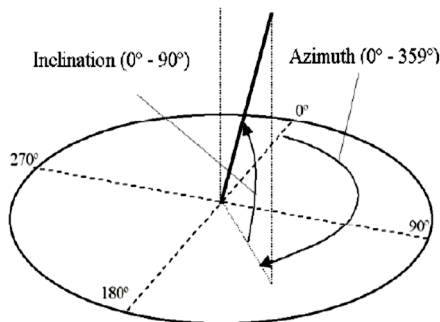


Figure 8. Azimuth and inclination angles of the pen respect to the plane of the graphic card GD-0405U from Wacom Graphire3 digitizing Tablet PC.

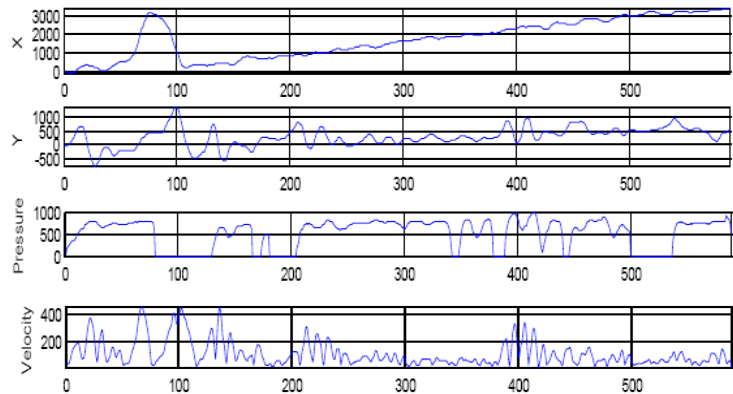


Figure 9. Signals ( $x$ -,  $y$ - position, pen pressure and velocity) of one signature fragment.

### 3. MULTIMODAL BIOMETRIC FUSION DECISION

The process of biometric user authentication can be outlined by the following steps [25]: a) acquisition of raw data, b) extraction of features from these raw data, c) computing a score for the similarity or dissimilarity between these features and a previously given set of reference features and d) classification with respect to the score, using a threshold. The results of the decision processing steps are *true* or *false* (or *accept/reject*) for verification purposes or the user identity for identification scenarios.

The fusion of different signals can be performed 1) at the raw data or the feature level, 2) at the score level or 3) at the decision level. These different approaches have advantages and disadvantages. For *raw data* or *feature level* fusion, the basis data have to be compatible for all modalities and a common matching algorithm (processing step c) must be used. If these conditions are met, the separate feature vectors of the modalities easily could be concatenated into a single new vector. This level of fusion has the advantage that only one algorithm for further processing steps is necessary instead of one for each modality. Another advantage of fusing at this early stage of processing is that no information is lost by previous processing steps. The main disadvantage is the demand of compatibility of the different raw data of features. The fusion at *score level* is performed by computing a similarity or dissimilarity (distance) score for each single modality. For joining of these different scores, normalization should be done. The straightforward and most rigid approach for fusion is the decision level. Here, each biometric modality results in its own decision; in case of a verification scenario this is a set of *true*s and *false*s. From this set a kind of voting (majority decision) or a logical *AND* or *OR* decision can be computed. This level of fusion is the least powerful, due to the absence of much information. On the other hand, the advantage of this fusion strategy is the easiness and the guaranteed availability of all single modality decision results. In practice, score level fusion is the best-researched approach, which appears to result in better improvements of recognition accuracy as compared to the other strategies.

### 3.1 Adaptive Bayesian Method Based Score Fusion

Let  $X = [X_1, X_2, \dots, X_K]$  denote the match scores of  $K$  different biometric matchers, where  $X_k$  is the random variable representing the match score of the  $k^{th}$  matcher,  $k = 1, 2, \dots, K$ . Let  $f_{gen}(x)$  and  $f_{imp}(x)$  be the conditional joint densities of the  $K$  match scores given the genuine and impostor classes, respectively, where  $x = [x_1, x_2, \dots, x_k]$ . Suppose we need to assign the observed match score vector  $X$  to genuine or impostor class. Let  $\Psi$  be a statistical test for testing  $H_0: X$  corresponds to an impostor against  $H_1: X$  corresponds to a genuine user. Let  $\Psi(x) = i$  imply that we decide in favor of  $H_i, i = 0, 1$ . The probability of rejecting  $H_0$  when  $H_0$  is true is known as the *false accept rate* (size or level of the test). The probability of correctly rejecting  $H_0$  when  $H_1$  is true is known as the *genuine accept rate*. The Neyman-Pearson theorem [21][22] states that:

- 1) For testing  $H_0$  against  $H_1$ , there exists a test  $\Psi$  and a constant  $\eta$  such that:

$$P(\Psi(X) = 1 | H_0) = \alpha \tag{16}$$

$$\Psi(x) = \begin{cases} 1, & \text{when } \frac{f_{gen}(x)}{f_{imp}(x)} \geq \eta \\ 0, & \text{when } \frac{f_{gen}(x)}{f_{imp}(x)} < \eta \end{cases} \tag{17}$$

- 2) If a test satisfies equations (16) and (17) for some  $\eta$ , then it is the *most powerful test* for testing  $H_0$  against  $H_1$  at level  $\alpha$ .

According to the Neyman-Pearson theorem, given the false accept rate (FAR)  $\alpha$ , the *optimal* test for deciding whether a score vector  $X$  corresponds to a genuine user or an impostor is the likelihood ratio test given by equation (17). For a fixed FAR, it can select a threshold  $\eta$  such that the likelihood ratio test maximizes the genuine accept rate (GAR). Based on the Neyman-Pearson theorem, we are guaranteed that *there does not exist any other decision rule with a higher GAR*. However, this optimality of the likelihood ratio test is guaranteed only when the underlying densities are known. In practice, it estimate the densities  $f_{gen}(x)$  and  $f_{imp}(x)$  from the training set of genuine and impostor match scores, respectively and the performance of likelihood ratio test will depend on the accuracy of these estimates [13][25].

**3.1.1 Estimation of Match Score Densities:** Gaussian mixture model (GMM) has been successfully used to estimate arbitrary densities and it is used for estimating the genuine and impostor score densities [8][24].

Let  $\Phi^K(x; \mu, \Sigma)$  be the  $K$ -variate Gaussian density with mean vector  $\mu$  and covariance matrix  $\Sigma$ , i.e.,

$\Phi^K(x; \mu, \Sigma) = (2\pi)^{-K/2} |\Sigma|^{-1/2} \exp(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu))$ . The estimates of  $f_{gen}(x)$  and  $f_{imp}(x)$  are obtained as a mixture of Gaussians as follows.

$$\hat{f}_{gen}(x) = \sum_{j=1}^{M_{gen}} P_{gen,j} \Phi^K(x; \mu_{gen,j}, \Sigma_{gen,j}) \tag{18}$$

$$\hat{f}_{imp}(x) = \sum_{j=1}^{M_{imp}} P_{imp,j} \Phi^K(x; \mu_{imp,j}, \Sigma_{imp,j}) \tag{19}$$

where  $M_{gen}$  ( $M_{imp}$ ) is the number of mixture components used to model the density of the genuine (impostor) scores,  $p_{gen,j}$  ( $p_{imp,j}$ ) is the weight assigned to the  $j^{th}$  mixture component in  $\hat{f}_{imp}(x)$  ( $\hat{f}_{gen}(x)$ ),  $\sum_{j=1}^{M_{gen}} P_{gen,j} = \sum_{j=1}^{M_{imp}} P_{imp,j} = 1$ . Selecting the appropriate number of components is one of the most challenging issues in mixture density estimation; while a mixture with too many components may result in over-fitting, a mixture with too few components may not approximate the true density well. The GMM fitting algorithm automatically estimates the number of components and the component parameters using an EM algorithms and the minimum message length criterion [8][24].

**Maximum Likelihood Parameter Estimation:** Given a set of observation data in a matrix  $X$  and a set of observation parameters  $\theta$  the ML parameter estimation aims at maximizing the likelihood  $L(\theta)$  or log likelihood of the observation data  $X = \{X_1, \dots, X_n\}$

$$\hat{\theta} = \arg \max_{\theta} L(\theta). \tag{20}$$

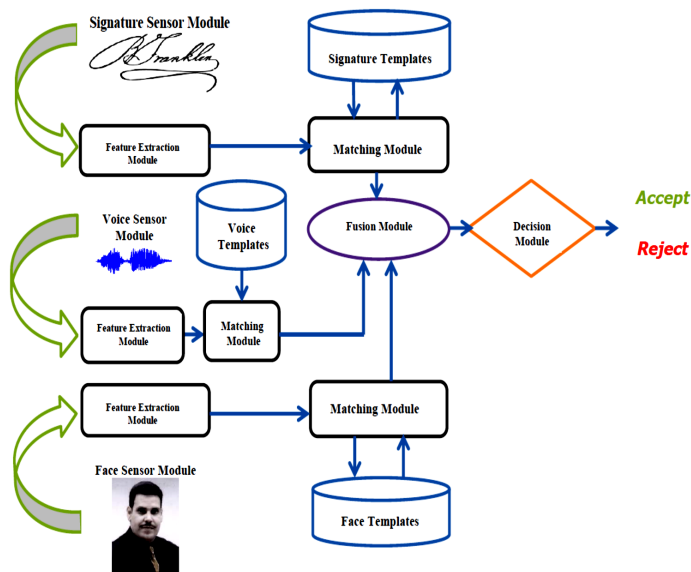


Figure 10. Score Fusion Based Multimodal Biometric Verification System

Assuming that it has independent, identically distributed data, it can write the above equations as:

$$L(\theta) = p(X|\theta) = p(X_1, \dots, X_n|\theta) = \prod_{i=1}^n p(X_i|\theta). \quad (21)$$

The maximum for this function can be find by taking the derivative and set it equal to zero, assuming an analytical function.

$$\frac{\partial}{\partial \theta} L(\theta) = 0. \quad (22)$$

The incomplete-data log-likelihood of the data for the mixture model is given by:

$$L(\theta) = \log(X|\theta) = \sum_{i=1}^N \log(x_i|\theta) \quad (23)$$

which is difficult to optimize because it contains the log of the sum. If it considers  $X$  as incomplete, however, and posits the existence of unobserved data items  $Y = \{y_i\}_{i=1}^N$  whose values inform us which component density generated each data item, the likelihood expression is significantly simplified. That is, it assume that  $y_i \in \{1 \dots K\}$  for each  $i$ , and  $y_i = k$  if the  $i$ -th sample was generated by the  $k$ -th mixture component. If it knows the values of  $Y$ , it obtains the complete-data log-likelihood, given by:

$$L(\theta, Y) = \log p(X, Y|\theta) \quad (24)$$

$$= \sum_{i=1}^N \log p(x_i, y_i|\theta) \quad (25)$$

$$= \sum_{i=1}^N \log(p(y_i|\theta)p(x_i|y_i, \theta)) \quad (26)$$

$$= \sum_{i=1}^N (\log p_{y_i} + \log g(x_i|\mu_{y_i}, \Sigma_{y_i})) \quad (27)$$

which, given a particular form of the component densities, can be optimized using a variety of techniques [23].

**EM Algorithm:** The expectation-maximization (EM) algorithm [24][25][26][27] is a procedure for maximum-likelihood (ML) estimation in the cases where a closed form expression for the optimal parameters is hard to obtain. This iterative algorithm guarantees the monotonic increase in the likelihood  $L$  when the algorithm is run on the same training database.

The probability density of the Gaussian mixture of  $k$  components in  $\mathbb{R}^d$  can be described as follows:

$$\Phi(x) = \sum_{i=1}^k \pi_i \phi(x|\theta_i) \quad \forall x \in \mathbb{R}^d, \quad (28)$$

where  $\phi(x|\theta_i)$  is a Gaussian probability density with the parameters  $\theta_i = (m_i, \Sigma_i)$ ,  $m_i$  is the mean vector and  $\Sigma_i$  is the covariance matrix which is assumed positive definite given by:

$$\phi(x|\theta_i) = \phi(x|m_i, \Sigma_i) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_i|^{\frac{1}{2}}} e^{-\frac{1}{2}(x-m_i)^T \Sigma_i^{-1}(x-m_i)}, \quad (29)$$

and  $\pi_i \in [0, 1] (i = 1, 2, \dots, k)$  are the mixing proportions under the constraint  $\sum_{i=1}^k \pi_i = 1$ . If it encapsulate all the parameters into one vector:  $\theta_k = (\pi_1, \pi_2, \dots, \pi_k, \theta_1, \theta_2, \dots, \theta_k)$ , then, according to Eq. (27), the density of Gaussian mixture can be rewritten as:

$$\Phi(x|\theta_k) = \sum_{i=1}^k \pi_i \phi(x|\theta_i) = \sum_{i=1}^k \pi_i \phi(x|m_i, \Sigma_i). \quad (30)$$

For the Gaussian mixture modeling, there are many learning algorithms. But the EM algorithm may be the most well-known one. By alternatively implementing the E-step to estimate the probability distribution of the unobservable random variable and the M-step to increase the log-likelihood function, the EM algorithm can finally lead to a local maximum of the log-likelihood function of the model. For the Gaussian mixture model, given a sample data set  $S = \{x_1, x_2, \dots, x_N\}$  as a special incomplete data set, the log-likelihood function can be expressed as follows:

$$\log p(S|\theta_k) = \log \prod_{t=1}^N \phi(x_t|\theta_k) = \sum_{t=1}^N \log \sum_{i=1}^k \pi_i \phi(x_t|\theta_i), \quad (31)$$

Which can be optimized iteratively via the EM algorithm as follows:

$$P(j|x_t) = \frac{\pi_j \phi(x_t|\theta_j)}{\sum_{i=1}^k \pi_i \phi(x_t|\theta_i)}, \quad (32)$$

$$\pi_j^+ = \frac{1}{N} \sum_{t=1}^N P(j|x_t), \quad (33)$$

$$\mu_j^+ = \frac{1}{\sum_{t=1}^N P(j|x_t)} \sum_{t=1}^N P(j|x_t) x_t, \quad (34)$$

$$\Sigma_j^+ = \frac{1}{\sum_{t=1}^N P(j|x_t)} \sum_{t=1}^N P(j|x_t) (x_t - \mu_j^+)(x_t - \mu_j^+)^T. \quad (35)$$

Although the EM algorithm can have some good convergence properties in certain situations, it certainly has no ability to determine the proper number of the components for a sample data set because it is based on the maximization of the likelihood.

#### 4. EXPERIMENTS AND RESULTS

The experiments were performed using still faces, signatures and audio database extracted from video, which is encoded in raw UYVY. AVI 640 x 480, 15.00 fps with uncompressed 16bit PCM audio; mono, 32000 Hz little endian. Uncompressed PNG files are extracted from the video files for feeding the face detection algorithms. The capturing devices for recording the video and audio data were: Allied Vision Technologies AVT marlin MF-046C 10 bit ADC, 1/2" (8mm) Progressive scan SONY IT CCD; and Shure SM58 microphone. Frequency response 50 Hz to 15000 Hz. Unidirectional (Cardiod) dynamic vocal microphones. Thirty subjects were used



for the experiments in which twenty-six are males and four are females. For each subject, 30 signatures (with dat header) are used. Each line of a (.dat files) consists of four comma separated integer values for the sampled x- and y-position of the pen tip, the pen pressure and the timestamp (in ms); the lines with values of -1 for x, y and pressure represent a pen-up/pen-down event; The device used for recording the handwriting data was a Wacom Graphire3 digitizing tablet. Size of sensing surface is 127.6mm x 92.8mm. With spatial resolution of 2032 lpi (lines per inch), able to measure 512 degrees of pressure. The signature data is acquired with a non-fixed sampling rate of about 100Hz. The audio is extracted as 16 bit PCM WAV file (with wav header), sampled at 16000 Hz, mono little endian. For the audio six multi-lingual (.wav files) of one minute each recording were used for each subject. The database obtained from eINTERFACE 2005 [6]. Thirty subjects were used for the experiments in which twenty-five are males and five are females. For face experts, ninety-six face images from a subject were randomly selected to be trained and projected into Eigen space, and the other twenty-four samples were used for the subsequent validation and testing. Similarly, four samples were used in speech experts for the modeling (training); two samples were used for the subsequent validation and testing. For signature experts, twenty four signatures from a subject were randomly selected for training, and the other six samples were used for the subsequent validation and testing. Three sessions of the face database, signature and speech database were used separately. Session one was used for training the speech and face experts. Each expert used ten mixture client models. To find the performance, Sessions two and three were used for obtaining expert opinions of known impostor and true claims.

**Performance Criteria:** The basic error measure of a verification system is false rejection rate (FRR) and false acceptance rate (FAR) as defined in the following equations:

**False Rejection Rate (FRR<sub>i</sub>):** is an average of number of falsely rejected transactions. If  $n$  is a transaction and  $x(n)$  is the verification result where 1 is falsely rejected and 0 is accepted and  $N$  is the total number of transactions then the personal False Rejection Rate for user  $i$  is

$$FRR_i = \frac{1}{N} \sum_{n=1}^N x(n) \quad (36)$$

**False Acceptance rate (FAR<sub>i</sub>)** is an average of number of falsely accepted transactions. If  $n$  is a transaction and  $x(n)$  is the verification result where 1 is a falsely accepted transaction and 0 is genuinely accepted transaction and  $N$  is the total number of transactions then the personal False Acceptance Rate for user  $i$  is

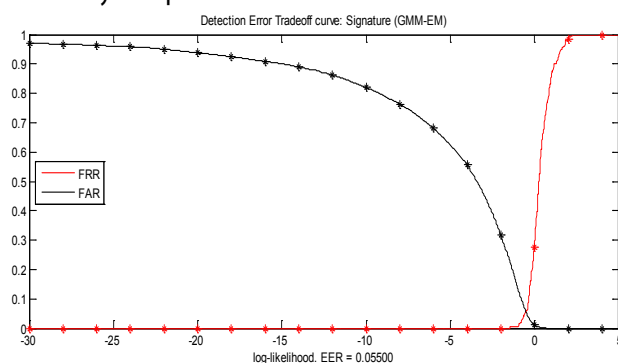
$$FAR_i = \frac{1}{N} \sum_{n=1}^N x(n) \quad (37)$$

Both FRR<sub>i</sub> and FAR<sub>i</sub> are usually calculated as averages over an entire population in a test. If  $P$  is the size of populations then these averages are

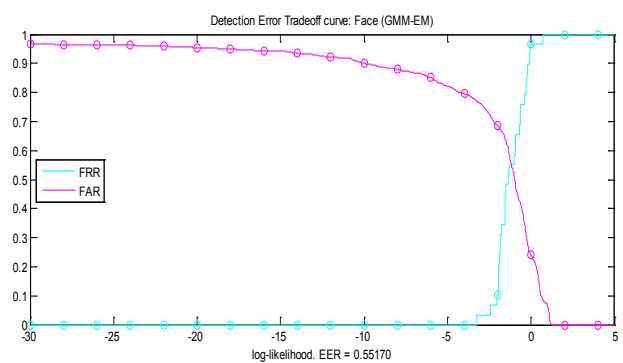
$$FRR = \frac{1}{P} \sum_i^P FRR_i \quad (38)$$

$$FAR = \frac{1}{P} \sum_i^P FAR_i \quad (39)$$

**Equal Error Rate (EER)**, is an intersection where FAR and FRR are equal at an optimal threshold value. This threshold value shows where the system performs at its best.



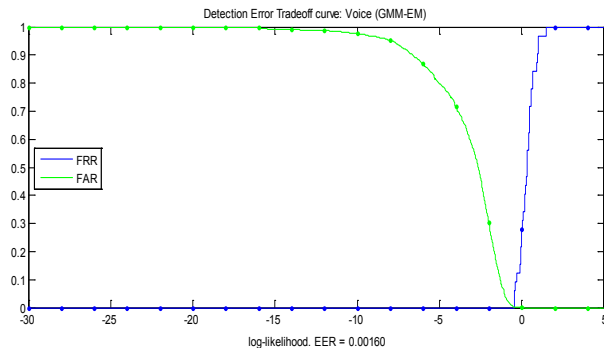
**Figure 11.** Detection error tradeoff curves: Signature (GMM-EM)



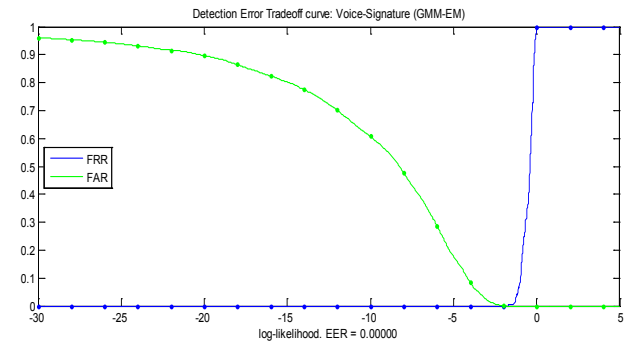
**Figure 12.** Detection error tradeoff curves: Face (GMM-EM)

As a common starting point, classifier parameters were selected to obtain performance as close as possible to EER on clean test data (following the standard practice in the face and speaker verification area of using EER as a measure of expected performance). A good

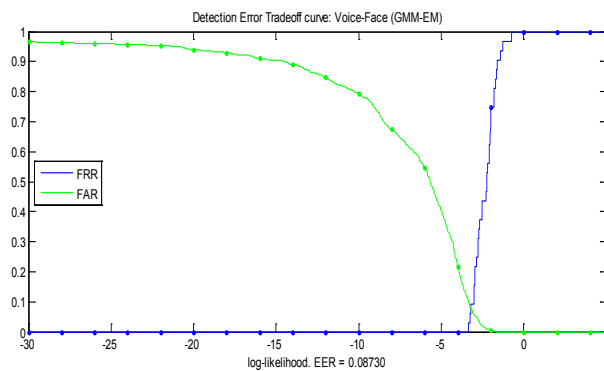
decision is to choose the decision threshold such as the false accept equal to the false reject rate. In this paper it uses the Detection Error Tradeoff (DET) curve to visualize and compare the performance of the system (see Figures 11-17).



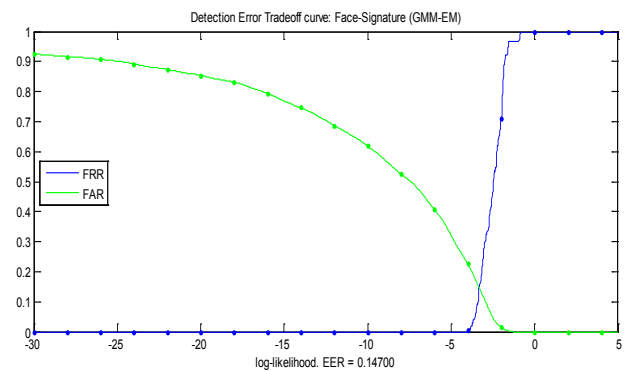
**Figure 13.** Detection error tradeoff curves: Voice (GMM-EM)



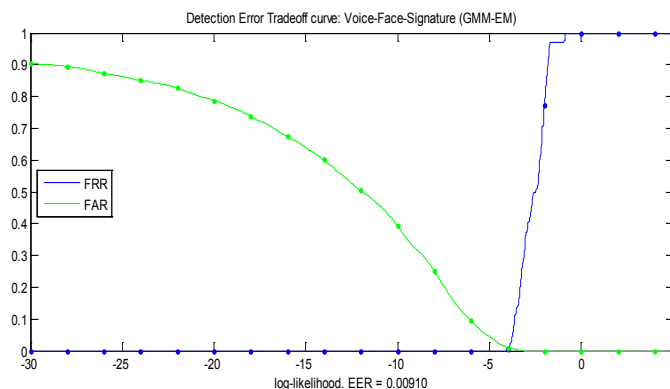
**Figure 14.** Detection error tradeoff curves: Voice-Signature (GMM-EM)



**Figure 15.** Detection error tradeoff curves: Voice-Face (GMM-EM)



**Figure 16.** Detection error tradeoff curves: Face-Signature (GMM-EM)



**Figure 17.** Detection error tradeoff curves: Voice-Face-Signature (GMM-EM)

## 5. CONCLUSIONS

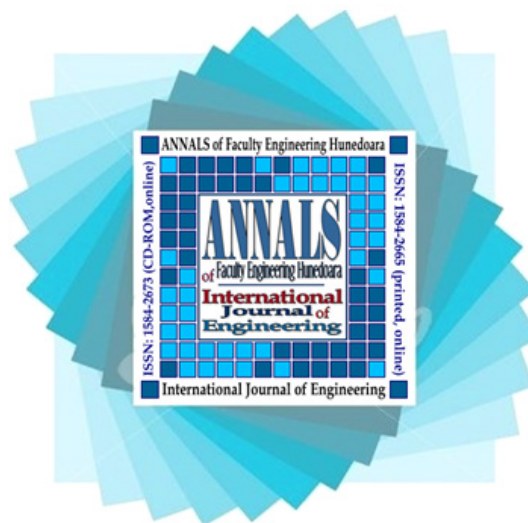
The paper has presented a human authentication method combined face, signature and speech information in order to improve the problem of single biometric authentication, since single biometric authentication has the fundamental problems of high FAR and FRR. It has presented a framework for fusion of match scores in multi-modal biometric system based on adaptive Bayesian method. The likelihood ratio based fusion rule with GMM-based density estimation achieves a significant recognition rates. As a result presented a combined authentication method can provide a stable authentication rate and it overcomes the limitation of a single mode system. Based on the experimental results, it has shown that EER can be reduced down significantly between the face, signature mode and a combined face-voice-signature mode.

## REFERENCES

- [1.] Sofia Gleni & Panagiotis Petratos, "DNA Smart Card for Financial Transactions" The ACM Student Magazine 2004, <http://www.acm.org>
- [2.] S. Davis and P. Mermelstein. Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences. *IEEE Transactions on Acoustics, Speech and Signal Processing*, (4):357–366, 1980.
- [3.] Girija Chetty and Michael Wagner, "Audio-Visual Multimodal Fusion for Biometric Person Authentication and Liveness Verification", 2006, Australian Computer Society, Inc. This paper appeared at the NICTA-HCSNet Multimodal User Interaction Workshop (MMUI2005),

- Sydney, Australia.
- [4.] Norman Poh and Samy Bengio, "Database, Protocol and Tools for Evaluating Score-Level Fusion Algorithms in Biometric Authentication", IDIAP RR 04-44, August 2004, a IDIAP, CP 592, 1920 Martigny, Switzerland.
  - [5.] Corradini, M. Mehta, N.O. Bernsen, J. C. Martin and S. Abrilian, "MULTIMODAL INPUT FUSION IN HUMAN-COMPUTER INTERACTION", On the Example of the NICE Project 2003.
  - [6.] Yannis Stylianou, Yannis Pantazis, Felipe Calderero, Pedro Larroy, Francois Severin, Sascha Schimke, Rolando Bonal, Federico Matta, and Athanasios Valsamakis, "GMM-Based Multimodal Biometric Verification", eNTerFACE 2005 The summer Workshop on Multimodal Interfaces July 18<sup>th</sup> – August 12<sup>th</sup>, Faculté Polytechnique de Mons, Belgium.
  - [7.] Lasse L Mølgaard and Kasper W Jørgensen, "Speaker Recognition: Special Course", IMM\_DTU December 14, 2005
  - [8.] Pekka Paalanen, "Bayesian classification using gaussian mixture model and EM estimation: implementation and comparisons", Information Technology Project, 2004, Lappeenranta, June 23, 2004, <http://www.it.lut.fi/project/gmmbayes/>
  - [9.] D.A. Reynolds, "Experimental Evaluation of Features for Robust Speaker Identification", IEEE Trans. Speech and Audio Processing 2 (4), 1994, 639-643.
  - [10.] M. Turk and A. Pentland, "Eigenfaces for Recognition", Journal of Cognitive Neuroscience, vol.3, no. 1, pp. 71-86, 1991.
  - [11.] Yossi Zana, Roberto M. Cesar-Jr, Rogerio S. Feris, and Matthew Turk, "Face Verification in Polar Frequency Domain: A Biologically Motivated Approach", G. Bebis et al. (Eds.): ISVC 2005, LNCS 3804, pp. 183–190, 2005. C\_Springer-Verlag Berlin Heidelberg 2005
  - [12.] Conrad Sanderson, Samy Bengio, Herve Bourlard, Johnny Mariéthoz, Ronan Collobert, Mohamed F. BenZeghiba, Fabien Cardinaux, and Sébastien Marcel, "Speech & face based biometric authentication at IDIAP", Dalle Molle Institute for Perceptual Artificial Intelligence (IDIAP). Rue du Simplon 4, CH-1920 Martigny, Switzerland.
  - [13.] Karthik Nandakumar, Yi Chen, Sarat C. Dass and Anil K. Jain, "Likelihood Ratio Based Biometric Score Fusion", IEEE Transactions on Pattern Analysis and Machine Intelligence, 2007
  - [14.] Claus Vielhauer, Sascha Schimke, Valsamakis Thanassis and Yannis Stylianou, "Fusion Strategies for Speech and Handwriting Modalities in HCI", Multimedia on Mobile Devices, edited by Reiner Creutzburg, Jarmo H. Takala, Proc. of SPIE-IS&T Electronic Imaging, Vol. 5684 © 2005
  - [15.] Souheil Ben-Yacoub, Yousri Abdeljaoued, and Eddy Mayoraz, "Fusion of Face and Speech Data for Person Identity Verification", IEEE Transactions on Neural Networks, vol. 10, no. 5, september 1999
  - [16.] H. J. Korves, L. D. Nadel, B. T. Ulery, D. M. Bevilacqua Masi, "Multi-biometric Fusion: From Research to Operations", MTS MitreTek Systems, sigma summer 2005, pp. 39-48, <http://www.mitretek.org/home.nsf/Publications/SigmaSummer2005>
  - [17.] Anil Jain<sup>a</sup>, Karthik Nandakumar<sup>a</sup>, Arun Ross<sup>b\*</sup>, "Score normalization in multimodal biometric systems\*", <sup>a</sup>Department of Computer Science and Engineering, Michigan State University, East Lansing, MI 48824, USA, <sup>b</sup>Department of Computer Science and Engineering, West Virginia University, Morgantown, WV 26506, USA, THE JOURNAL OF PATTERN RECOGNITION SOCIETY, ELSEVIER 2005
  - [18.] Robert Snelick, Umut Uludag, Alan Mink, Michael Indova and Anil Jain, "Large Scale Evaluation of Multimodal Biometric Authentication Using State-of-the-Art Systems", IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 27, No. 3, Mar 2005, pp 450-455
  - [19.] A. Teoh, S. A. Samad, and A. Hussain, "Nearest Neighbourhood Classifiers in Biometric Fusion", International Journal of The Computer, the internet and management Vol. 12#1 2004pp 23-36
  - [20.] Karthik Nandakumar, Yi Chen, Sarat C. Dass and Anil K. Jain, "Biometric Score Fusion: Likelihood Ratio, Matcher Correlation and Image Quality" March 2007 DRAFT.
  - [21.] Van Trees, Harry L., "Detection, Estimation, and Modulation Theory", Part I, John Wiley and Sons, 1968.
  - [22.] Qing Yan and Rick S. Blum, "Distributed Signal Detection under the Neyman-Pearson Criterion", EECS Department Lehigh University Bethlehem, PA 18015
  - [23.] Kittler, J., Hatef, M., Duin, R. P. W. and Matas, J., "On combining classifiers". IEEE Transactions on Pattern Analysis and Machine Intelligence, 20(3): 226–239. 1998
  - [24.] P. Paalanen, J.-K. Kamarainen, J. Ilonen, H. Kälviäinen, "Feature Representation and Discrimination Based on Gaussian Mixture Model Probability Densities: Practices and Algorithms", Department of Information Technology, Lappeenranta University of Technology, P.O.Box 20, FI-53851 Lappeenranta, Finland 2005
  - [25.] Kalyan Veeramachaneni, Lisa Ann Osadciw, and Pramod K. Varshney, "An Adaptive Multimodal Biometric Management Algorithm", IEEE TRANSACTIONS ON SYSTEMS, MAN, AND CYBERNETICS-PART C: APPLICATIONS AND REVIEWS, VOL. 35, NO. 3, AUGUST 2005
  - [26.] Van Trees, Harry L., "Detection, Estimation, and Modulation Theory", Part I, John Wiley and Sons, 1968.
  - [27.] Qing Yan and Rick S. Blum, "Distributed Signal Detection under the Neyman-Pearson Criterion", EECS Department Lehigh University Bethlehem, PA 18015
  - [28.] M. SOLTANE, N. DOGHMANE, N. GUERSI. "State of the Art: Signature Biometrics Verification", BRAIN. Broad Research in Artificial Intelligence and Neuroscience. Vol 1, N 2, Romania 2010. <http://www.edusoft.ro/brain>
  - [29.] M. SOLTANE, B. MIMEN, "Soft Decision Level Fusion Approach to a Combined Behavioral Speech Signature Biometrics Verification", International Journal of Signal Processing, Image Processing and Pattern Recognition– IJSIP, Vol.5, No. 5 South Korea (March 2013). [http://www.sersc.org/journals/IJSIP/vol5\\_no5.php](http://www.sersc.org/journals/IJSIP/vol5_no5.php)

- [30.] J. Richiardi, J. Fierrez-Aguilar, J. Ortiga-Garcia and A. Drygajlo, "On-line signature verification resilience to packet loss in IP networks", in Proc. 2<sup>nd</sup> COST-275 Workshop on Biometrics on the Internet: Fundamentals, Advances and Applications. University of Vigo, Vigo-Spain 25-26 March 2004.
- [31.] Youssef Elmir, Youssef Elmir, Zakaria Elberrichi and Réda Adjoudj, "Multimodal Biometric Using a Hierarchical Fusion of a Person's Face, Voice, and Online Signature". Journal of Information Processing Systems (JIPS), Volume:10, No:4, Page:555~567, Year:2014
- [32.] Mandeep Kaur, Akshay Girdhar and Manvjeet Kaur, "Multimodal Biometric System Using Speech and Signature Modalities", International Journal of Computer Applications (IJCA) Volume 5– No.12, August 2010



ANNALS of Faculty Engineering Hunedoara – International Journal of Engineering



copyright © UNIVERSITY POLITEHNICA TIMISOARA, FACULTY OF ENGINEERING HUNEDOARA,  
5, REVOLUTIEI, 331128, HUNEDOARA, ROMANIA  
<http://annals.fih.upt.ro>