

# Use of Identification Trial Statistics for Combination of Biometric Matchers

Sergey Tulyakov and Venu Govindaraju, *Fellow, IEEE*

**Abstract**—Combination functions typically used in biometric identification systems consider as input parameters only those matching scores which are related to a single person in order to derive a combined score for that person. We discuss how such methods can be extended to utilize the matching scores corresponding to all persons. The proposed combination methods account for dependencies between scores output by any single participating matcher. Our experiments demonstrate the advantage of using such combination methods when dealing with large number of classes, as is the case with biometric person identification systems. The experiments are performed on the NIST BSSR1 dataset and combination methods considered include likelihood ratio, neural network and weighted sum.

**Index Terms**—Combination of classifiers, biometric identification systems.

## I. INTRODUCTION

**B**IOMETRIC applications operate in two modes: verification (1:1) mode and identification (1:N) mode. Common approaches to combining biometrics for (1:N) identification applications are usually a simple iterative use of the (1:1) verification system. The combined score assigned to a particular enrolled person is obtained as a function of the scores assigned to that person by all the biometric matchers in either modes of operation. However, in the identification mode additional information is available for deriving the combined score for any person in the database of enrollees. This additional information is available from the matching scores returned for the enrollees other than the target person.

We consider  $M$  multiple biometric matchers used to produce  $MN$  matching scores (Figure 1), where  $N$  is the number of enrolled persons. We assume that  $M$  is small and  $N$  is large. Each biometric matcher in such a setting is equivalent to a classifier assigning matching scores to each of the  $N$  classes or persons. And the combination of biometric matchers can be viewed as a classifier combination problem with a large number of classes.

Combination methods can be categorized based on the construction properties of the combination functions  $f$ . When methods use a single common combination function, they are called *class generic* methods. When each class has its own combination function, so that the combined scores are calculated differently for different classes, the methods are called *class specific*.

*Local methods* take as parameters only the  $M$  match scores related to a particular class (single column in Figure 1)

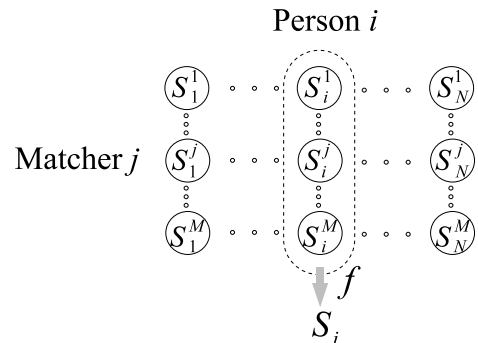


Fig. 1. The set of scores available for combinations in identification systems includes all  $MN$  matching scores from  $M$  matchers and assigned to all  $N$  persons. The combination functions  $f$  usually only utilize the set of scores related to one person  $i$  in order to calculate the combined matching score for this person.

whereas the *global methods* consider the whole set of  $MN$  match scores (all columns in Figure 1) to derive the combined score for any one class. In this paper we explore *global methods* whose combination functions use the additional information (all columns) when computing the integrated score for each person.

When classifiers deal with a small number of classes, the dependencies between the scores assigned to different classes can be learned and used for combination purposes. For example, Xu et al. [1] used class confusion matrices for deriving belief values and integrated these values into combination algorithms in the digit classification problem. This algorithm has class specific and global combination functions. It is the most general type of combination method allowing optimal performance. However, learning class dependencies requires significant number of training samples for each class. Such data might not be available for 1:N identification mode systems, where usually a single template is enrolled for each person. In addition, the database of enrolled persons can be frequently changed making learning class relationships infeasible.

As a consequence, combination approaches in 1:N identification systems have considered only the local methods even when all the  $MN$  scores are available. In this paper we investigate the question of whether it is possible to improve the performance of the identification system by using all the  $MN$  matching scores for deriving the combined score for each person [2], [3].

Manuscript received September 14, 2007, revised February 19, 2008

The authors are with the Center for Unified Biometrics and Sensors, State University of New York at Buffalo, USA (email: tulyakov@cubs.buffalo.edu; govind@cubs.buffalo.edu).

## II. PREVIOUS WORK IN IDENTIFICATION SYSTEM COMBINATIONS

Traditionally, two types of biometric person authentication systems are defined - verification (1:1) and identification (1:N) systems. It is usually implied that verification systems have only the matching scores related to one enrolled person available to the combination method. However, it is possible that a verification system additionally uses matching scores related to other persons. For example, in [4] authors performed 'identification based verification' by utilizing matching scores of other enrolled peoples while making verification decision on a particular person.

In order to avoid confusion, we define an identification system as a system which provides matching scores for all  $N$  enrolled persons. As in [4], such systems can operate in verification mode also. An identification system is operating in identification mode if its purpose is to classify an input as belonging to any of  $N$  classes or persons. We assume that the classification decision is performed by applying the  $\arg \max$  operator to the  $N$  combined scores:

$$C = \arg \max_{1 \leq i \leq N} S_i$$

The correct identification rate, that is the frequency of correctly finding the true class of the input, is the natural measure of performance in this case, and we will use it in our experiments. Note, that there could be other performance measures for identification mode operation, such as Rank Probability Mass, Cumulative Match Curve [5].

When an identification system operates in verification mode we can distinguish two classes: genuine and impostor verification attempts. The decision to accept is based on comparing a combined score of a claimed person identity  $i$ ,  $S_i$ , to some threshold  $\theta$ :  $S_i > \theta$ . The common way to describe the system performance in such two-class problems is to construct ROC curves showing the dependencies of errors on threshold  $\theta$  (or DET curve [5]).

If we have a combination algorithm for verification systems, it can be sequentially applied for all persons to operate in the identification mode [6]. However, this approach does not utilize dependencies between scores output by a single matcher, i.e. the dependencies between the scores along the rows in the score matrix of Figure 1. It is essentially a local method which considers only a single column of scores as input parameters to combination functions. Most combination algorithms used in biometric applications are of this type and sometimes are also user specific [7], [8].

We present here previous approaches which utilize score dependencies in the identification mode.

### A. Rank Based Combinations

T.K. Ho has used classifier combinations on the ranks of the scores instead of scores themselves by arguing that ranks provide more reliable information about a class being genuine [9], [10]. Thus, if the input image has low quality, then the genuine score, as well as the impostor scores will be low. Combining low score for genuine class with other scores could confuse a combination algorithm, but the rank of the

genuine class remains to be a stable statistic, and combining this rank with other ranks of the genuine class should result in true classification. Brunelli and Falavigna [11] considered a hybrid approach where traditional combination of matching scores is fused with the rank information in order to achieve identification decision. Hong and Jain [12] consider ranks, not for combination, but for modeling or normalizing the output score of a classifier. Behavior-Knowledge Space combination methods [13] are also based on ranks. Saranli and Demirekler [14] provide additional references for rank based combination methods.

The problem with rank based methods is that the score information is lost. Indeed, the best score can be much better than second best score, or it could be only slightly better, but score ranks do not reflect this difference. It would be desirable to have a combination method which retains the score information as well as the rank information.

### B. Score normalization approaches

Usually score normalization [15] refers to transformation of scores based on a classifier's score model learned during training, and each score is transformed individually using such a model. Thus the other scores output by a matcher during the same identification trial (rows in the score matrix of Figure 1) are not taken into consideration. If these normalized scores are later combined class-wise (column-wise), then score dependence is not accounted for by the combination algorithm.

Some score normalization techniques can use a set of identification trial scores output by the classifier. For example, Kittler et al. [16] normalize each score by the sum of all the other scores before combination. Similar normalization techniques are used in Z(zero)- and T(test)- normalizations [17], [18]. Z-normalization uses impostor matching scores to produce a class specific normalization. Z-normalization does not include the set of identification trial scores (rows in Figure1), and thus does not utilize the score dependency. On the other hand, T-normalization uses a set of scores produced during a single identification trial by utilizing statistics of mean and variance of the identification score set. Note that T-normalization is a predetermined routine with no training. Still, using this simple kind of score modeling turns out to be quite useful; for example, [19] argued for applying T-normalizations in face verification. There is also a counterargument [20] that useful classification information gets lost during such normalizations.

Score normalization techniques have been well developed in the speaker identification literature. Cohort normalizing method [21], [22] considers a subset of enrolled persons close to the current test person in order to normalize the score for that person by a log-likelihood ratio of the genuine (current person) and impostor (cohort) score density models. Auckenthaler et al.[17] separated cohort normalization methods into cohorts found during training (constrained) and cohorts dynamically formed during testing (unconstrained cohorts). Normalization by constrained cohorts utilizes only one matching score of each classifier and thus does not consider score dependencies. On the other hand, normalization

by unconstrained cohorts potentially uses all scores of all classifiers.

### III. COMPLEXITY TYPES OF CLASSIFIER COMBINATIONS

This section describes four types of combination methods and their requirements of training data. Ultimately, the problem characteristics and the availability of training scores determine the type of combination method which is appropriate for a particular problem.

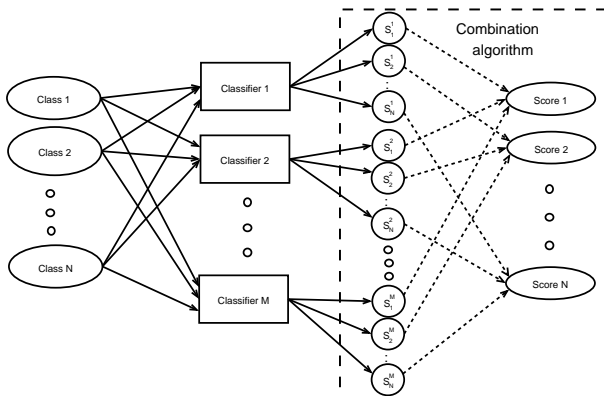


Fig. 2. Classifier combination takes a set of  $s_i^j$  - score for class  $i$  by classifier  $j$  and produces combination scores  $S_i$  for each class  $i$ .  $i$  is the index for the  $N$  classes and  $j$  is the index for the  $M$  classifiers

Figure 2 gives a different view of the problem of integrating scores in identification systems, for the purpose of formally categorizing the combination methods. The combination functions of local methods have a reduced parameter set (as connections in Figure 2 show), and many well known combination methods (e.g. weighted sum of scores) fall into this category. A fully connected artificial neural network accepting  $MN$  input parameters and having  $N$  output parameters would present an example of the most general, class specific and global combination function algorithm [1], [23]. The disadvantage of this more general approach is that it requires very large amount of training data, which might not be always available in identification systems.

#### A. Types of Combination Methods

We develop here a formal framework for combination methods further categorizing the *local* and *global* combination functions that are required to be trained. The first two categories correspond to *local* and the remaining two correspond to *global* methods.

- 1) Low complexity methods:  $S_i = f(\{s_i^j\}_{j=1, \dots, M})$ . Methods of this type require only one combination function to be trained, and the combination function takes as input scores for one particular class as parameters. It represents class generic and reduced parameter set combination functions.
- 2) Medium complexity I methods:  $S_i = f_i(\{s_i^j\}_{j=1, \dots, M})$ . Methods of this type have separate score combining functions for each class and each such function takes, as input parameters, only the scores related to its class.

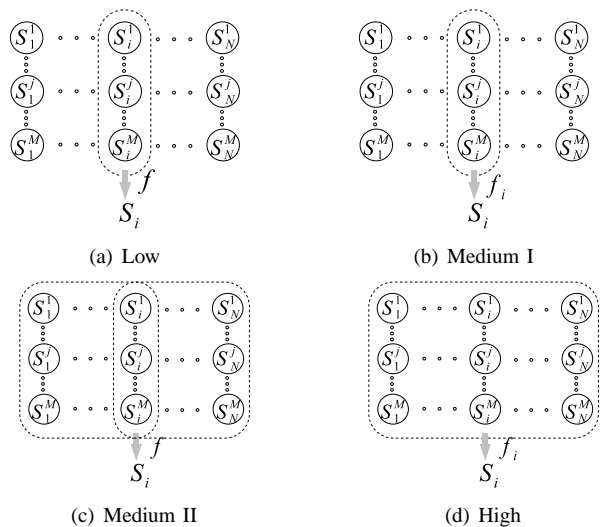


Fig. 3. The range of scores considered by each combination type and combination functions.

It represents class specific and reduced parameter set combination functions.

- 3) Medium complexity II methods:

$$S_i = f(\{s_i^j\}_{j=1, \dots, M}, \{s_k^j\}_{j=1, \dots, M; k=1, \dots, N, k \neq i}).$$

Methods of this type take as parameters not only the scores related to the same class, but all output scores of classifiers. Combination scores for each class are calculated using the same function, but scores for class  $i$  are given a special place in the parameter list. Applying function  $f$  for different classes effectively means permutation of the function's parameters. These combination functions are class generic and use the whole parameter set.

- 4) High complexity methods:

$S_i = f_i(\{s_k^j\}_{j=1, \dots, M; k=1, \dots, N})$ . Functions calculating final scores are different for all classes, and they take as parameters all the scores output by the base classifiers. This represents class specific and whole parameter set combination functions.

We can illustrate the different combination types using the matrix score representation (Figure 1) as well. Each row corresponds to a set of scores output by a particular classifier, and each column corresponds to scores assigned by classifiers to a particular class. The illustration of each combination type functions is given in Figure 3. In order to produce the combined score  $S_i$  for class  $i$ , low complexity methods (Figure 3a) and medium I complexity (Figure 3b) combinations consider only those classifier scores which are assigned to class  $i$  (column  $i$ ), reflecting the property of local combination functions. Medium II (Figure 3c) and high complexity (Figure 3d) methods consider all the scores output by classifiers for calculating a combined score  $S_i$  for class  $i$ , reflecting the property of global combination functions.

Low (Figure 3a) and medium II (Figure 3c) complexity methods have the same *class generic* combination functions  $f$  irrespective of the class for which the score is calculated. Note that medium II complexity type methods have scores related

to a particular class in a special consideration as indicated by the second ellipse around these scores. We can think of these combinations as taking two sets of parameters - scores for a particular class, and all other scores. The important property in these methods is that the combination function  $f$  is the same for all classes, but the combined scores  $S_i$  differ, since we effectively permute function inputs for different classes. Medium I (Figure 3b) and high (Figure 3d) complexity methods have *class specific* combination functions  $f_i$  trained differently for different classes.

It is interesting to compare our combinations types with previous categorization of combination methods by Kuncheva et al.[24], who refer to the score matrix as ‘decision profile’ and ‘intermediate feature space’. Kuncheva’s work also separates combinations into ‘class-conscious’ set which corresponds to the union of ‘low’ and ‘medium I’ complexity types, and ‘class-indifferent’ set which corresponds to the union of ‘medium II’ and ‘high’ complexity types. The continuation of this work [25] gave an example of the weighted sum rule having three different numbers of trainable parameters (and accepting different numbers of input scores), which correspond to ‘low’, ‘medium I’ and ‘high’ complexity types.

In contrast to Kuncheva’s work, our categorization of combination methods is more general since we are not limiting ourselves to simple combination rules like the weighted sum rule. Further, we consider an additional category of ‘medium II’ type. An example of ‘medium II’ combination is the two step combination algorithm where in the first step the scores produced by a particular classifier are normalized (with possible participation of all scores of this classifier), and in the second step, scores are combined by a function from the ‘low’ complexity type. Thus scores in each row are combined first, and then the results are combined columnwise in the second step.

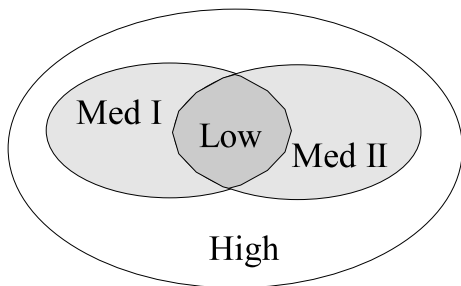


Fig. 4. The relationship diagram of different combination complexity types.

Figure 4 illustrates the relationships between the types of combination methods. Medium complexity types are subsets of high complexity type, and the set of low complexity methods is exactly the intersection of sets of medium I and medium II combination methods. In order to avoid a confusion in terminology we will henceforth assume that a combination method belongs to a particular type only if it belongs to this type and does not belong to the more specific type.

In [26] we provide a description of these complexity types using the concept of VC (Vapnik-Chervonenkis) dimension [27]. The ability to use VC dimension for characterization

of different combination types justifies our usage of term ‘complexity types’.

Higher complexity methods can potentially produce better classification results since more information is used. However, the availability of training samples limits the types of possible combinations to choose from. Thus the choice of combination method in any particular application is a trade-off between classifying capabilities of the combination functions and the availability of sufficient training samples. Different generic classifiers such as neural networks, decision trees, etc., can be used for combination within each complexity class. From the perspective of our framework, the main effort in solving the classifier combination problem consists of identifying the complexity type and modifying the generic classifiers (if needed) to compensate for a mismatched function complexity type for reasons of inadequate training data.

The biometric person authentication systems we experimented with in this research have a high number of enrolled classes (persons)  $N$  and a small number of classifiers (biometric face and fingerprint matchers)  $M$ . Most combinations methods described in the literature for biometric applications are of low complexity type. In this work we are interested in exploiting higher complexity combinations. We will derive combinations rules of medium II complexity type which are analogous to the traditional likelihood ratio, neural network and weighted sum combinations of the low complexity type. Our experiments on large biometric score sets confirm that the medium II complexity combinations have better performance than their counterparts of low complexity.

Both identification and verification modes of operation can utilize combinations of all four complexity types. In our experiments we compare combination methods of low and medium II complexity types and report performance for both, verification and identification, modes of operation.

#### IV. DERIVATION OF COMBINATION RULES USING IDENTIFICATION TRIAL STATISTICS

In this section we present different combination methods of medium II complexity type utilizing the statistics of the identification trial score set for normalization purposes. Our goal is to theoretically derive an optimal combination algorithm with the assumption that the joint densities of the scores and score set statistics are known. We will also discuss the application of the ‘background model’ (used in speaker identification [18]) and its relation to our approach.

##### A. Likelihoods With the Use of Identification Trial Score Set Statistics

Suppose we combine  $M$  independent classifiers, and each classifier outputs  $N$  dependent scores. The optimal combination algorithm is the Bayesian classifier which accepts these  $NM$  scores and chooses the class which maximizes the posterior class probability. Thus the goal of the optimal combination method is to find

$$\arg \max_k P(C_k | \{s_i^j\}_{i=1, \dots, N; j=1, \dots, M})$$

Term  $C_k$  refers to the fact that the class  $k$  is the genuine class. By Bayes theorem

$$P(C_k | \{s_i^j\}_{i=1,\dots,N;j=1,\dots,M}) = \frac{p(\{s_i^j\}_{i=1,\dots,N;j=1,\dots,M} | C_k) P(C_k)}{p(\{s_i^j\}_{i=1,\dots,N;j=1,\dots,M})}$$

and since the denominator is the same for all classes, our goal is to find

$$\arg \max_k p(\{s_i^j\}_{i=1,\dots,N;j=1,\dots,M} | C_k) P(C_k)$$

or, assuming all classes have the same prior probability,

$$\arg \max_k p(\{s_i^j\}_{i=1,\dots,N;j=1,\dots,M} | C_k)$$

Given the assumption that classifiers are independent, which means that any subset of scores produced by one classifier is statistically independent from any other subset of scores produced by another classifier, our problem is to find

$$\arg \max_k \prod_j p(\{s_i^j\}_{i=1,\dots,N} | C_k) \quad (1)$$

The goal is to reliably estimate the densities  $p(\{s_i^j\}_{i=1,\dots,N} | C_k)$ , which is a hard task given that the number  $N$  of classes is large and we do not have many samples of each class for training. Since we do not want to construct a class specific combination method, the class indexes can be permuted. Thus all the training samples pertaining to different genuine classes can be used to train only one density,  $p(s_k^j, \{s_i^j\}_{i=1,\dots,N,i \neq k} | C_k)$ . Now  $s_k^j$  is a score belonging to a genuine match, and all other scores  $\{s_i^j\}_{i=2,\dots,N}$  are from impostor matches. In order to keep the problem tractable, instead of  $p(s_k, \{s_i^j\}_{i=1,\dots,N,i \neq k} | C_k)$  we can consider  $p(s_k^j, t_k^j | C_k)$ , where  $t_k^j$  is some statistics of all the other scores besides  $s_k^j$ . The final combination rule for this method is to find

$$\arg \max_k \prod_j p(s_k^j, t_k^j | C_k) \quad (2)$$

As our previous experiments have shown [3] this algorithm does not perform as well as the traditional likelihood ratio combination:

$$\arg \max_k \prod_j \frac{p(s_k^j | C_k)}{p(s_k^j | \overline{C_k})} \quad (3)$$

One reason for the lower performance could be that the score set statistics  $t_k^j$  does not fully reflect the background information for score  $s_k^j$ , whereas the term  $p(s_k^j | \overline{C_k})$  contains such information. For example, the genuine matching scores  $s_k^j$  can be very strong, but located in the region of low probability (both  $p(s_k^j | C_k)$  and  $p(s_k^j, t_k^j | C_k)$  are small), whereas  $p(s_k^j | \overline{C_k})$  could be even smaller, and the likelihood ratio can still succeed. In the next section we will derive a combination rule which combines the use of the score set statistics and background models [21].

## B. Likelihood Ratios with the Use of Identification Trial Score Set Statistics

We consider the posterior class probability, apply Bayes formula as before, but now use the independence of classifiers to decompose the denominator:

$$\begin{aligned} & P(C_k | \{s_i^j\}_{i=1,\dots,N;j=1,\dots,M}) \\ &= \frac{p(\{s_i^j\}_{i=1,\dots,N;j=1,\dots,M} | C_k) P(C_k)}{p(\{s_i^j\}_{i=1,\dots,N;j=1,\dots,M})} \\ &= \frac{\prod_j p(\{s_i^j\}_{i=1,\dots,N} | C_k) P(C_k)}{\prod_j p(\{s_i^j\}_{i=1,\dots,N})} \\ &= P(C_k) \prod_j \frac{p(\{s_i^j\}_{i=1,\dots,N} | C_k)}{p(\{s_i^j\}_{i=1,\dots,N})} \end{aligned} \quad (4)$$

The next step is similar to the step in deriving the algorithm for the background speaker model [18]. We consider the class  $\overline{C_k}$  to mean that some other class is genuine, and decompose  $p(\{s_i^j\}_{i=1,\dots,N}) = P(C_k) p(\{s_i^j\}_{i=1,\dots,N} | C_k) + P(\overline{C_k}) p(\{s_i^j\}_{i=1,\dots,N} | \overline{C_k})$ . By assuming that  $N$  is large and  $P(\overline{C_k}) \gg P(C_k)$ , we can discard the first term and represent 4 as:

$$\frac{P(C_k)}{P(\overline{C_k})^M} \prod_j \frac{p(\{s_i^j\}_{i=1,\dots,N} | C_k)}{p(\{s_i^j\}_{i=1,\dots,N} | \overline{C_k})}$$

Assuming that all classes have the same probability of occurring ( $P(C_k) = P(C_i)$  and  $P(\overline{C_k}) = P(\overline{C_i})$ ) we obtain the following classifier decision:

$$\arg \max_k \prod_j \frac{p(\{s_i^j\}_{i=1,\dots,N} | C_k)}{p(\{s_i^j\}_{i=1,\dots,N} | \overline{C_k})} \quad (5)$$

In comparison with decision 1 of the previous section we have an additional density  $p(\{s_i^j\}_{i=1,\dots,N} | \overline{C_k})$ . This density can be viewed as a background of impostors for the genuine class  $C_k$ . As research in speaker identification suggests [21], utilizing such a background model is beneficial for system performance.

We estimate the ratios of equation 5 by additional modeling of  $p(\{s_i^j\}_{i=1,\dots,N} | \overline{C_k})$ . We use an approach similar to the previous section to estimate this density as  $p(s_k^j, t_k^j | \overline{C_k})$  with  $t_k^j$  - the joint density of impostor scores  $s_k^j$  and corresponding identification trial statistics  $t_k^j$ . The final combination rule is then,

$$\arg \max_k \prod_j \frac{p(s_k^j, t_k^j | C_k)}{p(s_k^j, t_k^j | \overline{C_k})} \quad (6)$$

The use of the identification trial score set statistics considers  $p(s_k^j, t_k^j | C_k)$  and  $p(s_k^j, t_k^j | \overline{C_k})$  instead of  $p(s_k^j | C_k)$  and  $p(s_k^j | \overline{C_k})$ , and the background model considers  $p(s_k^j, t_k^j | \overline{C_k})$  or  $p(s_k^j | \overline{C_k})$  in addition to  $p(s_k^j, t_k^j | C_k)$  or  $p(s_k^j | C_k)$ . Thus, the use of the identification trial score statistics differs from the background model in being able to account for dependencies of scores in identification trials by using the statistic  $t_k^j$ .

Note, that traditional likelihood ratio (Eq. 3) is the optimal combination method for low complexity combinations operating in verification mode (see [28]). Thus, its extension by Eq. 6 should provide a good combination method of medium II complexity type for verification mode operations.

### C. Statistics of Identification Trial Scores

The important question which we have to decide is what particular identification trial score statistics  $t_k^j$  will be most suitable to replace the set of scores  $\{s_i^j\}_{i=1,\dots,N,i\neq k}$ . The likelihood ratio incorporating score statistics (Eq. 6) will be more discriminating than the traditional likelihood ratio (Eq. 3) if  $t_k^j$  provides an information on whether the score  $s_k^j$  is genuine or impostor. We use the term "identification model" to denote a particular way of choosing identification trial score set statistics  $t_k^j$  and using this statistics together with scores  $s_k^j$ .

One of the identification models we previously presented is the second best score model [29], where statistics  $t_k^j = sbs(s_k^j)$  is calculated as the best score in the set  $\{s_i^j\}_{i=1,\dots,N,i\neq k}$  ("second best" after  $s_k^j$ ). We can reason that if second best score is big (e.g. bigger than current score  $s_k^j$ , so  $s_k^j$  is not the best score), then we have less confidence that  $s_k^j$  is a genuine score, and more confidence that this is impostor score. And if it is small (so  $s_k^j$  is big relative to all other scores), we have more confidence that  $s_k^j$  is a genuine. Originally, we used  $sbs(s_k^j)$  for accepting first-choice decisions in open-set identification systems [29]. In this case,  $sbs(s_k^j)$  exactly coincides with second best score of the identification trial score set.

T-normalization can be considered as another identification model. It is expressed as a transformation of all matching scores  $s_k^j$  by the formula

$$s_k^j(l) \rightarrow \frac{s_k^j(l) - \mu^j(l)}{\sigma^j(l)}$$

where  $\mu^j(l)$  and  $\sigma^j(l)$  are correspondingly the mean and the standard deviation of the set of scores produced by matcher  $j$  during the identification trial  $l$ . In contrast to second best score model, T-normalization uses different statistics -  $\mu^j$  and  $\sigma^j$  which are the same for all scores  $s_k^j$  in the current identification trial, and it does predetermined transformation using these statistics.

Clearly, there might be many variations on calculating statistics  $t_k^j$  - it may or may not be dependent on  $k$ , it might include mean, variance, n-th ranked score or any other statistics of a score set. It seems that for different applications the most useful statistics will be different, and it would be desirable to have an automatic way of determining useful score statistics. In our experiments we limited ourselves to only using second best score statistics and T-normalization.

One approach to choose a best statistics of identification trial score sets is to look at the dependence between genuine and impostor scores. In order to verify the dependence of match scores we measured the correlation between the genuine score and different statistics of the sets of impostor scores. Table I contains a small part of measured correlations corresponding to  $first_{imp}$  - 1st ranked impostor score,  $second_{imp}$  - 2nd ranked impostor score, and  $mean_{imp}$  - the mean of impostor scores. As results of Table I show, the scores produced by real life classifiers are indeed dependent.

The correlations between genuine and impostor set statistics indicate the usefulness of a given statistics - bigger correlation means that this statistics is better able to predict whether the

Matchers	$first_{imp}$	$second_{imp}$	$mean_{imp}$
li	0.3164	0.3400	0.2961
ri	0.3536	0.3714	0.3626
C	0.1419	0.1513	0.1440
G	0.1339	0.1800	0.1593

TABLE I  
CORRELATIONS BETWEEN  $s_{gen}$  AND DIFFERENT STATISTICS OF THE IMPOSTOR SCORE SETS PRODUCED DURING IDENTIFICATION TRIALS IN NIST BSSR1 DATA.

score is genuine or not. So we might want to calculate such correlations for many different statistics and choose statistics with bigger correlations. Second best score statistics seems to provide a good prediction on the strength of genuine score, and this is additional reason we used it in our experiments. Note, that  $sbs(s_k^j)$  used in our experiments is calculated with respect to  $s_k^j$  and if  $s_k^j$  is an impostor score it might not be  $first_{imp}$  or  $second_{imp}$ . During testing we do not know what the exact set of impostors is, so instead of  $first_{imp}$  or  $second_{imp}$  we are forced to use  $sbs(s_k^j)$ .

### D. Combinations of Dependent Classifiers

The combination algorithms presented in the previous two sections deal with independent classifiers. How should we address dependent classifiers?

By looking at the combination equations 1 and 6 we can see that each classifier contributes terms  $p(\{s_i^j\}_{i=1,\dots,N}|C_k)$  and  $\frac{p(\{s_i^j\}_{i=1,\dots,N}|C_k)}{p(\{s_i^j\}_{i=1,\dots,N}|C_k)}$  correspondingly to the combination algorithm. Thus one can conclude that it is possible to model the same terms for each classifier with the help of identification trial score statistics,  $p(s_k^j, t_k^j|C_k)$  and  $\frac{p(s_k^j, t_k^j|C_k)}{p(s_k^j, t_k^j|C_k)}$ , and then combine them by some other trainable function.

Note that any relationships between scores  $s_{i_1}^{j_1}$  and  $s_{i_2}^{j_2}$  where  $i_1 \neq i_2$  and  $j_1 \neq j_2$  will be essentially discarded. This seems to be inevitable for the current complexity type of combinations - medium II. If we wanted to account for such relationships, we would need class-specific combination functions, and thus higher complexity combination methods.

Another way to construct combinations of medium II complexity type for dependent classifiers is presented in the next section.

### E. Normalizations Followed by Combinations and Single Step Combinations

Figure 5 represents in graphical form the type of combinations we have presented thus far. All these combinations consist of two steps. In the first step, each score is normalized by using other scores output by the same matcher. In the second step, normalized scores are combined using a predetermined or trained combination function.

Score normalization based on modeling the joint densities of scores and statistics,  $p(s_k^j, t_k^j|C_k)$  and  $\frac{p(s_k^j, t_k^j|C_k)}{p(s_k^j, t_k^j|C_k)}$ , might not correctly represent original terms of Eqs. 1 and 5,  $p(\{s_i^j\}_{i=1,\dots,N}|C_k)$  and  $\frac{p(\{s_i^j\}_{i=1,\dots,N}|C_k)}{p(\{s_i^j\}_{i=1,\dots,N}|C_k)}$ . Approximating densities might also be unreliable if few statistics are used.

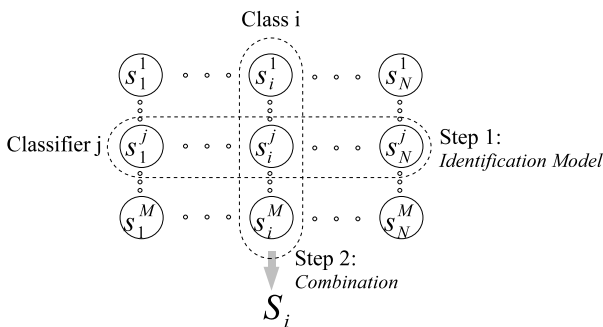


Fig. 5. 2-step combination method utilizing identification model.

However, it is not necessary to have the two steps for combinations. The contribution of the particular classifier  $j$  to the whole combination algorithm's output for class  $i$  is calculated only from score  $s_i^j$  and statistic  $t_i^j$ . Figure 6 illustrates how scores and statistics from all the participating classifiers could be combined in a single combination step.

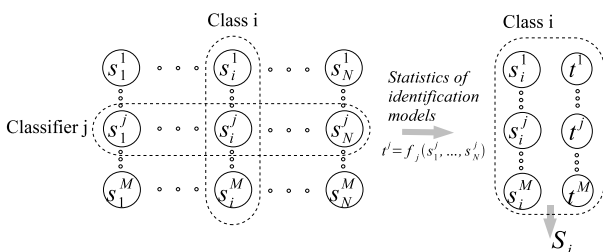


Fig. 6. 1-step combination method utilizing identification model.

In the algorithm described by Figure 6 the statistics  $t^j$  are extracted for each classifier  $j$  using its output scores by a predetermined and non-trainable algorithm. The scores related to a particular class and statistics are combined together by a trainable function. This combination function is not class-specific and is easily trainable. This type of combinations are of medium II complexity type. In comparison, for the low complexity type combinations only the scores for a particular class are combined, and statistics from other classes are not considered.

#### F. Neural Network and Weighted Sum Combinations Using Second-Best Score

As an example of single step combinations, we consider two combination methods incorporating the second best score statistics: the neural network and the weighted sum rule.

Traditional neural network corresponding to low complexity combination type can be represented as a function  $S_i = f(s_i^1, \dots, s_i^M)$ . Following the diagram of Figure 6, the neural network combination of medium II complexity type will have the form  $S_i = f(s_i^1, sbs(s_i^1), \dots, s_i^M, sbs(s_i^M))$ . We used multilayer perceptron trained by a traditional backpropagation method and optimizing MSE.

The traditional weighted sum combination without the use of second-best score ('weighted sum local') is a low complexity combination which combines  $M$  scores from  $M$  biometric

matchers assigned to a particular class  $i$ :

$$S_i = w_1 s_i^1 + \dots + w_M s_i^M \quad (7)$$

The weighted sum rule with the sbs model ('weighted sum global') combines scores as well as statistics of score sets:

$$S_i = w_1 s_i^1 + w_2 sbs(s_i^1) + \dots + w_{2M-1} s_i^M + w_{2M} sbs(s_i^M) \quad (8)$$

Weighted sum rule can be specifically trained to maximize the correct identification rate in identification mode of operation[28]. However, it is not optimal for the verification mode. Thus, we will test the performance of the weighted sum rule with and without the second-best score model modification for the identification mode operation only. The neural network, on the other hand, might not be optimal for identification mode due to MSE minimization criteria, but gives an output approximating likelihood ratio. We give the performance of neural network method for both identification and verification modes.

## V. EXPERIMENTS

We have used the biometric matching score set BSSR1 distributed by NIST[30]. This set contains matching scores for a fingerprint matcher and two face matchers 'C' and 'G'. Fingerprint matching scores are given for the left index 'li' finger matches and right index 'ri' finger matches. For each combination method we performed six sets of experiments on combining any two pairs of scores: 'C' & 'G', 'li' & 'ri', 'li' & 'C', 'li' & 'G', 'ri' & 'C', and 'ri' & 'G'.

Although the BSSR1 score set has a subset of scores obtained from the same physical individuals, this subset is rather small - 517 identification trials with 517 enrolled persons. We use bigger subsets of fingerprint and face matching scores of BSSR1 by creating virtual persons; the fingerprint scores of a virtual person come from one physical person and the face scores come from a different physical person. Note, that for pairs of face scores and for pairs of fingerprint scores, we retain the correspondence of scores to real persons as specified in the database. The scores are not reused, and thus we are limited to the maximum number of identification trials - 6000 and the maximum number of classes, or enrolled persons, - 3000. Some enrollees and some identification trials had to be discarded due to enrollment errors. We use a bootstrap testing procedure: for 100 iterations, we randomly split the data in two parts - 2991 identification trials with 2991 enrolled persons in each part used as separate training and testing sets. The results of 100 training/testing iterations are averaged at the end.

In order to achieve good performance of training algorithms all the scores were normalized using simple min-max algorithm to interval [0,1]. When we used T-normalization, additional min-max normalization was performed after it.

#### A. Description of Used Algorithms

The goal of our experiments is to compare three general architectures for classifier combination - traditional low complexity combinations which do not use any identification model, medium II complexity combinations based on T-normalized scores and medium II complexity combinations

using second best score model. Three types of learning algorithms are used in the experiments - likelihood ratio, neural network and weighted sum. In order to make the comparison objective we utilize each learning algorithm in each of three architectures. Each classifier in traditional and T-normalization methods supply only a single score and the learning function depends on two input parameters:  $f(s^1, s^2)$ . On the other hand, second best score model has additional score statistics  $sbs(s^1)$  and  $sbs(s^2)$  and the learning function depends on four input parameters:  $f(s^1, sbs(s^1), s^2, sbs(s^2))$ .

For likelihood ratio combinations we estimate score densities using the Parzen window method with Gaussian kernels. The kernel width is determined by the maximum likelihood method. We use only 1000 identification trial scores for reconstructing densities, and the remainder of the training set (2991-1000 trials) is used for validating kernel widths. Note, that for each identification trial there is 1 genuine score and 2990 impostor scores. In order to make our implementation faster, we only used a single random impostor score from a trial for training. We did not utilize the statistical independence of data when combining matchers of different modalities, and in all experiments we approximated either two dimensional densities of genuine and impostor scores -  $p_{gen}(s^1, s^2)$  and  $p_{imp}(s^1, s^2)$ , or four dimensional densities -  $p_{gen}(s^1, sbs(s^1), s^2, sbs(s^2))$  and  $p_{imp}(s^1, sbs(s^1), s^2, sbs(s^2))$ .

The neural network is multilayer perceptron trained by backpropagation algorithm. The neural network has two hidden layers with 8 and 9 nodes and an output layer with 1 node in all cases. The input layer has 2 nodes for traditional training (no identification model) and T-normalized training, and 4 nodes for training with second best score model. As for likelihood ratio method, we used 1000 training samples (1 genuine and one random impostor score from identification trial) for backpropagation training and remaining 2991 – 1000 training samples for validation. The training was stopped when the MSE on the validation set achieved minimum.

For the weighted sum methods we need to find the optimal weights maximizing the number of correct identification trials on the training sets. Though there exists previous research proposing solutions for this problem (e.g. [31], [32]), it deals with the case of small number of classes and is not directly applicable to our case. The key idea of learning algorithms minimizing classification error is to replace the discrete misclassification cost function with some smooth approximation in order to be able to take a derivative of the cost function and perform gradient descent optimization. For our experiments we implemented a simpler approach of random modification of weights and accepting new weights if classification performance improves. Though our approach takes more training time than gradient descent method would have taken, it does not depend on the smoothing parameters and it is sufficiently fast.

### B. Performance in Identification Operating Mode

Table II shows the obtained correct identification rate for experiments with neural network and weighted sum combination methods. The correct identification means that the genuine

Matchers	NN	NN+T	NN+sbs	WS	WS+T	WS+sbs
C & G ( $\sigma$ )	83.49 (0.65)	83.59 (0.84)	<b>83.86</b> (0.62)	84.51 (0.50)	84.53 (0.50)	<b>84.85</b> (0.50)
li & ri ( $\sigma$ )	95.12 (0.30)	95.11 (0.30)	<b>95.17</b> (0.29)	95.11 (0.29)	<b>95.13</b> (0.32)	95.02 (0.32)
li & C ( $\sigma$ )	96.44 (0.93)	<b>97.13</b> (0.24)	96.21 (0.78)	97.15 (0.23)	97.17 (0.23)	<b>97.19</b> (0.25)
li & G ( $\sigma$ )	95.38 (0.35)	94.65 (0.80)	<b>95.73</b> (0.43)	95.38 (0.30)	95.28 (0.26)	<b>96.12</b> (0.29)
ri & C ( $\sigma$ )	97.51 (0.63)	<b>98.10</b> (0.17)	97.39 (0.41)	98.11 (0.16)	98.10 (0.17)	<b>98.16</b> (0.22)
ri & G ( $\sigma$ )	96.69 (0.29)	96.09 (0.54)	<b>97.03</b> (0.26)	96.85 (0.23)	96.76 (0.21)	<b>97.29</b> (0.25)

TABLE II  
CORRECT IDENTIFICATION RATES OF COMBINATIONS IN IDENTIFICATION SYSTEMS. THE STANDARD DEVIATIONS OF THESE RATES ESTIMATED FROM BOOTSTRAP SAMPLES ARE GIVEN IN PARENTHESES.

combined score was better than 2990 impostor combined scores (there is a total of 2991 enrollees). In this table, ‘NN’ is the traditional neural network combination method of low complexity type, ‘NN+T’ is the neural network operating on T-normalized scores and ‘NN+sbs’ is the neural network augmented with the second-best score model. Similarly, ‘WS’ is the traditional weighted sum combination of Eq. 7, ‘WS+T’ is the weighted sum operating on T-normalized scores and ‘WS+sbs’ is the weighted sum combination augmented with the second-best score model of Eq. 8.

We also provided the CMC graphs showing the performance of ‘WS’, ‘WS+T’ and ‘WS+sbs’ methods in Figure 7. As we discussed in section IV-F, neural network training is not optimized for best rank performance and we chose do not include similar CMC graphs for it.

We can see that in all cases, the addition of either the T-normalization or the second-best score statistic into the corresponding low complexity algorithm results in performance improvement. The weighted sum has generally better performance than neural network combination method, and second best score statistics mostly outperforms T-normalization.

### C. Performance in Verification Operating Mode

Although there are examples where score normalization techniques with background models have been used for identification tasks [11], even more applications use such techniques for identification systems operating in verification mode [21], [33], [18]. We also applied the combinations utilizing identification models for biometric person verification tasks. The drawback of using either the background models or the second-best score statistic in verification tasks is that we have to produce not only one match per person and per matcher, but also some set of matching scores for other persons enrolled in the system (or some artificially modeled persons).

Figures 8 and 9 contain the results of experiments when operating in the verification mode for likelihood ratio and neural network combination methods. The ROC performance curves were constructed using combinations of  $2991 \times 100$  (test trials  $\times$  iterations) genuine and impostor score sets. Note that only a single random impostor was used from each test identification trial.



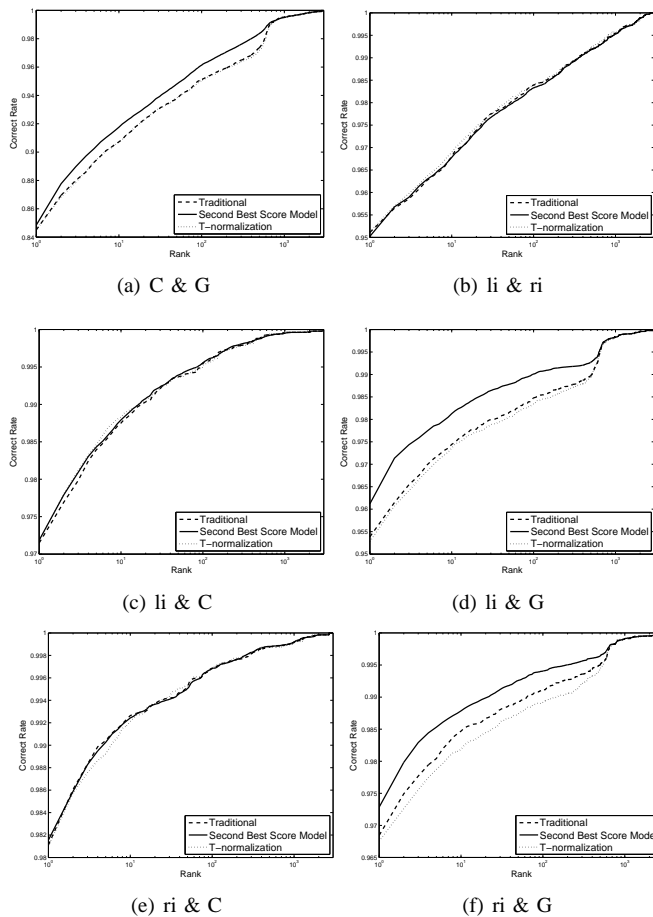


Fig. 7. CMC curves for weighted sum combinations utilizing and not utilizing identification models in identification mode.

We were able to achieve significant improvement in the verification task performance as well, by utilizing the second-best score statistic. The T-normalization is also beneficial, but to a smaller extent in these experiments.

#### D. Dependence of the Performance on the Number of Training Samples

Since the use of second best score model requires learning combination functions with bigger number of parameters, the errors associated with the learning algorithm might increase and negate the benefits of additional model information. In order to clarify the impact of additional training demand on proposed methods, we conducted experiments with different numbers of training samples supplied to the learning algorithm. Figures 10 and 11 present results of these experiments for neural network and likelihood ratio combination methods. Same numbers of training and validation samples are chosen here - 8, 16, ..., 512.

Figure 10 presents the correct identification rate together with 90% confidence intervals estimated from bootstrap samples (extreme 5% of bootstrap samples were discarded from each end) for neural network combinations. The performance results agree with the results presented in Table II - combinations involving 'C' are well handled by T-normalization

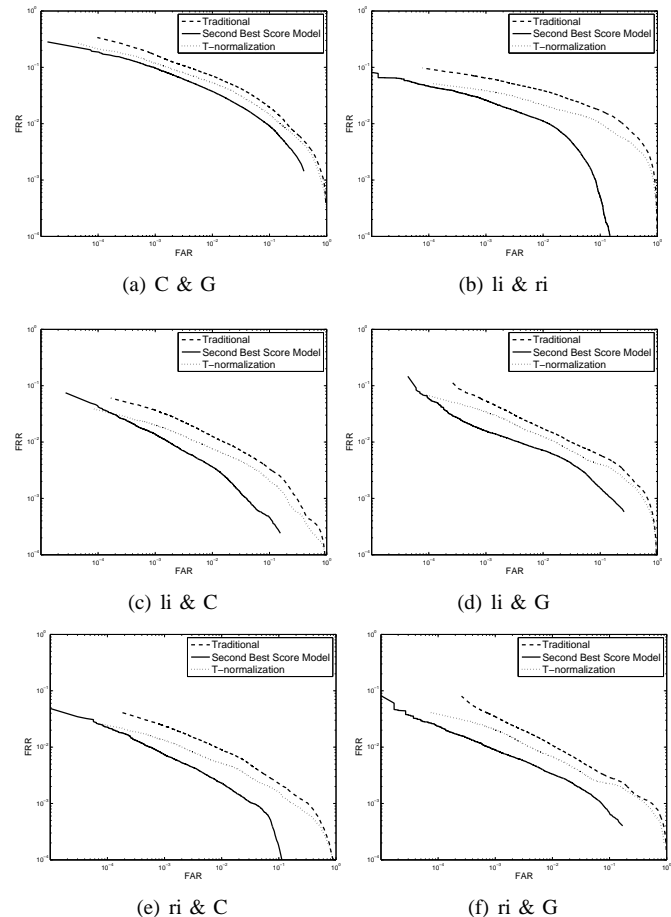


Fig. 8. ROC curves for likelihood ratio combinations utilizing and not utilizing identification models in verification mode.

method, and combinations involving 'G' have better performance when using second best score model. The size of training and validation sets have little impact on average correct identification rate, though it tends to slightly increase with the increasing size of these sets. The impact on spread of rate measurements is more significant. If we want to avoid the accidental bad performance of a particular learned algorithm, we need to ensure that a sufficient number of training samples is used.

Figure 11 presents the equal error rates together with 90% confidence intervals for likelihood ratio combination methods. The increase of the training sample size has big impact on the spread of error rates, and lesser impact on the average error rate. In some cases, the second best score model has worse performance than T-normalization when the number of training samples is small. For bigger number of training samples, second best score model overtakes T-normalization. This observation confirms that learning 4-dimensional score densities for second best score model can result in a worse performance than the approaches requiring learning 2-dimensional densities, such as T-normalization. When the number of training samples is sufficiently large (more than 100 in this case), the density approximations for second best score method are good enough to outperform T-

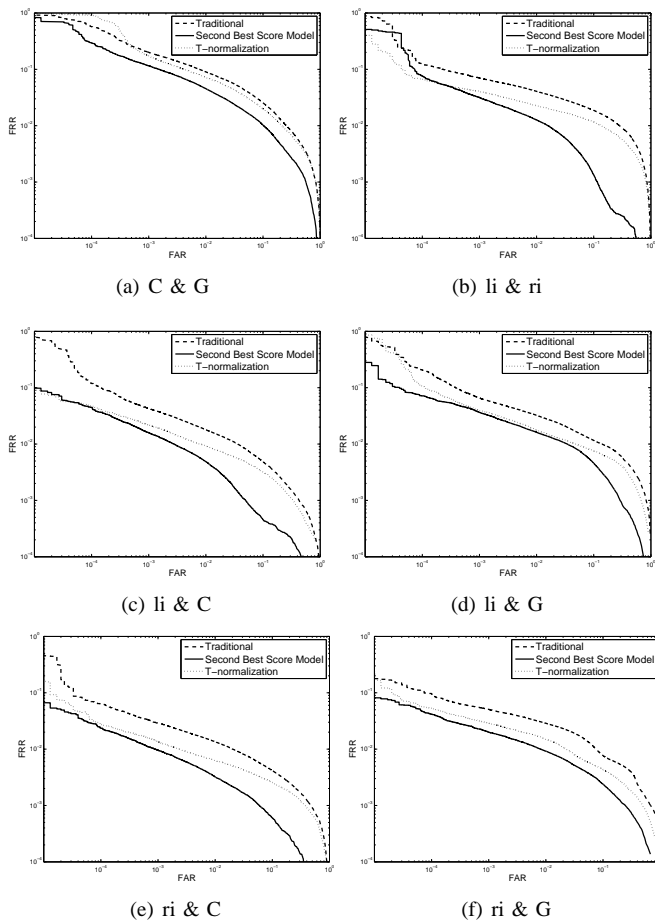


Fig. 9. ROC curves for neural network combinations utilizing and not utilizing identification models in verification mode.

normalization.

## VI. CONCLUSION

We have presented four complexity combination types that originate naturally from the structure of the constructed combination method. We showed the usefulness of differentiating these four combination types for better understanding the problem of classifier combination and for constructing well-performing combination algorithms. We observe that often the algorithms used for combining matchers in biometric identification systems only utilize the scores related to one class to produce the final combination score. Combination algorithms of low complexity type discard the dependency information between scores assigned to all classes by any single classifier. Instead of using low complexity combination algorithms in identification systems, we describe the use of medium II complexity type combinations, which utilize all the available scores and require the training of only a single combination function.

In order to use the relationships between scores assigned by one classifier to different classes, we have introduced the concept of the second-best score statistic. It is a way of score normalization where the normalization depends on all the scores output by a classifier in any one identification trial,

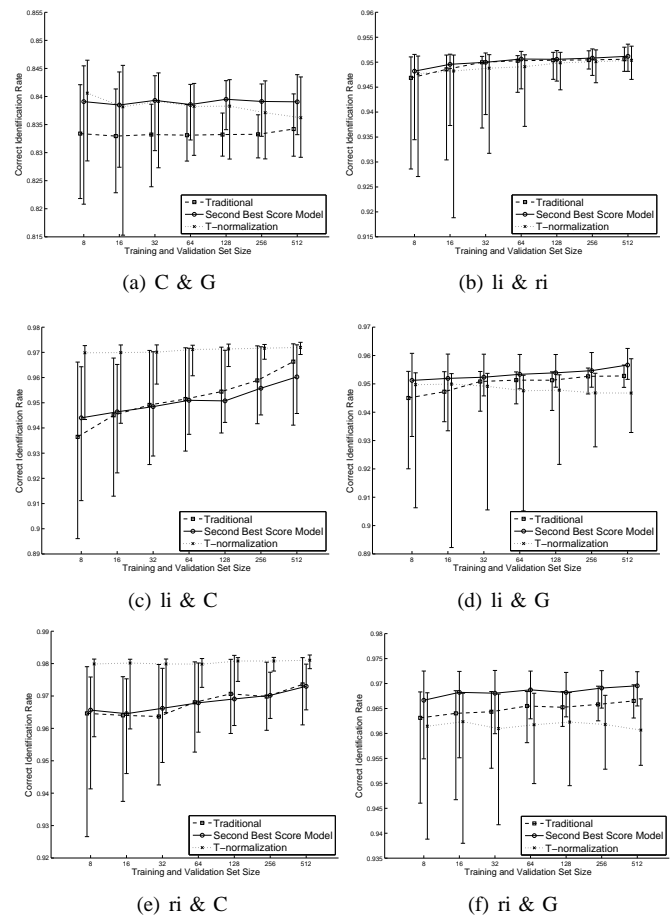


Fig. 10. Correct identification rates together with 90% bootstrap confidence estimates for different numbers of training and validation samples in neural network combination methods.

and the method is the same for all classes. This approach has less complexity than previous attempts of normalization [34], [35]. In these previous attempts normalizations were class specific and required huge amount of training data. The combinations utilizing such normalizations are similar to Behavior Knowledge Space combination [36], and belong to the high complexity combination type. Biometric identification problems can have a large number of enrolled persons, and such combinations are not feasible due to the lack of training data. By restricting ourselves to non-class-specific normalizations we are able to concentrate on combinations of medium II complexity type. Such combinations have significantly lower complexity, and result in efficient algorithms for identification systems.

## REFERENCES

- [1] L. Xu, A. Krzyzak, and C. Y. Suen, "Methods for combining multiple classifiers and their applications to handwriting recognition," *IEEE transactions on System, Man, and Cybernetics*, vol. 23, no. 3, pp. 418–435, 1992.
- [2] S. Tulyakov and V. Govindaraju, "Classifier combination types for biometric applications," in *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2006), Workshop on Biometrics*, New York, USA, 2006.
- [3] —, "Identification model for classifier combinations," in *Biometrics Consortium Conference*, Baltimore, MD, 2006.

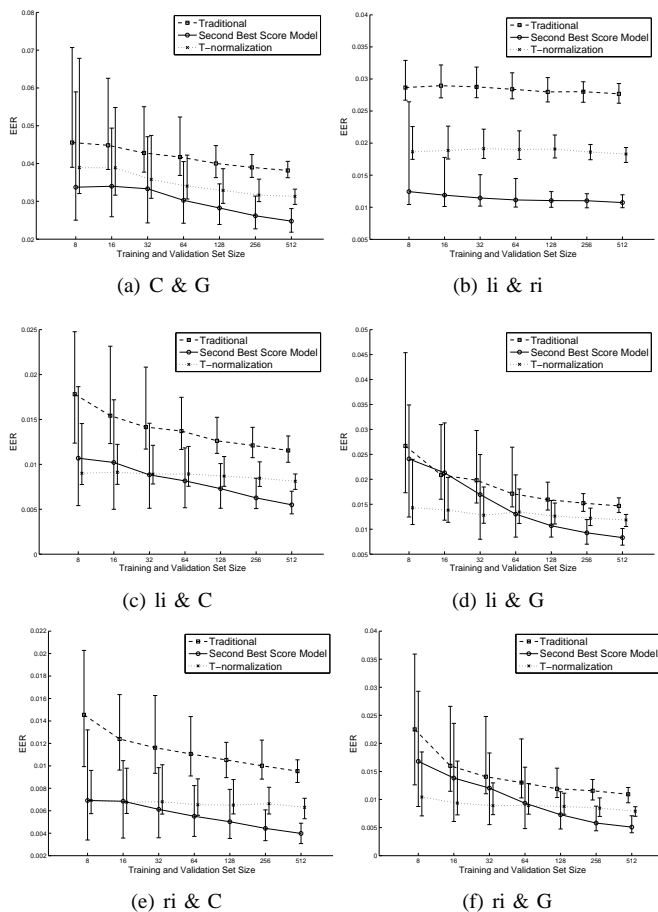


Fig. 11. Equal error rates (EER) together with 90% bootstrap confidence estimates for different numbers of training and validation samples in likelihood ratio combination methods.

[4] N. Fan, J. Rosca, and R. Balan, "Speaker verification with combined threshold, identification front-end, and ubm," in *Automatic Identification Advanced Technologies, 2005. Fourth IEEE Workshop on*, 2005, pp. 112–117.

[5] R. M. Bolle, J. H. Connell, S. Pankanti, N. K. Ratha, and A. W. Senior, *Guide To Biometrics*. New York: Springer, 2004.

[6] Y. Lee, K. Lee, H. Jee, Y. Gil, W. Choi, D. Ahn, and S. Pan, "Fusion for multimodal biometric identification," in *Audio- and Video-Based Biometric Person Authentication, 2005*, pp. 1071–1079.

[7] A. Jain and A. Ross, "Learning user-specific parameters in a multibiometric system," in *Image Processing, 2002. Proceedings. 2002 International Conference on*, vol. 1, 2002, pp. I-57–I-60 vol.1.

[8] J. Fierrez-Aguilar, D. Garcia-Romero, J. Ortega-Garcia, and J. Gonzalez-Rodriguez, "Bayesian adaptation for user-dependent multimodal biometric authentication," *Pattern Recognition*, vol. 38, no. 8, pp. 1317–1319, 2005.

[9] T. K. Ho, "A theory of multiple classifier systems and its application to visual word recognition," Ph.D Thesis, SUNY Buffalo, 1992.

[10] T. K. Ho, J. Hull, and S. Srihari, "Decision combination in multiple classifier systems," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 16, no. 1, pp. 66–75, 1994.

[11] R. Brunelli and D. Falavigna, "Person identification using multiple cues," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 17, no. 10, pp. 955–966, 1995.

[12] L. Hong and A. Jain, "Integrating faces and fingerprints for personal identification," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 20, no. 12, pp. 1295–1307, 1998.

[13] Y. Huang and C. Suen, "A Method of Combining Multiple Experts for Recognition of Unconstrained Handwritten Numerals," *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 17, no. 1, pp. 90–94, 1995.

[14] A. Saranli and M. Demirekler, "A statistical unified framework for rank-

based multiple classifier decision combination," *Pattern Recognition*, vol. 34, no. 4, pp. 865–884, 2001.

[15] A. Jain, K. Nandakumar, and A. Ross, "Score normalization in multimodal biometric systems," *Pattern Recognition*, vol. 38, no. 12, pp. 2270–2285, 2005.

[16] J. Kittler, M. Hatef, R. Duin, and J. Matas, "On combining classifiers," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 20, no. 3, pp. 226–239, 1998.

[17] R. Auckenthaler, M. Carey, and H. Lloyd-Thomas, "Score normalization for text-independent speaker verification systems," *Digital Signal Processing*, vol. 10, no. 1-3, pp. 42–54, 2000.

[18] J. Mariethoz and S. Bengio, "A unified framework for score normalization techniques applied to text independent speaker verification," *IEEE Signal Processing Letters*, vol. 12, 2005.

[19] P. Grother, "Face recognition vendor test 2002 supplemental report," NIST, Tech. Rep. NISTIR 7083, 2004.

[20] H. Altincay and M. Demirekler, "Undesirable effects of output normalization in multiple classifier systems," *Pattern Recognition Letters*, vol. 24, no. 9-10, pp. 1163–1170, 2003.

[21] A. Rosenberg and S. Parthasarathy, "Speaker background models for connected digit password speaker verification," in *Acoustics, Speech, and Signal Processing, 1996. ICASSP-96. Conference Proceedings., 1996 IEEE International Conference on*, vol. 1, 1996, pp. 81–84 vol. 1.

[22] J. Colombi, J. Reider, and J. Campbell, "Allowing good impostors to test," in *Signals, Systems & Computers, 1997. Conference Record of the Thirty-First Asilomar Conference on*, vol. 1, 1997, pp. 296–300 vol.1.

[23] D.-S. Lee, *Theory of Classifier Combination: The Neural Network Approach*. SUNY at Buffalo: Ph.D Thesis, 1995.

[24] L. I. Kuncheva, J. C. Bezdek, and R. P. W. Duin, "Decision templates for multiple classifier fusion: an experimental comparison," *Pattern Recognition*, vol. 34, no. 2, pp. 299–314, 2001.

[25] L. I. Kuncheva, *Combining Pattern Classifiers: Methods and Algorithms*. Wiley InterScience, 2004.

[26] S. Tulyakov, "A complexity framework for combination of classifiers in verification and identification systems," Ph.D. dissertation, State University of New York at Buffalo, 2006.

[27] V. Vapnik, *Statistical Learning Theory*. New York: Wiley, 1998.

[28] S. Tulyakov, V. Govindaraju, and C. Wu, "Optimal classifier combination rules for verification and identification systems," in *7th International Workshop on Multiple Classifier Systems*, Prague, Czech Republic, 2007.

[29] S. Tulyakov and V. Govindaraju, "Combining matching scores in identification model," in *8th International Conference on Document Analysis and Recognition (ICDAR 2005)*, Seoul, Korea, 2005.

[30] "Nist biometric scores set. <http://www.nist.gov/biometricscores/>."

[31] B.-H. Juang and S. Katagiri, "Discriminative learning for minimum error classification [pattern recognition]," *Signal Processing, IEEE Transactions on [see also Acoustics, Speech, and Signal Processing, IEEE Transactions on]*, vol. 40, no. 12, pp. 3043–3054, 1992.

[32] D. Miller, A. Rao, K. Rose, and A. Gersho, "A global optimization technique for statistical classifier design," *Signal Processing, IEEE Transactions on [see also Acoustics, Speech, and Signal Processing, IEEE Transactions on]*, vol. 44, no. 12, pp. 3108–3122, 1996.

[33] A. Schlappbach and H. Bunke, "Using hmm based recognizers for writer identification and verification," in *9th Intl Workshop on Frontiers in Handwriting Recognition (IWFHR-9 2004)*, 2004.

[34] D. Bouchaffra, V. Govindaraju, and S. Srihari, "A methodology for mapping scores to probabilities," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 21, no. 9, September 1999.

[35] K. Ianakiev, *Organizing Multiple Experts for Efficient Pattern Recognition*. SUNY at Buffalo: Ph.D Thesis, 2000.

[36] L. Lam and C. Y. Suen, "Optimal combinations of pattern classifiers," *Pattern Recognition Letters*, vol. 16, no. 9, pp. 945–954, 1995.