



# Combining flat and structured representations for fingerprint classification with recursive neural networks and support vector machines

Yuan Yao<sup>a</sup>, Gian Luca Marcialis<sup>b</sup>, Massimiliano Pontil<sup>a,c,\*</sup>, Paolo Frasconi<sup>d</sup>,  
Fabio Roli<sup>b</sup>

<sup>a</sup>*Department of Mathematics, City University of Hong Kong, 83 Tat Chee Ave, Kowloon, Hong Kong*

<sup>b</sup>*DIEE, University of Cagliari, Italy*

<sup>c</sup>*DII, University of Siena, Italy*

<sup>d</sup>*DSI, University of Florence, Italy*

Received 21 December 2001

## Abstract

We present new fingerprint classification algorithms based on two machine learning approaches: support vector machines (SVMs) and recursive neural networks (RNNs). RNNs are trained on a structured representation of the fingerprint image. They are also used to extract a set of distributed features of the fingerprint which can be integrated in the SVM. SVMs are combined with a new error-correcting code scheme. This approach has two main advantages: (a) It can tolerate the presence of ambiguous fingerprint images in the training set and (b) it can effectively identify the most difficult fingerprint images in the test set. By rejecting these images the accuracy of the system improves significantly. We report experiments on the fingerprint database NIST-4. Our best classification accuracy is of 95.6 percent at 20 percent rejection rate and is obtained by training SVMs on both FingerCode and RNN-extracted features. This result indicates the benefit of integrating global and structured representations and suggests that SVMs are a promising approach for fingerprint classification.

© 2002 Published by Elsevier Science Ltd on behalf of Pattern Recognition Society.

**Keywords:** Fingerprint classification; Error-correcting output codes; Recursive neural networks; Support vector machines

## 1. Introduction

The pattern-recognition problem studied in this paper consists of classifying fingerprint images into one out of five categories: whorl (W), right loop (R), left loop (L), arch (A), and tented arch (T). These categories were defined during early investigations of fingerprint structure [1] and have been used extensively since then. The task is important because classification can be employed as a preliminary step for reducing the complexity of database search in the

problem of automatic fingerprint matching [2,3]: If a query image can be classified with high accuracy, the subsequent matching algorithm only needs to compare stored images belonging to the same class.

Many approaches to fingerprint classification have been presented in the literature and this research topic is still very active. Overviews of the literature can be found in Refs. [4–7]. These approaches can be coarsely divided into two main categories: “flat” and “structural” approaches. Flat approaches are characterized by the use of “decision-theoretic” (or statistical) techniques for pattern classification, namely, a set of characteristic measurements, called feature vector, is extracted from fingerprint images and used for classification. Methods in this category include detection of singular points [8], connectionist algorithms such as self-organizing feature

\* Corresponding author. DII, University of Siena, Italy. Fax: +39-0577-233602.

E-mail addresses: mayyao@cityu.edu.hk (Y. Yao), pontil@dii.unisi.it (M. Pontil).

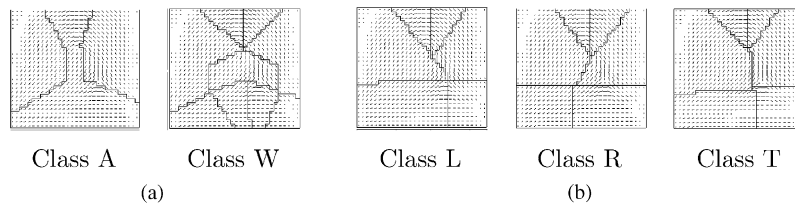


Fig. 1. (a) Examples of segmented fingerprint images related to the A and W classes. It is easy to see that structural information is very useful for distinguishing such fingerprint classes. (b) Examples of segmented fingerprint images related to the L, R, and T classes. It is easy to see that structural information is not very useful for distinguishing such fingerprint classes.

maps [9], neural networks [10,11], and hidden Markov models [12]. Structural approaches presented in the literature rely on syntactic or structural pattern-recognition techniques. In this case, fingerprints are described by production rules [13] or relational graphs. Parsing or (dynamic) graph matching algorithms [4,14] are used for classification. The structural approach has received less attention until now. However, a simple visual analysis of the “structure” of fingerprint images allows one to observe that structural information can be very useful for distinguishing some fingerprint classes (e.g., for distinguishing fingerprints belonging to class A from the ones of class W—see Fig. 1). Accordingly, we will attempt to “integrate” flat and structured representations and evaluate the practical benefits of their combination. Concerning this issue, very few works investigated the potentialities of such an integration [15,16].

In this paper, we propose new fingerprint classification algorithms based on two machine learning approaches: support vector machines (SVMs) and recursive neural networks (RNNs). SVM is a relatively new technique for pattern classification and regression that is well founded in statistical learning theory [17]. One of the main attractions of using SVMs is that they are capable of learning in *sparse, high-dimensional spaces* with very few training examples. They have been successfully applied to various classification problems (see Ref. [18] and references therein<sup>1</sup>). An RNN is a connectionist architecture designed for solving the supervised learning problem when the instance space is comprised of labeled graphs [19]. This architecture can exploit structural information in the data, which, as explained above, may help in discriminating between certain classes. In this paper, RNNs are also used to extract a distributed vectorial representation of the relational graph associated with a fingerprint. This vector is regarded as an additional set of features subsequently used as inputs for the SVM classifier.

An important issue in fingerprint classification is the problem of ambiguous examples: some fingerprints are assigned to two classes simultaneously, i.e. they have double labels (these images are also called “cross-referenced”). In order to address this issue, we designed an error-correcting code

(ECC) [20] scheme of SVM classifiers based on a new type of decoding distance. This method presents two main advantages. First, it allows a more accurate use of ambiguous examples because each SVM is in charge of generating only one codebit, whose value discriminates between two disjoint sets of classes. Then, if a fingerprint has labels all belonging to the same set for a particular codebit, we can retain this example in the training set without introducing any labeling noise. The second advantage of our system is its capability to deal with rejection problems. This is due to the concept of *margin* inherent to the SVM, which is incorporated in the decoding distance.

The system is validated on the NIST database 4 [21]. We present three series of experiments. As a base fingerprint representation we use FingerCode features,<sup>2</sup> a flat representation scheme proposed in Ref. [3]. In the first set of experiments, FingerCode features are combined with a structural representation of fingerprints based on relational graphs. In this case, a connectionist architecture that integrates flat and structural representations achieves 87.9 percent accuracy at 1.8 percent rejection rate. In the second experiment, SVMs are trained on FingerCode [3] preprocessed images, achieving 89.1 percent accuracy with 1.8 percent rejection. This result is only 0.9 percent worse than the accuracy obtained in Ref. [3] using the same features and a two stages  $k$ -NN/MLP classifier. Interestingly the SVM’s accuracy is much better than separate accuracies of both  $k$ -NN and MLP. Finally, the SVM is trained on both FingerCode and RNN-extracted features. In so doing, the performance is improved to 90.0 percent at 1.8 percent rejection rate. By allowing 20 percent rejection rate, accuracy increases to 95.6 percent. This result indicates the benefit of integrating global and structural representations for fingerprint classification.

The paper is organized as follows: In Section 2, we briefly describe SVM’s theory and introduce our new ECC classification method. Section 3 discusses the method we designed for the extraction of the structural features. The RNN is presented in Section 4. In Section 5 we report the experimental results. Finally, we report our conclusions in Section 6.

<sup>1</sup> For an updated list of SVM applications see [www.clopinet.com/isabelle/Projects/SVM/appllist.html](http://www.clopinet.com/isabelle/Projects/SVM/appllist.html).

<sup>2</sup> A short introduction on FingerCode feature is given in Section 5.1.

## 2. Support vector machines

In this section, we briefly overview the main concepts of SVMs [17] for pattern classification. More detailed accounts are given in Refs. [17,18,22,23].

### 2.1. Binary classification

SVMs perform pattern recognition for two-class problems by determining the separating hyperplane<sup>3</sup> with maximum distance to the closest points of the training set. These points are called *support vectors*. If the data are not linearly separable in the input space, a non-linear transformation  $\Phi(\cdot)$  can be applied which maps the data points  $\mathbf{x} \in \mathbb{R}^n$  into a high (possibly infinite) dimensional space  $\mathcal{H}$  which is called feature space. The data in the feature space are then separated by the optimal hyperplane as described above.

The mapping  $\Phi(\cdot)$  is represented in the SVM classifier by a kernel function  $K(\cdot, \cdot)$  which defines an inner product in  $\mathcal{H}$ , i.e.  $K(\mathbf{x}, \mathbf{t}) = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{t})$ . The decision function of the SVM has the form

$$f(\mathbf{x}) = \sum_{i=1}^{\ell} \alpha_i c_i K(\mathbf{x}_i, \mathbf{x}), \quad (1)$$

where  $\ell$  is the number of data points and  $c_i \in \{-1, 1\}$  is the class label of training point  $\mathbf{x}_i$ . Coefficients  $\alpha_i$  in Eq. (1) can be found by solving a quadratic programming problem with linear constraints. The support vectors are the nearest points to the separating boundary and are the only ones for which  $\alpha_i$  in Eq. (1) can be non-zero.

An important family of admissible kernel functions are the Gaussian kernel

$$K(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2 / 2\sigma^2),$$

with  $\sigma$  the variance of the Gaussian, and the polynomial kernels

$$K(\mathbf{x}, \mathbf{y}) = (1 + \mathbf{x} \cdot \mathbf{y})^d,$$

with  $d$  the degree of the polynomial. For other important examples of kernel functions used in practice see Refs. [17,23].

Let  $\rho$  be the distance of the support vectors to the hyperplane. This quantity is called *margin* and it is related to the coefficients in Eq. (1),

$$\rho = \left( \sum_{i=1}^{\ell} \alpha_i \right)^{1/2}. \quad (2)$$

The margin is an indicator of the separability of the data. More precisely, the expected error probability of the SVM is bounded by the average (with respect to the training set)

<sup>3</sup> SVM theory also includes the case of non-separable data, see Ref. [17].

of  $R^2 / \ell \rho^2$ , where  $R$  is the radius of the smallest sphere containing the support vector in the feature space [17]. In the case where the Gaussian kernel is used,  $R$  can be upper bounded by 1. Then, in this case the effect of the radius on the expected error probability of the SVM can be discarded.

### 2.2. Multiclass classification with ECCs

Many real-world classification problems involve more than two classes. Attempts to solve  $q$ -class problems with SVM have involved training  $q$  SVMs, each of which separates a single class from all remaining classes [17], or training  $q(q-1)/2$  machines, each of which separates a pair of classes [24–26]. The first type of classifiers are usually called *one-vs.-all*, while classifiers of the second type are called *pairwise* classifiers. When the one-vs.-all classifiers are used, a test point is classified into the class whose associated classifier has the highest score among all classifiers. In the case of pairwise classifiers, a test point is classified in the class which gets most votes among all the possible classifiers [25].

Classification schemes based on training one-vs.-all and pairwise classifiers have two extreme approaches: the first uses all the data, the second the smallest portion of the data. In practice, it can be more effective to use intermediate classification strategies in the style of ECCs [20,27]. In this case, each classifier is trained to separate a subset of classes from another disjoint subset of classes (the union of these two subsets does not need to cover all the classes). For example, the first set could consist of classes A and T and the second of classes R, L, and W. By doing so, we associate each class with a row of the “coding matrix”  $M \in \{-1, 0, 1\}^{q \times s}$ , where  $s$  denotes the number of classifiers.  $M_{ij} = -1$  or  $1$  indicates that points in class  $i$  are regarded as negative or positive examples for training the  $j$ th classifier.  $M_{ij} = 0$  indicates that points in class  $i$  are not used for training the  $j$ th classifier.

Three sets of SVM classifiers were used to construct the coding matrix: 5 one-vs.-all classifiers, 10 two-vs.-three classifiers, and 10 pairwise classifiers. This allows a more accurate use of ambiguous examples, since each SVM is only in charge of generating one codebit, whose value discriminates between two disjoint sets of classes. If a fingerprint has labels all belonging to the same set for a particular codebit, then clearly we can retain this example in the training set without introducing any labeling noise. As an example, consider fingerprints with double labels A and T. These examples are discarded by the one-vs.-all classifiers of classes A and T, and the pairwise classifier A-vs.-T. However, they are used to train the classifier AT-vs.-RLW.

A test point is classified in the class whose row in the coding matrix has minimum distance to the vector of outputs of the classifiers. Let  $\mathbf{m}$  be a row of the coding matrix,  $\mathbf{f}$  the vector of outputs of the classifiers, and  $\gamma_i$  the margin of the  $i$ th classifier. The simplest and most commonly used distance is the Hamming distance  $d(\mathbf{f}, \mathbf{m}) = \sum_{i=1}^s 1 - \text{sign}(f_i m_i)$ . We will also use two other distance measures which take into

account the margin of the classifiers: The margin-weighted Euclidean distance,  $d(\mathbf{m}, \mathbf{f}) = \sum_{i=1}^s (\gamma_i f_i - m_i)^2$ , and the soft margin distance proposed in Ref. [27],  $d(\mathbf{f}, \mathbf{m}) = \sum_{i=1}^s |1 - f_i m_i|_+$ . In the latter expression, the function  $|x|_+$  is equal to  $x$ , when  $x > 0$ , and zero otherwise.

In Section 5, we will show that our ECC of SVM achieves a performance of 89.1 percent when fingerprints are represented with FingerCode features [3] and a performance of 90.0 percent when we use both FingerCode and the structural features below. At the same time we found that, when only the first label is used as a target (like standard systems do), the performance drops down to 88.4 and 89.4 percent, respectively. This result indicates the advantage of our ECC-based approach to better handling of cross-referenced fingerprints during training.

### 3. Extraction of structural features

In this section, we discuss the algorithm we have used to extract the structured representation of the fingerprint image.

#### 3.1. Flow diagram of the structural classification module

As pointed out in the Introduction, a visual analysis of fingerprint images allows one to note that a fingerprint's "structure" can be extracted by segmenting the fingerprint into regions characterized by homogeneous ridge directions. To explain this aspect, consider the sample-segmented directional images in Fig. 1 [4]. These images clearly tell us that structural information can be very useful for distinguishing fingerprint classes A and W (see Fig. 1a). On the other hand, structural information is not effective for distinguishing fingerprints belonging to classes L, R, and T (see Fig. 1b). The fingerprint structure can be characterized by local attributes of the image regions (area, average directional value, etc.) and by the geometrical and spectral relations among adjacent regions (relative positions, differences among directional average values, etc.). Therefore, the developed module extracts and represents the structural information of fingerprints by segmenting the related directional images and by converting such segmented images into relational graphs whose nodes correspond to regions extracted by the segmentation algorithm. Graph nodes are then characterized by local characteristics of regions and by the geometrical and spectral relations among adjacent regions.

The main steps of our structural approach to fingerprint classification are as follows:

- Computation of the directional image of the fingerprint. This directional image is a  $28 \times 30$  matrix. Each matrix element represents the ridge orientation within a given block of the input image. The directional image was computed using the algorithm proposed in Ref. [28].
- Segmentation of the directional image into regions containing ridges with similar orientations. To this end, the

segmentation algorithm described in Ref. [4] was used.

- Representation of the segmented image by a directed positional<sup>4</sup> acyclic graph (DPAG). The representation of fingerprint structure is completed by characterizing each graph node with a numerical feature vector containing local characteristics of the related image region (area, average directional value, etc.) and some geometrical and spectral relations with respect to adjacent regions (relative positions, differences among directional average values, etc.).
- Classification of the above DPAG by an RNN made up of two multilayer perceptrons (MLPs) neural nets. This neural network model is briefly described below (see Ref. [19] for more details).

#### 3.2. Image segmentation and relational graph construction

In order to segment directional fingerprint images, we used an algorithm explicitly designed for such a task [28]. This algorithm implements a very sophisticated "region growing" process. It starts from the central element of the directional image and scans the image according to a square spiral strategy. At each step, the segmentation algorithm uses a cost function to decide about the creation of a new region. The resulting segmentation is related to a minimum of such a cost function. Details about this segmentation algorithm can be found in Ref. [4]. The construction of relational graphs from segmented images is performed by the following main steps:

- The image region containing the "core" point is selected as the starting region for the graph construction, that is, a graph node associated to such a region is initially created;
- The regions that are adjacent to such a core region are then evaluated for the creation of new nodes. Nodes are created for adjacent regions which are located in one of the following spatial positions with respect to the core region: North, North East, East, South East, South, South West, West, and North West;
- The above process is repeated for each of the new nodes until all the regions of the segmented images have been considered;
- The graph nodes created by the above algorithm are finally characterized by a numerical feature vector containing local characteristics of the related image regions (area, average directional value, etc.) and some geometrical and spectral relations with respect to adjacent regions (relative positions, differences among directional average values, etc.).

<sup>4</sup> Positional here means that for each vertex  $v$ , a bijection  $P : \mathcal{E} \rightarrow \mathbb{N}$  is defined on the edges leaving from  $v$ .

---

**Algorithm 1.** rgExtract( $R$ ) Construct relational graphs from segmented directional images.

---

**Input:**  $R = [R_1, R_2, \dots, R_n]$  is the set of regions of the segmented image.

**Output:**  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , the relational graph with vertex set  $\mathcal{V}$  and edge set  $\mathcal{E}$ .

**Notation:** For  $k = 0, \dots, 7$  we shall denote the “reference angles” for each position by  $\beta_k = k\pi/4$ , and the “reference interval” by  $I_k = [\beta_k - \pi/8, \beta_k + \pi/8]$ . An edge in a DPAG is denoted  $(u, v, k)$  indicating that  $v$  is the  $k$ th child of  $u$ . Finally, for each region  $R$ ,  $Xbar[R]$  and  $Ybar[R]$  denote the x-y coordinates of the region’s center of mass.

$\mathcal{V} \leftarrow R$

$\mathcal{E} \leftarrow \emptyset$

**for**  $j \leftarrow 1, \dots, n$  **do**

    Color[ $R_j$ ]  $\leftarrow$  WHITE;

**end for**

Let  $Q$  be a FIFO queue of regions (initially empty);

Let  $R_i$  be the region containing the core of the fingerprint;

Enqueue( $Q, R_i$ );

**while**  $Q$  is not empty **do**

$R_c \leftarrow$  Dequeue( $Q$ );

    Color[ $R_c$ ]  $\leftarrow$  BLACK;

$R'_c \leftarrow \{R_j : R_j \text{ is adjacent to } R_c\}$

**for all**  $R_j \in R'_c$  **do**

**if** Color[ $R_j$ ] = WHITE **then**

            Enqueue( $Q, R_j$ );

$\alpha_j \leftarrow$  slope of the vector from  $(Xbar[R_c], Ybar[R_c])$  to  $(Xbar[R_j], Ybar[R_j])$ ;

            Note that  $\alpha_j \in [0, 2\pi]$ .

**for**  $k \leftarrow 0, \dots, 7$  **do**

**if**  $\alpha_j \in I_k$  **then**

$d_{jk} \leftarrow |\alpha_j - \beta_k|$ ;

**else**

$d_{jk} \leftarrow \infty$ ;

**end if**

**end for**

**end if**

**end for**

**for**  $k \leftarrow 0, \dots, 7$  **do**

$h \leftarrow \operatorname{argmin}_j \{d_{jk}\}$ ;

$\mathcal{E} \leftarrow \mathcal{E} \cup \{(R_c, R_h, k)\}$ ;

**end for**

**end while**

**return**  $\mathcal{G}$

---

Algorithm 1 is used to derive relational graphs from segmented directional images. Fig. 2 depicts a relational graph that the algorithm derived from a segmented directional image.

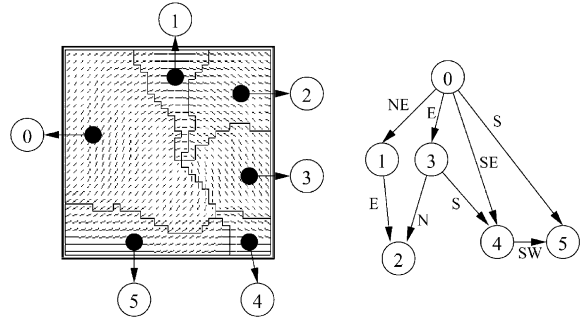


Fig. 2. Left: segmented fingerprint image and right: relational graph.

#### 4. Recursive neural networks

Research in connectionist models capable of representing and learning structured (or hierarchically organized) information began in the early 1990s with recursive auto-associative memories (RAAMs) [29]. Since then, several other architectures have been proposed, including holographic reduced representations (HRRs) [30], and RNNs [19,31,32]. A selection of papers in this research area recently appeared in Ref. [33]. RNNs are capable of solving the supervised learning problems such as classification and regression when output prediction is conditioned on a hierarchical data structure, like the structural representation of fingerprints described above. The input to the network is a labeled DPAG  $U$ , where the label  $U(v)$  at each vertex  $v$  is a real-value feature vector associated with a fingerprint region, as described in Section 3.2. A hidden state vector  $X(v) \in R^n$  is associated with each node  $v$ . This vector contains a distributed representation of the subgraph dominated by  $v$  (i.e. all the vertices that can be reached starting a directed path from  $v$ ). The state vector is computed by a state transition function which combines the state vectors of  $v$ 's children with a vector encoding of the label of  $v$ :

$$X(v) = f(X(w_1), \dots, X(w_k), U(v)), \quad (3)$$

$\{w_1, \dots, w_k\}$  being the ordered set of  $v$ 's children. Computation proceeds recursively from the frontier to the supersource (the vertex dominating all other vertices). The base step for Eq. (3) is  $X(w) = 0$  if  $w$  is a missing child. Transition function  $f$  is computed by a multilayer perceptron (MLP), which is replicated at each node in the DPAG, sharing weights among replicas. Note that the state vector is similar to the context layer used in simple recurrent networks [34], except that it encodes its context as a recursive data structure, and not simply a sequence. Pollack's RAAM [29] used similar distributed compositional representations but can only encode and decode trees and cannot solve supervised learning problems.

In the case of classification we assume that the input graph possesses a supersource  $s$  (i.e. a node from which every other node can be reached by a directed path). Classification with recurrent neural networks is then performed by adding an

output function  $g$  that takes as input the hidden state vector  $X(s)$  associated with the supersource  $s$ :

$$Y = g(X(s)). \quad (4)$$

Function  $g$  is also implemented by an MLP. The output layer in this case uses the *softmax* function (normalized exponentials), so that  $Y$  can be interpreted as a vector of conditional probabilities of classes given the input graph, i.e.  $Y_i = P(C = i|U)$ ,  $C$  being a multinomial class variable. Training relies on maximum likelihood. The training set consists of  $T$  pairs  $\mathcal{D} = \{(U_1, c_1), \dots, (U_t, c_t), \dots, (U_T, c_T)\}$  where  $c_t$  denotes the class of the  $t$ th fingerprint in the dataset. According to the multinomial model, the log-likelihood has the form

$$l(\mathcal{D}; \theta) = \sum_t \log Y_{c_t}, \quad (5)$$

where  $\theta$  denotes the set of trainable weights and  $t$  ranges over training examples. Optimization is performed with a gradient descent procedure, where gradients are computed by the back-propagation through structure algorithm [31,32].

We remark that the state vector  $X(s)$  at the supersource is a distributed representation of the entire input DPAG and encodes features of the input DPAG deemed to be relevant for discrimination amongst classes. In the subsequent experiments, the components of the state vector at the supersource are thus used as additional features that may help to discriminate between some class pairs.

## 5. Experimental results

### 5.1. Dataset and FingerCode features

Our system was validated on the NIST Database 4 [21]. This database consists of 4000 images analyzed by a human expert and labeled with one *or more* of the five structural classes W, R, L, A, and T (more than one class is assigned in cases where ambiguity could not be resolved by the human expert). Previous works on the same dataset either rejected ambiguous examples in the training set (loosing in this way part of the training data), or used the first label as a target (potentially introducing output noise).

Fingerprints were represented with the structured representation discussed in Section 3 as well as with FingerCode features. FingerCode is a representation scheme described in Ref. [3] and consists of a vector of 192 real features computed in three steps. First, the fingerprint core and center are located. Then the algorithm separates the number of ridges present in four directions ( $0^\circ$ ,  $45^\circ$ ,  $90^\circ$ , and  $135^\circ$ ) by filtering the central part of a fingerprint with a bank of Gabor filters. Finally, standard deviations of grayscale values are computed on 48 disk sectors, for each of the four directions.

As in Ref. [3], a few fingerprint images were rejected<sup>5</sup> both during training (1.4 percent) and testing (1.8 percent).

<sup>5</sup> This is due to the failure of FingerCode in reliably locating the fingerprint core.

Thus, when not explicitly said, the evaluation performance of the methods discussed below is intended to be at 1.8 percent rejection rate. For the same reason, we were not able to compute the performance at zero percent rejection rate, which is not expected to increase. However, this does not reduce the interest of the discussion. Our main goal here is to show the novel methodology offered by our machine learning approach: (a) the advantage of ECC to better exploit cross-referenced (i.e. double labeled) images for training, (b) the strength of the SVM in improving accuracy versus rejection curves, and (c) the benefit of integrating flat and structural representations. In particular, as we explained in Section 2.2, our ECC scheme provides a novel approach to deal with cross-referenced images during training which can be easily implemented/extended to work on new image representations.

### 5.2. Results using RNN and integration of flat and structural classifiers

In order to investigate the potentialities of the combination of flat and structural methods, we coupled our structural approach (Section 3) with the vector-based approach proposed in Ref. [3]. Several strategies for combining classifiers were evaluated in order to combine the flat and structural approaches [35–37]. We firstly assessed the performance of simple and commonly used combination rules, namely, the majority voting rule and the linear combination. Such classifier combination rules work well if the classifiers exhibit similar accuracies and make “independent” errors. As the flat and the structural classifiers considered exhibit very different accuracies (86.0 percent vs. 71.5 percent; see Section 5.2.1) and their classifications were correlated in a complex way (namely, correlation of classifications strongly varies in the feature space), such simple combination rules performed poorly. Accordingly, we combined classifiers adopting the so-called “metaclassification” (or “stacked”) approach, which uses a trainable rule or an additional classifier for the combination. In particular, after various experiments with different trainable combination rules proposed in the literature [37], a  $k$ -nearest neighbor classifier was selected for the combination.

#### 5.2.1. Comparison between vector-based and structural classification

We have trained an MLP using the FingerCode feature vector as input. The best performance on the test set was obtained with an MLP architecture with 28 hidden units. Table 1a shows the corresponding confusion matrix. The overall accuracy is 86.0 percent.

Table 1b shows the confusion matrix of the RNN trained on the structural features discussed in Section 3. In this case the overall accuracy is 71.5 percent.



Table 1  
Confusion matrices

	W	R	L	A	T
(a) <i>MLP using Finger Code</i>					
W	354	24	13	4	1
R	8	349	0	6	27
L	8	2	349	6	15
A	0	8	4	347	72
T	55	8	12	2	292
(b) <i>RNN using relational graphs</i>					
W	323	28	36	5	1
R	24	305	16	6	64
L	44	11	266	7	57
A	3	13	12	333	68
T	21	44	49	41	179
(c) <i>Combination of MLP and RNN</i>					
W	359	19	10	3	2
R	6	345	0	7	29
L	4	2	342	5	18
A	0	4	0	361	65
T	0	8	9	46	312

Rows denote the true class, columns the assigned class.

Tables 1a and b point out that the accuracy of the structural classifier is much lower than the one of the flat classifier. This is mainly due to the large degree of confusion among L, R, and T classes. On the other hand, as expected, the best performance of the structural classifier is related to the discrimination between A and W classes.

Afterwards, we analyzed the degree of complementarity between the two classifiers discussed above. To this end, we computed the performance of an ideal “oracle” that, for each input fingerprint, always selects one of two classifiers, if any, that correctly classifies such a fingerprint. Such an oracle applied to the two classifiers provided an overall accuracy of 92.5 percent. This accuracy value obviously represents a very tight upper bound for any combination method applied to the two classifiers. Yet, it points out the potential benefit of the combination of the flat and structural classifiers.

#### 5.2.2. Combined flat and structural classification

A  $k$ -nearest neighbor classifier (with a value of the  $k$  parameter equal to 113) was used for combining the flat and structural classifiers. Such a metaclassifier takes the outputs of the two classifiers as inputs and provides the final fingerprint classification as output. Table 1c depicts the confusion matrix for this experiment. Note that such a combination outperforms the best single classifier (i.e. the MLP classifier using the FingerCode representation; see Table 1a), thus pointing out that the exploitation of structural information allows increasing classification performances. The accuracy of this classifier is 87.9 percent. In particular, we note that such a combination improves the performances related to A and W classes, confirming the intuition suggested by Fig. 1.

Table 2  
Confusion matrices

	W	R	L	A	T
(a) <i>One-vs.-all SVM</i>					
W	356	23	14	3	1
R	4	344	1	7	33
L	4	2	356	6	13
A	0	2	5	371	55
T	0	7	7	48	302
(b) <i>Pairwise SVM</i>					
W	359	24	15	3	1
R	4	341	1	6	30
L	5	0	356	6	15
A	0	2	4	363	58
T	0	7	9	38	317
(c) <i>ECC SVM with margin-weighted Euclidean decoding</i>					
W	362	22	11	3	0
R	4	350	3	8	27
L	7	2	357	5	11
A	0	3	3	398	32
T	0	13	9	51	283

Rows denote the true class, columns the assigned class.

#### 5.3. Results with SVM

We used the three types of multiclass classification schemes discussed in Section 2.1 which are based on the combination of binary SVMs. SVMs have been trained using the SVMFu code<sup>6</sup> on a 550 MHz Pentium-II PC. Training on 2000 examples takes about 10 s for pairwise classifiers and 20 s for one-vs.-all classifiers.

*One-vs.-all SVM.* We trained five one-vs.-all SVM classifiers using both Gaussian kernels and polynomials of degree between 2 and 6. The best result was obtained with the Gaussian kernel with  $\sigma = 1$ :88.0 percent at 1.8 percent rejection rate; the confusion matrix is reported in Table 2a. The best polynomial SVM was of degree 3 and achieved a performance of 84.5 percent only. Then, in the remaining experiments we used only the Gaussian kernel.

*Pairwise SVM.* We trained the ten pairwise SVMs. The test set accuracy increases to 88.4 percent at 1.8 percent rejection rate, improving the MLP accuracy reported in Ref. [3] of 2 percent. The confusion matrix is reported in Table 2b.

*Error-correction SVM scheme.* Three sets of SVM classifiers were used to construct the coding matrix: 5 one-vs.-all classifiers, 10 two-vs.-three classifiers, and 10 pairwise classifiers. The three kinds of decoding distances discussed in Section 2 were compared: (i) Hamming distance: 88.0 percent, (ii) Margin-weighted Euclidean distance: 89.1 percent, and (iii) Soft margin distance: 88.8 percent (all results are

<sup>6</sup> This software can be downloaded at <http://five-percent-nation.mit.edu/SvmFu>.

Table 3

Accuracy vs. rejection rate for the ECC of SVM trained on FingerCode features

Rejection rate (%):	1.8	8.5	20.0	32.5
Five classes (%):	89.1	90.6	93.9	96.2
Five classes (%):	93.7	95.4	97.1	98.4

at 1.8 percent rejection rate). These results confirm the advantage of taking into account the margin of the classifiers inside the decoding distance. The confusion matrix for the margin-weighted Euclidean distance is reported in Table 2c.

We have also trained the ECC of SVM for the four classes task (classes A and T merged together) using the margin-weighted Euclidean distance. The obtained accuracy is of 93.7 percent at 1.8 percent rejection rate. For comparison, the accuracy reported in Ref. [3] for the sole MLP is 92.1 percent, while the cascade of  $k$ -NN and MLP yields 94.8 percent.

### 5.3.1. Analysis of the margin

We have measured the margin as in Eq. (2) and the number of support vectors of each SVM classifier used in our experiments (the training error of each individual classifier was always equal to zero). The number of support vector ranges between 1/5 and 1/2 of the number of training points. As expected, the margin decreases for those classifiers which involve difficult pairs of classes. Among the pairwise classifiers, the A–T classifier has the smallest margin. The margin of the T-vs.-all and A-vs.-all is also small. However, the margin of the AT-vs.-RLW classifier increases, which might explain why our error-correcting strategy works well.

### 5.3.2. Rejection vs. accuracy

Let  $d_1$  be the minimum distance of the vector of outputs of the classifiers from the coding row, and  $d_2$  the second minimum distance. Rejection can be decided by looking at the difference  $\Delta = d_2 - d_1$ . This quantity can be seen as the margin of the multiclass classifier. A large value of  $\Delta$  indicates high confidence in classification; when  $\Delta$  is smaller than a given threshold we reject the data. The rejection rate is controlled by this threshold. Table 3 shows the accuracy–rejection tradeoff obtained in our experiments. Note that the accuracy of the system increases sharply with the rejection rate. At 20 and 32.5 percent rejection, the system shows a moderate improvement over the best results in Ref. [3].

Table 6

Summary of the results of the different proposed methods

Method:	RNN	MLP	$k$ -NN	SVM1	SVM2	SVM3	SVM and RNN
Performance (%):	71.5	86.0	87.9	88.0	88.4	89.1	90.0

We have used the following notation: SVM1 for one-vs.-all SVM, SVM2 for pairwise SVM, and SVM3 for the ECC of SVM.

Table 4

Confusion matrix for the ECC of SVM trained on both FingerCode and RNN-extracted features

	W	R	L	A	T
W	366	18	8	2	0
R	5	354	0	7	29
L	6	1	357	2	13
A	0	2	2	396	33
T	1	8	12	48	294

Table 5

Accuracy vs. rejection rate for the ECC of SVM trained on the union of FingerCode and RNN-extracted features

Rejection rate (%):	1.8	8.5	20.0	32.5
Five classes (%):	90.0	92.2	95.6	97.6
Four classes (%):	94.7	96.6	98.4	99.2

### 5.4. Using both vectorial and structured features with SVM

We have trained SVM on both FingerCode and RNN-extracted features and used the ECC scheme with margin-weighted Euclidean decoding. The confusion matrix is summarized in Table 4. The performance is improved to 90.0 percent at 1.8 percent rejection rate. If we compare this performance to the performance obtained with FingerCode features only (89.1 percent), we observe the benefit of integrating global and structural representations.

This effect is especially clear in the accuracy vs. rejection rate results. As shown in Table 5, the accuracy sharply increases with the rejection rate, improving significantly over the results obtained with FingerCode features only (see Table 3).

Finally, Table 6 summarizes the performance results of the different methods presented above. From there, one can clearly observe the advantage offered by the ECC approach as well as the effect of integrating flat and structural features.

## 6. Conclusions

In this paper we have studied the combination of flat and structured representations for fingerprint classification. An algorithm for extracting a structural representation of



fingerprint images was presented. RNNs were used to process this structural representation and to extract a distributed vectorial representation of the fingerprint. This vectorial representation was integrated with other global representation, showing significant improvement over global features only. Experiments were performed on the NIST Database 4. The best performance was obtained with an error-correcting code of SVM. This method can tolerate the presence of ambiguous examples in the training set and is shown to be precise to identify difficult test images, then sharply improving the accuracy of the system at a higher rejection rate.

A number of questions and future research directions are still open. The main one is how to train SVM directly on structured representation.

### Acknowledgements

We wish to thank Anil Jain for providing us with the dataset of preprocessed NIST-4 fingerprints. This work was partially supported by CERG Grant 9040457.

### References

- [1] E.R. Henry, *Classification and Uses of Finger Prints*, Routledge, London, 1900.
- [2] R.S. Germain, A. Califano, S. Colville, Fingerprint matching using transformation parameter clustering, *IEEE Comput. Sci. Eng.* 4 (4) (1997) 42–49.
- [3] A.K. Jain, S. Prabhakar, L. Hong, A multichannel approach to fingerprint classification, *Pattern Anal. Mach. Intell.* 21 (4) (1999) 348–359.
- [4] R. Cappelli, A. Lumini, D. Maio, D. Maltoni, Fingerprint classification by directional image partitioning, *Trans. Pattern Anal. Mach. Intell.* 21 (5) (1999) 402–421.
- [5] A.K. Jain, R. Bolle, S. Pankajanti, *Biometrics: Personal Identification in Networked Society*, Kluwer Academic Publishers, Norwell, MA, 1999.
- [6] L.C. Jain, U. Halici, I. Hayashi, S.B. Lee, Tsutsui (Eds.), *Intelligent Biometric Techniques in Fingerprint and Face Recognition*, CRC Press, Boca Raton, FL, USA, 1999.
- [7] A.K. Jain, S. Pankanti, *Fingerprint classification and recognition*, *The Image and Video Processing Handbook*, Academic Press, New York, 2000.
- [8] K. Karu, A.K. Jain, Fingerprint classification, *Pattern Recognition* 29 (3) (1996) 389–404.
- [9] U. Halici, G. Ongun, Fingerprint classification through self-organizing feature maps modified to treat uncertainty, *Proc. IEEE* 84 (10) (1996) 1497–1512.
- [10] K. Moscinska, G. Tyma, Neural network based fingerprint classification, in: *Proceedings of the Third International Conference on Neural Networks*, 1993, pp. 229–232.
- [11] H.V. Neto, D.L. Borges, Fingerprint classification with neural networks, in: *Proceedings of the Fourth Brazilian Symposium on Neural Networks*, IEEE Press, New York, 1997, pp. 66–72.
- [12] A. Senior, A hidden markov model fingerprint classifier, in: *Proceedings of the 31st Asilomar Conference on Signal, Systems, and Computers*, 1997, pp. 305–310.
- [13] B. Moayer, K.S. Fu, A syntactic approach to fingerprint pattern recognition, *Pattern Recognition* 7 (1975) 1–23.
- [14] D. Maio, D. Maltoni, A structural approach to fingerprint classification, in: *Proceedings of the 13th International Conference on Pattern Recognition*, Vol. 3, 1996, pp. 578–585.
- [15] R. Cappelli, D. Maio, D. Maltoni, Combining fingerprint classifiers, in: *Proceedings of the First International Workshop on Multiple Classifier Systems (MCS2000)*, Cagliari, 2000, pp. 351–361.
- [16] G.L. Marcialis, F. Roli, P. Frasconi, Fingerprint classification by combination of flat and structural approaches, *Proceedings of the 3rd International Conference on Audio- and Video-Based Biometric Person Authentication*, 2001.
- [17] V.N. Vapnik, *Statistical Learning Theory*, Wiley, New York, 1998.
- [18] N. Cristianini, J. Shawe-Taylor, *An Introduction to Support Vector Machines*, Cambridge University Press, Cambridge, 2000.
- [19] P. Frasconi, M. Gori, A. Sperduti, A general framework for adaptive processing of data structures, *IEEE Trans. Neural Networks* 9 (5) (1998) 768–786.
- [20] T.G. Dietterich, G. Bakiri, Solving multiclass learning problems via error-correcting output codes, *J. Artif. Intell. Res.* 1995.
- [21] C.I. Watson, C.L. Wilson, *Fingerprint Database*, National Institute of Standards and Technology, April 1992, Special Database 4, FPDB.
- [22] C. Burges, A tutorial on support vector machines for pattern recognition, in: *Data Mining and Knowledge Discovery*, Vol. 2, Kluwer Academic Publishers, Boston, 1998.
- [23] T. Evgeniou, M. Pontil, T. Poggio, Regularization networks and support vector machines, *Adv. Comput. Math.* 13 (2000) 1–50.
- [24] J. Platt, N. Cristianini, J. Shawe-Taylor, Large margin dags for multiclass classification, in: *Advances in Neural Information Processing Systems*, Denver, CO, 2000.
- [25] J.H. Friedman, Another approach to polychotomous classification, Technical Report, Department of Statistics, Stanford University, 1997.
- [26] M. Pontil, A. Verri, Support vector machines for 3-d object recognition, *IEEE Trans. Pattern Anal. Mach. Intell.* (1998) 637–646.
- [27] R.E. Schapire, Y. Singer, Erin Lee Young, Reducing multiclass to binary: a unifying approach for margin classifiers, Technical Report, AT&T Research, 2000.
- [28] M.J. Donahue, S.I. Rokhlin, On the use of level curves in image analysis, *Image Understanding* 57 (3) (1993) 185–203.
- [29] J.B. Pollack, Recursive distributed representations, *Artif. Intell.* 46 (1–2) (1990) 77–106.
- [30] T.A. Plate, Holographic reduced representations, *IEEE Trans. Neural Networks* 6 (3) (1995) 623–641.
- [31] C. Goller, A. Küchler, Learning task-dependent distributed structure-representations by backpropagation through structure, in: *IEEE International Conference on Neural Networks*, 1996, pp. 347–352.
- [32] A. Sperduti, A. Starita, Supervised neural networks for the classification of structures, *IEEE Trans. Neural Networks* 8 (3) (1997).

- [33] P. Frasconi, M. Gori, A. Sperduti, Special section on connectionist models for learning in structured domains, *IEEE Trans. Knowledge Data Eng.* 13 (2) (2001).
- [34] J.L. Elman, Finding structure in time, *Cognitive Sci.* 14 (1990) 179–211.
- [35] G. Giacinto, F. Roli, Ensembles of neural networks for soft classification of remote sensing images, in: *Proceedings of the European Symposium on Intelligent Techniques*, 1997, pp. 166–170.
- [36] G. Giacinto, F. Roli, L. Bruzzone, Combination of neural and statistical algorithms for supervised classification of remote-sensing images, *Pattern Recognition Lett.* 21 (5) (2000).
- [37] J. Kittler, F. Roli (Eds.), *Proceedings of the First International Workshop on Multiple Classifier Systems*, *Lecture Notes in Computer Science*, Vol. 1857, Springer, New York, 2000.

**About the Author**—YUAN YAO received his B.S. and M.S. in Control Engineering from Harbin Institute of Technology in 1996 and 1998, respectively. He is currently a Ph.D. student in Department of Mathematics, City University of Hong Kong and on leave at Department of Mathematics, University of California at Berkeley. From 1995 to 1999, he worked on Robust Control Theory and Cellular Neural Networks in Professor G.-X. Wang's group at HIT. From 1999, he has been working in Professor S. Smale's group, where his main research interests include mathematical foundation of learning, computational vision, and bioinformatics.

**About the Author**—GIAN LUCA MARCIALIS is a Ph.D. Student in Electronic and Computer Science Engineering at the Department of Electrical and Electronic Engineering of the University of Cagliari (Italy). He received the Laurea Degree (M.S.) in Electronic Engineering in June 2000 from the University of Cagliari. His research interests are in the fields of biometric applications, data complexity, design of multiple classifier systems. Gian Luca Marcialis is a Student Member of the International Association for Pattern Recognition (IAPR) and a Student Member of the Italian Association for Artificial Intelligence (AI\*IA).

**About the Author**—MASSIMILIANO PONTIL was born on 14 August 1970 in Genova, Italy. He received his Bachelor and Ph.D. in Theoretical Physics from the University of Genova in 1994 and 1999. Since October 2001 he has been with the Department of Information Engineering, University of Siena, Italy. His research interests involve machine learning, pattern recognition, and statistics. Previously, he has been a visiting student at MIT in 1998, a Post-doctoral Fellow in 1999 and 2000 at the Center for Biological and Computational Learning at MIT, and a Research Fellow at City University of Hong Kong in 2001. He has also spent some time as a visiting researcher at RIKEN Brain Science Institute, Tokyo, and AT&T Labs, USA. He has published about 30 papers in international journals and conferences on different aspects of learning theory, machine learning, and computer vision.

**About the Author**—PAOLO FRASCONI received his M.Sc. degree in Electronic Engineering in 1990, and his Ph.D. degree in Computer Science in 1994, both from the University of Florence, Italy. Since 2000 he is an Associate Professor of Computer Science with the Department of Systems and Computer Science (DSI) at the University of Florence. In 1999, he was an Associate Professor at the University of Cagliari, Italy. In 1998 he was a Visiting Lecturer with the School of Information Technology and Computer Science at the University of Wollongong, Australia. From 1995 to 1998 he was an Assistant Professor at the University of Florence. In 1992, he was a Visiting Scholar in the Department of Brain and Cognitive Science at Massachusetts Institute of Technology.

His current research interests are in the area of machine learning with connectionist models and belief networks, with particular emphasis on problems involving learning about sequential and structured information. Application fields of his interest include bioinformatics, natural language processing, pattern recognition, and document processing. Dr. Frasconi serves as an Associate Editor for the *IEEE Transactions on Neural Networks* and for the *IEEE Transactions on Knowledge and Data Engineering*. In 2001, he co-directed the NATO Advanced Studies Institute "Artificial Intelligence and Heuristic Methods for Bioinformatics". He is a member of the ACM, the IAPR, the IEEE, and the AI\*IA.

**About the Author**—FABIO ROLI obtained his M.S. degree and Ph.D. degree in Electronic Engineering from the University of Genoa, Italy, in 1988 and 1993, respectively. He was with the research group on Image Processing and Understanding of the Department of Biophysical and Electronic Engineering, University of Genoa, Italy, from 1988 to 1994. He was an adjunct professor at the University of Trento, Italy, in 1993 and 1994. Since 1995, he is with the Department of Electrical and Electronic Engineering of the University of Cagliari, Italy. He is full professor of computer engineering and leads the research activities of the department in the areas of pattern recognition and computer vision. His main area of expertise is the development of pattern recognition algorithms for applications like remote sensing, biometrics personal identification, video surveillance, and intrusion detection in computer networks. Prof. Roli's current research activity is focused on the theory and the applications of multiple classifier systems. He organized and co-chaired the two editions of the International Workshop on Multiple Classifier Systems ([www.diee.unica.it/mcs](http://www.diee.unica.it/mcs)). In his research fields, Prof. Roli has published more than 80 conference and journal papers, and he regularly acts as a reviewer for international journals.