

# High-Performance Memory-based Face Recognition for Visitor Identification

Terence Sim<sup>1</sup> Rahul Sukthankar<sup>1,2</sup> Matthew D. Mullin<sup>2</sup> Shumeet Baluja<sup>1,2</sup>

<sup>1</sup>The Robotics Institute  
Carnegie Mellon Univ.  
Pittsburgh, PA 15213

<sup>2</sup>Just Research  
4616 Henry Street  
Pittsburgh, PA 15213

{tsim, rahuls, mdm, baluja}@justresearch.com

## Abstract

*We show that a simple, memory-based technique for view-based face recognition, motivated by the real-world task of visitor identification, can outperform more sophisticated algorithms that use Principal Components Analysis (PCA) and neural networks. This technique is closely related to correlation templates; however, we show that the use of novel similarity measures greatly improves performance. We also show that augmenting the memory base with additional, synthetic face images results in further improvements in performance. Results of extensive empirical testing on two standard face recognition datasets are presented, and direct comparisons with published work show that our algorithm achieves comparable (or superior) results. This paper further demonstrates that our algorithm has desirable asymptotic computational and storage behavior, and is ideal for incremental training. Our system is incorporated into an automated visitor identification system that has been operating successfully in an outdoor environment for several months.*

## 1 Introduction

The problem of visitor identification consists of the following: a security camera monitors the front door of a building, acquiring images of people as they enter; an automated system extracts faces from these images and quickly identifies them using a database of known individuals. The system must easily adapt as people are added or removed from its database, and the system must be able to recognize individuals in near-frontal photographs. This paper focuses on the face recognition technology that is required to address this real-world task.

Face recognition has been actively studied [6, 11], particularly over the last few years [8]. The research effort has focused on the subproblem of frontal face recognition, with limited variance in illumination and facial expression. In this domain, techniques based on Principal Components Analysis (PCA) [9] popularly termed *eigenfaces* [25, 15],

have demonstrated excellent performance. This paper introduces a simple, memory-based algorithm for face recognition, termed ARENA, that satisfies the requirements outlined above and also significantly outperforms PCA-based methods on two standard face recognition datasets.

The remainder of the paper is structured as follows. Section 2 presents the ORL and FERET datasets. Section 3 describes the ARENA face recognition algorithm. Section 4 reviews Principal Components Analysis (PCA) and outlines two standard PCA-based algorithms for face recognition. Section 5 presents a variety of experiments designed to analyze ARENA’s behavior on the ORL and FERET datasets, and compares ARENA with the established baselines. Section 6 introduces a technique for further improving memory-based face recognition systems by augmenting the training set with synthetic images and examines its effects on both PCA-based techniques and ARENA. Section 7 summarizes numerous additional comparisons with other face recognition techniques. Section 8 examines the computational complexity and storage requirements for the algorithms and demonstrates that ARENA is competitive according to these metrics. Section 9 advances some hypotheses for ARENA’s surprising successes and places face recognition in the context of our visitor identification application. Section 10 concludes by presenting promising directions for future research.

## 2 Image Datasets and Preprocessing

Our results use human face images from two standard datasets: Olivetti-Oracle Research Lab (ORL) [21] and FERET [16, 18]. ORL consists of 400 frontal faces: 10 tightly-cropped images of 40 individuals with only minor variations in pose ( $\pm 20^\circ$ ), illumination and facial expression. The faces are consistently positioned in the image frame, and very little background is visible. FERET contains over 1100 faces; however many of them are unsuitable for our experiments since they are partial or full profiles, or the individuals were only photographed twice.



Figure 1: Top row: Two sample images each, of two subjects from ORL (left), and FERET (right). Note the difference in facial orientation, expression and accessories between the two images of the same individual. FERET images tend to exhibit greater variation in appearance (including hairstyle and clothing). Bottom row: the corresponding ARENA reduced-resolution images ( $16 \times 16$  pixels).

Therefore, from FERET, we selected the subset of images that satisfied the following two constraints: (1) near-frontal poses; (2) images of individuals with more than five such images (our tests require several images for each person). The resulting 275 images consist of 40 individuals, with greater variation in pose and lighting than in the ORL dataset. For instance, many of these images were taken over different days and display significant differences in hairstyles, eyewear, and illumination. Unlike the ORL images, the FERET faces are of non-uniform size and do not always appear in the same location of the image. We perform no explicit face extraction in the FERET images to explore the potential limitations of our template-based face recognition technique. The only preprocessing consists of simple intensity stretching.<sup>1</sup> Figure 1 shows two images for each of two individuals from the two datasets.

### 3 The ARENA Face Recognition Algorithm

ARENA is a memory-based [1] algorithm that employs reduced-resolution images ( $16 \times 16$ ) and the  $L_0^*$  similarity measure (described below). The reduced-resolution images are created by simply averaging over non-overlapping rectangular regions in the image. The distance from the query image to each of the stored images in the database is computed, and the label of the best match is returned.

#### 3.1 $L_p^*$ Similarity Measures

Our results show that the obvious choice for ARENA’s similarity measure, the Euclidian distance, performs poorly. In this section we present alternatives. The  $L_p$  norm is defined as:  $L_p(\vec{a}) \equiv (\sum |a_i|^p)^{\frac{1}{p}}$ . Thus, the Euclidian distance is simply:  $L_2(\vec{x} - \vec{y})$ . Note that since we are not interested in the actual distances, but only in the or-

<sup>1</sup>Intensity stretching, also termed intensity normalization, consists of scaling and shifting intensity values (by a constant amount) so that the intensities in the output image span the entire, available range (0 to 255).

dering, we can equivalently employ the similarity measure  $L_p^*(\vec{a}) \equiv (\sum |a_i|^p)$ .

Robust statistics literature shows that  $L_2^*$ , despite its convenient analytic properties, overly penalizes outliers [10]. For this reason, the  $L_1^*$  similarity measure is often used in noisy environments. For ARENA, we also explore the  $L_0^*$  similarity measure, defined as  $L_0^*(\vec{a}) \equiv \lim_{p \rightarrow 0^+} L_p^*(\vec{a})$ . Intuitively,  $L_0^*(\vec{x} - \vec{y})$  counts the number of components in  $\vec{x}$  and  $\vec{y}$  that differ in value.

In our application, each reduced-resolution image is converted into a vector,  $\vec{x}$ , where each pixel in the image is represented as a component of the vector. In practice, since individual pixel intensities are noisy, we relax the definition of  $L_0^*$  to be:

$$L_0^*(\vec{x} - \vec{y}) \equiv \sum_{|x_i - y_i| > \delta} 1$$

where  $\delta$  is a threshold, such that pixels whose intensities differ by less than  $\delta$  are considered equivalent. The experiments with different norms, presented in Section 5.1, indicate that the best performance on this task is achieved with  $p < 2$ .

### 4 Principal Components Analysis (PCA)

The most widely used baseline for face recognition, *eigenfaces* [25, 15] employs Principal Components Analysis (PCA), which is based on the discrete Karhunen-Loève (K-L), or Hotelling Transform [9], is the optimal linear method for reducing redundancy, in the least mean squared reconstruction error sense. Points in  $\mathcal{R}^d$  are linearly projected into  $\mathcal{R}^m$ , (where  $m \leq d$ , and typically  $m \ll d$ ). PCA has become popular for face recognition with the success of *eigenfaces* [25]. For face recognition, given a dataset of  $N$  training images (full-resolution originals, each with  $d$  pixels), we create  $N$   $d$ -dimensional vectors,  $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_N$ , where each pixel is a unique dimension.

The principal components of this set of vectors is computed as described in [9, 25] to obtain a  $d \times m$  projection matrix,  $W$ .

Now, the image  $\vec{x}_i$  may be compactly represented as *weights*,  $\vec{\theta}_i = (\theta_{i1}, \theta_{i2}, \dots, \theta_{im})^T$ , such that  $\vec{z}_i = \vec{\mu} + W\vec{\theta}_i$  approximates the original image, where  $\vec{\mu}$  is the mean of the  $\vec{x}_i$  and this reconstruction is perfect when  $m = d$ . The columns of  $W$  form an orthonormal basis for the space spanned by the training images.

Two variants of PCA for face recognition are evaluated in this paper, termed PCA-1, and PCA-2. For both algorithms, each training image is first projected into the eigenspace, and represented as a weight vector  $\vec{\theta}_i$ :

$$\vec{\theta}_i = W^T(\vec{x}_i - \vec{\mu}) \quad (1)$$

In PCA-1, the centroid of the weight vectors for each person’s images in the training set is computed and stored [25] — PCA-1 assumes that each person’s face images will be clustered in the weight space, so the actual training data is not needed. In PCA-2, a memory-based variant of PCA, each of the weight vectors is individually stored [12] — requiring more storage space, but providing PCA-2 with a richer representation. When a test image is presented to the system, it is first projected into the eigenspace (by Equation 1), and its weight vector  $\vec{\theta}_{\text{new}}$  is computed.  $\vec{\theta}_{\text{new}}$  is then compared against the stored weight vectors,  $\Theta$ , and the  $\vec{\theta}_k$  that is closest  $\vec{\theta}_{\text{new}}$  is located:

$$\vec{\theta}_{\text{best}} = \arg \min_{\vec{\theta}_k \in \Theta} L_2^*(\vec{\theta}_{\text{new}} - \vec{\theta}_k)$$

The label of  $\vec{\theta}_{\text{best}}$  is returned as the identity of the face represented by  $\vec{\theta}_{\text{new}}$ .

## 5 Baseline Comparisons

In the experiments described in this paper,  $n$  randomly-selected images for each individual in the dataset were placed in the training set, and the remaining images were used for testing. Multiple runs for each  $n$  were performed with different, random partitions between training and testing images, and the results were averaged.<sup>2</sup> The experiments were performed on both ORL and FERET images, and the results are reported separately so that they may be directly compared with other published results.<sup>3</sup>

<sup>2</sup>In testing ARENA, we exploit the fact that the distances between any two images in the dataset are independent of the test/train split, and contain sufficient information to efficiently enumerate the number of test/train splits that result in a correct identification for each image in the dataset. This allows us to effectively compute the average performance of ARENA over *all* possible test/train splits, without suffering the combinatorial explosion (as detailed in [13]).

<sup>3</sup>FERET is available from <jonathon@nist.gov>. ORL is available at <www.cam-orl.co.uk/facedatabase.html>. The list of FERET images used in our experiments, as well as Matlab code for our algorithms is available by contacting the authors.

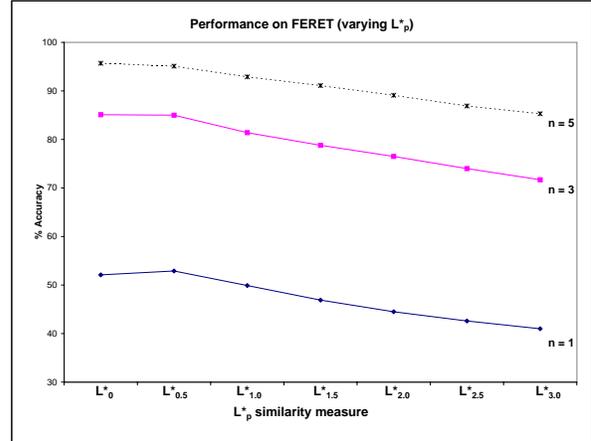


Figure 2: This figure shows how ARENA performs for different  $L_p^*$ . For  $L_0^*$ ,  $\delta$  was set to 10 (pixels range in value from 0 to 255). Note that the  $p \leq 1$  norms perform significantly better regardless of the number of training images ( $n$ ). Experiments conducted on the FERET database with the original images subsampled to  $16 \times 16$  images.

### 5.1 Experiments with similarity measures

In this section, we present experiments with different similarity measures of the form  $L_p^*$ , for  $0 \leq p \leq 3$ , on a variety of training set sizes, where  $n \in \{1, 3, 5\}$ , is the number of training images per person (see Figure 2). Note that the Euclidian metric (equivalent to  $p = 2$ ) does not perform well in comparison with the  $p \leq 1$  similarity measures. Therefore, in the remainder of the paper, we show results only with  $p \in \{0, 1\}$ . Similar experiments performed with PCA-1 and PCA-2 reveal that the change from  $L_2^*$  to  $L_1^*$  similarity measures does not improve the overall classification performance. In the cases where PCA accuracy was improved using  $L_1^*$ , it was still inferior to the comparable ARENA.

### 5.2 Experiments with different resolutions

Here, we examine how ARENA’s performance changes as the dimension of reduced-resolution images is varied. Each original face image is reduced to  $s \times s$  using simple local averaging. Figure 3 shows experiments with  $s \in \{2^k | k = 0, \dots, 5\}$ , and for  $92 \times 112$  full resolution images. Performance improves rapidly as  $s$  increases (over all training set sizes) and shows no significant improvement beyond  $s = 16$ . In the remainder of this paper, we present results with  $16 \times 16$  ARENA images.<sup>4</sup>

<sup>4</sup>Low-resolution images of similar dimensions are commonly used in the neural network literature [19].

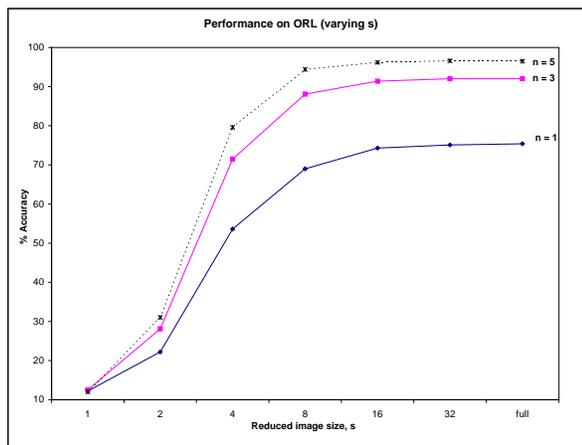


Figure 3: ARENA performance on the ORL dataset as  $s$ , the size of the reduced-resolution images is varied. “Full” indicates that the full-resolution image was used. ARENA’s performance improves rapidly with  $s$ , and plateaus by  $s = 16$ .

### 5.3 Comparisons with PCA-based techniques

PCA performance depends on the number of eigenvectors,  $m$ , that are stored. If  $m$  is too low, important information about the identity is likely to be lost. However, if  $m$  is too high, the weights corresponding to small eigenvalues will be noisy. This is analogous to selecting the appropriate subsampling ratio in ARENA. Similarly, ARENA’s performance is related to  $s$ , the dimensions of the reduced-resolution image. We aim to select an  $s$  such that computational complexity and storage are minimized without sacrificing classification accuracy.

Figure 4 shows that ARENA ( $p \in \{0, 1\}$ ) outperforms both PCA-1 and PCA-2, as the number of dimensions,  $m$  is varied. It is interesting to observe that, on the FERET dataset, the accuracy for PCA-1 ( $m = 10$ ) drops when the number of training images,  $n$  is increased from 3 to 5. This may be because the training faces for a given individual are not well-represented by the centroid of a single cluster. The memory-based techniques (PCA-2 and ARENA) are not adversely affected.

## 6 Augmenting the Training Set with Synthetic Images

Because we wish to perform recognition with the fewest number of training images per person, we augment the training set with additional, synthetically-generated face images. Since the task addressed in this paper is near-frontal face recognition, these images can be synthesized with simple geometric transformations (i.e., translation, rotation and scaling); more complex transformations

to account for out-of-plane rotations have been explored in [4, 27]. Incorporating synthetic training images into a memory-based model generally improves performance because it can increase the likelihood that an unknown query image will be matched to a correct instance in the memory base. Note that methods such as normalized correlation [5, 17] automatically account for translation, but do not address either rotation or scaling.

A number of synthetic images are generated from each raw training image by making small, random perturbations to the original image: rotation (up to  $\pm 5^\circ$ ); scaling (by a factor between 95% and 105%); and translation (up to  $\pm 2$  pixels, in each direction).<sup>5</sup> The process is similar in concept to the supplementary images used for neural network training for autonomous navigation [19] and automatic digit recognition [2, 22].

Figure 5 shows the improvement in ARENA’s performance (on both ORL and FERET datasets) when the training data for each person is augmented with 10 synthetic images per original. All instances of ARENA display some improvement due to the augmented memory base. Figure 6 shows the results of performing PCA on an augmented version of the ORL dataset. Interestingly, the synthetic data does *not* improve the PCA algorithms, with the exception of PCA-2 with  $m = 40$ .

## 7 Additional Experiments

In addition to the comparisons described above, with standard PCA techniques, we have extensively compared ARENA with other state-of-the-art face recognition algorithms. Due to space limitations, only a brief summary of these experiments is presented.

### 7.1 Standard PCA Variants

Recently, many modifications to the standard eigenface algorithm have been proposed and have been shown to work better in limited situations. We have duplicated two common such variants.

The first variant uses Mahalanobis distance [7] rather than standard Euclidian distance: PCA is initially used to reduce dimensionality by discarding eigenvectors corresponding to the lowest-magnitude eigenvalues (these are assumed to be noise). The remaining eigenvectors are then scaled such that their contributions to the distance are effectively equal. Unfortunately, in our experiments on this task, Mahalanobis-PCA does not consistently improve performance: Mahalanobis-PCA-1 is inferior to PCA-1 for

<sup>5</sup>It is important to correctly handle the border pixels in the synthetic images. For instance, a translation to the right uncovers pixels on the left, which should be filled with reasonable values. If these are zeroed (or assigned some other arbitrary value), the nearest-neighbor algorithm will be adversely affected. In ARENA, the border pixel values from the original image were replicated before perturbation, preserving the overall intensity of the border and preventing artifacts in the synthetic images.

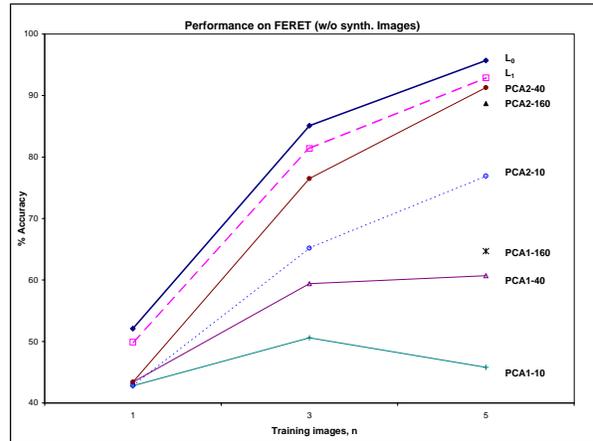
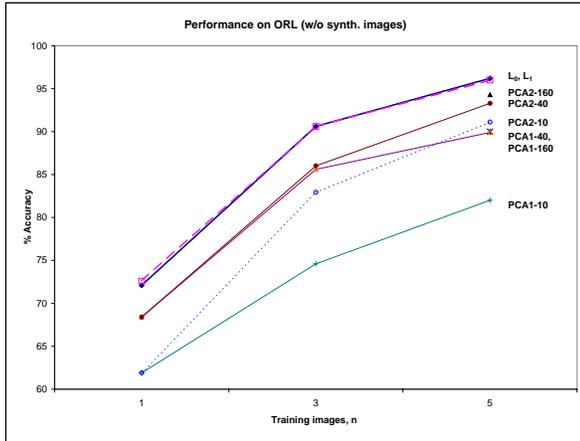


Figure 4: Performance of ARENA ( $p \in \{0, 1\}$ ,  $16 \times 16$  images), compared to PCA-nearest-centroid (PCA-1) and PCA-nearest-neighbor (PCA-2). For PCA algorithms, number of eigenvectors  $m \in \{10, 40, 160\}$ . Since  $m$  is limited by the rank of the training set matrix,  $m = 160$  can only be used when there are more than 160 training images (more than 4 training images for each of the 40 individuals,  $n \geq 4$ ). **Left:** ORL; **Right:** FERET.

low  $n$  or  $m$ , but slightly better in other cases; Mahalanobis-PCA-2 is uniformly inferior to PCA-2.

The second variant is motivated by the observation that the eigenvectors corresponding to the greatest eigenvalues often encode variations in illumination rather than the identity of the individual [3]. Consequently, if these eigenvectors are discarded (typically the top three [14]), then projecting the query image along the remaining eigenvectors should result in weights that do not encode these illumination effects. While there is some variation in lighting in both ORL and FERET datasets, we have found that PCA performance drops with this variant. It appears that the top three eigenvectors are (at least partially) encoding important information (supported by [3]).

In our experiments, even the best PCA algorithm,<sup>6</sup> which achieved an accuracy result of 94.8% in its best run, was outperformed (in identical experiments) by the average ARENA ( $L_0^*$  without synthetic images: 96.2%, with synthetic images, 97.1%).

## 7.2 Comparisons with other algorithms

We have also duplicated a set of experiments reported in [12]. They examined the performance of four algorithms, “Eigenfaces - average per class” (identical to PCA-1), “Eigenfaces - one per image” (identical to PCA-2), “PCA+CN” (PCA combined with a convolutional network classifier), and “SOM+CN” (Self-Organizing Map combined with a CN), on the ORL dataset with  $n$  ranging from 1 to 5. Table 1 summarizes these re-

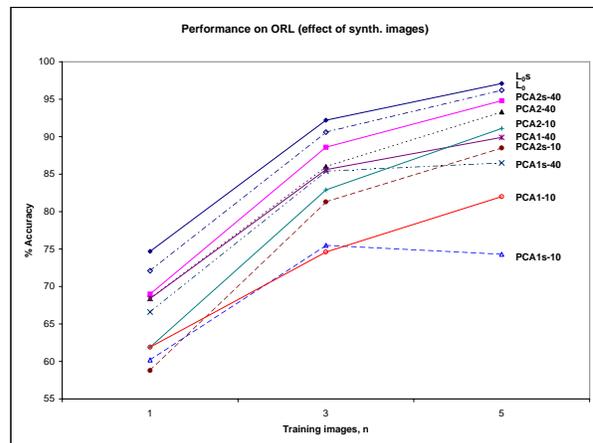


Figure 6: This figure shows the effect of adding 10 synthetic training images generated from each original to PCA-1 and PCA-2 ( $m \in \{10, 40\}$ ). Results for ARENA ( $p = 0$ ) are given for comparison. In all cases, PCA is outperformed by ARENA (even without synthetic images). Tests with synthetic images are marked with an ‘s’ above.

<sup>6</sup>This was PCA-2 with Euclidian distance,  $m = 40$ ,  $n = 5$  (without dropping top eigenvectors), and 10 synthetic images when tested on the ORL dataset.

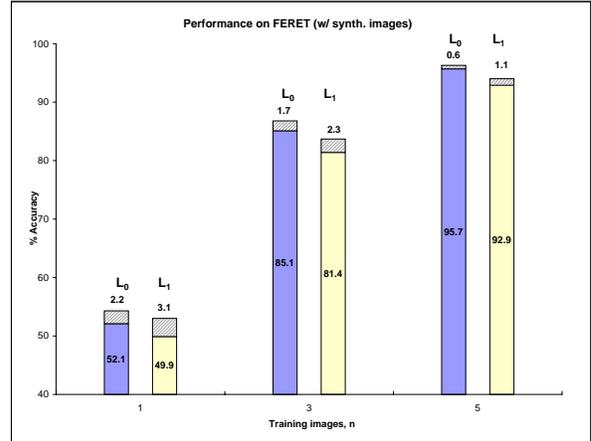
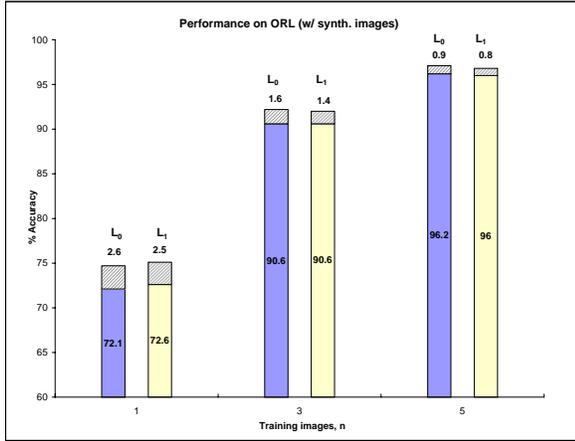


Figure 5: This figure shows how ARENA’s performance improves on the ORL (**left**) and FERET (**right**) when the memory base is augmented with 10 synthetic images generated from each original, for  $p \in \{0, 1\}$  and  $n \in \{1, 3, 5\}$ .

Images per person	1	3	5
Eigenface - avg per class	61.4%	71.1%	74.0%
Eigenface - one per img	61.4%	81.8%	89.5%
PCA+CN	65.8%	76.8%	92.5%
SOM+CN	70.0%	88.2%	96.5%
ARENA ( $p = 0, s = 16$ )	74.7%	<b>92.2%</b>	<b>97.1%</b>
ARENA ( $p = 1, s = 16$ )	<b>75.1%</b>	92.0%	96.8%

Table 1: Comparison of ARENA with results reported in [12].

sults. The last two rows of the table present results obtained with ARENA, augmented with the synthetic images ( $p \in \{0, 1\}$ ). Both variants of ARENA outperform all of the reported results.

Face recognition using Hidden Markov Models (HMM) on the ORL database is reported in [21]. Their best algorithm, with  $n = 5$ , obtained an accuracy of 88%, putting it between Lawrence’s implementations of PCA-1 and PCA-2. This is inferior to any ARENA variant.

## 8 Computational Complexity and Storage

In this section, we examine the computational complexity of PCA-1, PCA-2, and ARENA, and compare their storage requirements. We also discuss an important practical consideration: whether the algorithms can support incremental updates. We define the following terms:

- $c$  The number of people in the training set.  $c = 40$  for both ORL and FERET datasets.
- $n$  The number of training images per person. In our tests,  $n \in \{1, 3, 5\}$ .
- $N$  The total number of training images:  $N = cn$ .
- $d$  Each image is represented as a point in  $\mathcal{R}^d$ , where  $d$  is the dimensionality of the image.  $d = 10304$  for ORL and  $d = 98304$  for FERET.
- $m$  The dimension of the reduced representation: number of stored weights (PCA), or number of pixels ( $s^2$ ) in reduced-resolution ARENA. Normally,  $d \gg m$ .

The asymptotic behavior of the various algorithms is summarized in Table 2. The following observations are noteworthy. First, the training time for ARENA scales linearly with  $N$ , while both PCA-1 and PCA-2 training times scale poorly (due to the eigenvector computations inherent in the PCA algorithm). Second, the classification times for PCA-2 and ARENA are asymptotically slower than the corresponding time for PCA-1; however, ARENA avoids the  $dm$  term which is required for both PCA algorithms. Third, the storage space for ARENA is typically smaller than that of either PCA algorithm: ARENA always requires substantially less storage than PCA-2, and unless  $N$  is very large, ARENA also requires less storage than PCA-1. This is because ARENA performs all computations in the reduced dimensional space, and does not need to store any vectors of size  $d$  whereas any variant of PCA must store the projection matrix ( $m$  vectors of dimension  $d$ ). ARENA achieves further savings in storage space by quantizing the pixel values in the reduced-resolution image to a single byte (compared to the 8-byte double-precision values used for every element in the PCA methods). Our

experiments show that this quantization does not significantly reduce ARENA’s accuracy. From this, we conclude that not only does the ARENA algorithm perform better on the standard face recognition datasets, it is also faster and requires far less storage space.

Finally, standard PCA-based techniques cannot be trained incrementally (an important consideration in many applications) since the projection matrix,  $W$ , must be re-computed when new images are added to the system — an expensive operation (see Table 2). This can be avoided by assuming that the new images do not have a significant impact on the eigenspace: just compute the weights for the new images using the old projection matrix. Unfortunately, this is valid only if the system was initially trained on a very large set of individuals. Even with the shortcut of leaving  $W$  unchanged, PCA-based techniques do not scale as well as ARENA (see Table 2).

## 9 Discussion

Why does ARENA perform so well? Let us consider the behavior of the ARENA algorithm from two perspectives: (1) ARENA is performing a dimensionality reduction which, although non-domain-specific, is well-suited for face recognition since it reduces noise and compensates for small changes in the image; (2) it is performing a variant of template-matching on reduced-resolution images.

The first viewpoint indicates that ARENA is transforming high-dimensional points into a space that is manageable for nearest-neighbor algorithms. ARENA uses local averaging, which unlike PCA, is more robust to small image registration errors. The synthetic images further help the nearest-neighbor algorithm by populating the space with positive instances. Using synthetic images for standard PCA-based methods is expensive because the order-of-magnitude expansion in the training set results in high memory usage during the training phase (not to mention training time, as shown in Table 2).

From the second viewpoint, ARENA uses a large number of static templates (from several training images, augmented by the synthetic images). Template-matching has been used in early face recognition research [11], for facial-feature-detection; many recent approaches to face recognition can also be considered to be sophisticated versions of template-matching [5]. The FERET-96 test [17] includes a normalized correlation algorithm and shows that it is outperformed by several face recognition techniques. There are several significant differences between these correlation algorithms and ARENA. First, the correlation algorithms use high-resolution images and are therefore sensitive to small details in the image. Second, the use of the  $L_2^*$  similarity measure further exacerbate this sensitivity to unimportant differences. Consequently, two images of the same person with slightly different orientation or facial ex-

pression may be difficult to match. One of the strengths of the FERET-96 is that it tests the recognition ability of algorithms on subjects photographed over several sittings (spread over a year apart). Preliminary results using the publicly available FERET images on this task show no degradation in ARENA’s performance.

ARENA has recently been integrated into a visitor identification system [23]. The system obtains images from a security camera that monitors the front door of a building. Faces are extracted from these images using a neural-network-based face tracker [20], histogram-equalized and sent to ARENA. ARENA attempts to recognize the visitor and the system notifies interested parties of the visitor’s arrival. ARENA is particularly well-suited for this application because it supports incremental training: a human operator can label incorrect guesses and these are immediately incorporated into the training set. This system has been operational (24 hours a day) for several weeks. Note that the images gathered often display significant out-of-plane rotation, occlusion and extreme lighting conditions (half-faces); therefore, the images for a given individual can look very different. However, by acquiring many images for each common visitor, the system is able to robustly recognize these individuals in a variety of situations. Under these challenging conditions, we are pleased to report overall accuracies of 55% for an image set containing 50 individuals (more than 1000 training images, added incrementally over a period of several months). To test the system further, we added 1500 “distractor” images of faces collected from the web and tagged them with the single label “stranger”. There has been no noticeable drop in classification performance of known visitors, but unknown visitors are often correctly classified as “stranger”. Detailed performance statistics on the visitor identification task are forthcoming.

## 10 Conclusions and Future Work

This paper demonstrates that ARENA, a very simple algorithm, can significantly outperform established face recognition algorithms on standard datasets. Unlike the standard PCA-based algorithms, ARENA easily handles incremental updates to the face recognition database and has been shown to scale well. Given the algorithm’s simplicity, ARENA’s high-performance is somewhat surprising. We invite other researchers to independently confirm our findings, and plan to enter ARENA (or its latest variant) in the next FERET test.

We are extending the work described here in several directions. First, we are comparing ARENA against Fisher Discriminant Analysis (FDA) [3, 24] approaches to frontal face recognition. Most FDA methods require a dimensionality-reduction step, traditionally performed using PCA. We are exploring whether reduced-resolution im-

Method	Training time	Classification time	Storage space	Incremental update cost	
				Recomputed $W$	Unchanged $W$
PCA-1	$O(N^3 + N^2d)$	$O(cm + dm)$	$O(cm + dm)$	$O(N^3 + N^2d)$	$O(md)$
PCA-2	$O(N^3 + N^2d)$	$O(Nm + dm)$	$O(Nm + dm)$	$O(N^3 + N^2d)$	$O(md)$
ARENA	$O(Nd)$	$O(Nm + d)$	$O(Nm)$		$O(d)$

Table 2: Comparison of asymptotic behavior. ARENA displays clear advantages over both PCA-based techniques.

ages (as used in ARENA) can perform this role. We also plan to investigate the effects of using different  $L_p^*$  similarity measures in such algorithms. Other experiments with wavelet-based schemes for dimensionality reduction in combination with support vector machines [26] are in progress.

Finally, a useful byproduct of the visitor identification system is the collection of a labelled face dataset, with unposed images captured in the natural lighting of an outdoor environment. This dataset will be made available on our website to researchers in the near future.

## Acknowledgments

The visitor identification system was developed in collaboration with Robert Stockton. We would like to thank Henry Rowley, Henry Schneiderman, and Gita Sukthakar for valuable feedback on this paper. Thanks also to Keiko Hasegawa and Erich Greene for ensuring that the statistical tests were rigorous.

## References

- [1] C. Atkeson, W. Moore, and S. Schaal. Locally weighted learning. *AI Review*, 11, 1997.
- [2] S. Baluja. Making templates rotationally invariant: An application to rotated digit recognition. In *Advances in Neural Information Processing Systems*, 1998.
- [3] P. Belhumeur, J. Hespanha, and D. Kriegman. Eigenfaces vs. fisherfaces: Recognition using class specific linear projection. *IEEE Transactions on PAMI*, 19(7), 1997.
- [4] D. Beymer and T. Poggio. Face recognition from one example view. In *Proceedings of ICCV*, 1995.
- [5] R. Brunelli and T. Poggio. Face recognition: Features versus templates. *IEEE Transactions on PAMI*, 15(10), 1993.
- [6] R. Chellappa, C. Wilson, and S. Sirohey. Human and machine recognition of faces: A survey. *Proceedings of the IEEE*, 83(5), 1995.
- [7] R. Duda and P. Hart. *Pattern Classification and Scene Analysis*. John Wiley and Sons, 1973.
- [8] T. Fromherz. Face recognition: A summary of 1995–1997. Technical Report ICSI TR-98-027, International Computer Science Institute, Berkeley, 1998.
- [9] R. Gonzales and R. Woods. *Digital Image Processing*. Addison-Wesley, 1992.
- [10] P. Huber. *Robust Statistics*. Wiley, 1981.
- [11] T. Kanade. *Picture Processing by Computer Complex and Recognition of Human Faces*. PhD thesis, Kyoto University, 1973.
- [12] S. Lawrence, C. Giles, A. Tsoi, and A. Back. Face recognition: A hybrid neural network approach. Technical Report UMIACS-TR-96-16, University of Maryland, 1996.
- [13] M. Mullin and R. Sukthakar. An efficient technique for calculating exact nearest-neighbor accuracy. Technical report, Just Research, 1999. (forthcoming).
- [14] A. O’Toole, H. Abdi, K. Deffenbacher, and D. Valentin. Low dimensional representation of faces in high dimensions of the space. *Journal of the Optical Society of America A*, 10, 1993.
- [15] A. Pentland, B. Moghaddam, and T. Starner. View-based and modular eigenspaces for face recognition. In *Proceedings of Computer Vision and Pattern Recognition*, 1994.
- [16] P. Phillips, H. Moon, P. Rauss, and S. Rizvi. The FERET evaluation methodology for face-recognition algorithms. In *Proceedings Computer Vision and Pattern Recognition*, 1997.
- [17] P. Phillips, H. Moon, P. Rauss, and S. Rizvi. The FERET september 1996 database and evaluation procedure. In *Proceedings of Audio and Video-based Biometric Person Authentication*, 1997.
- [18] P. Phillips, H. Wechsler, J. Huang, and P. Rauss. The FERET database and evaluation procedure for face-recognition algorithms. *Image and Visual Computing*, 16(5), 1998.
- [19] D. Pomerleau. *Neural Network Perception for Mobile Robot Guidance*. PhD thesis, Carnegie Mellon University, February 1992.
- [20] H. Rowley, S. Baluja, and T. Kanade. Neural network-based face detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(1), 1998.
- [21] F. Samaria and A. Harter. Parametrisation of a stochastic model for human face identification. In *Proceedings of IEEE Workshop on Applications on Computer Vision*, 1994. ORL database is available at: <www.cam-orl.co.uk/facedatabase.html>.

- [22] B. Schölkopf, C. Burgess, and V. Vapnik. Incorporating invariances in support vector learning machines. In *Artificial Neural Networks — ICANN'96*, 1996.
- [23] R. Sukthankar and R. Stockton. Argus: An automated multi-agent visitor identification system. In *Proceedings of AAAI-99*, 1999.
- [24] D. Swets and J. Weng. Using discriminant eigenfeatures for image retrieval. *Transactions on PAMI*, 18(8), 1996.
- [25] M. Turk and A. Pentland. Eigenfaces for recognition. *Journal of Cognitive Neuroscience*, 3(1), 1991.
- [26] V. Vapnik. *Statistical Learning Theory*. Wiley, 1998.
- [27] T. Vetter. Synthesis of novel views from a single face image. *International Journal of Computer Vision*, 28(2), 1998.