

Multi-Layer Architectures for Facial Action Unit Recognition

Tingfan Wu, Nicholas J. Butko, Paul Ruvolo, Jacob Whitehill, Marian S. Bartlett, Javier R. Movellan
Machine Perception Laboratory, University of California San Diego

{ting,nick,paul,jake,marni,movellan}@mplab.ucsd.edu

Abstract—In expression recognition and many other computer vision applications, the recognition performance is greatly improved by adding a layer of non-linear texture filters between the raw input pixels and the classifier. The function of this layer is typically known as feature extraction. Popular filter types for this layer are Gabor Energy Filters (GEFs) and Local Binary Patterns (LBPs). Recent work [39] suggests that adding a second layer of non-linear filters on top of the first layer may be beneficial. However, it is unclear what is the best architecture of layers and selection of filters. In this paper, we present a thorough empirical analysis of the performance of single-layer and dual-layer texture based approaches for AU recognition. For the single hidden layer case, GEFs perform consistently better than LBPs, which may be due to their robustness to jitter and illumination noise as well as to their ability to encode texture at multiple resolutions. For dual-layer case, we confirm that, while small, the benefit of adding this second layer is reliable and consistent across datasets. Interestingly for this second layer LBPs appear to perform better than GEFs.

Index Terms—Gabor energy filters, local binary patterns, facial expression recognition, action unit recognition



1 INTRODUCTION

Expression recognition systems are organized as a serial pipeline. The first layer in this pipeline detects and normalizes (registers) face images. The input to this layer is an image and the output is a collection of patches with the detected faces normalized to a common size. The second layer applies a bank of non-linear filters to the image patches found by the previous layer. The goal of the second layer is to make the mapping from images to expression categories easier to separate by a classifier and to filter out irrelevant information, *e.g.* jitter or illumination noise. The final layer consists of a statistical classifier that converts the output of the filter bank into numbers representing the presence or absence of target facial expressions, *e.g.* Multinomial Logistic Regression (MLR) or a Support Vector Machine (SVM), .

Systems differ in the type of filters used in the intermediate layer. Geometric approaches utilize image pixel values exclusively as a means to extract the relative location and shape of a set of facial landmarks (*e.g.* corners of the mouth, brows) [28], [35], [1], [27], [42]. Texture-based approaches [25], [37], [48] use filters that describe the local texture information in the face image patches, without explicitly representing landmark locations.

Historically, the computer vision community focused first on geometric approaches perhaps because they closely matched our intuitions of how humans perceive faces, and because it was thought they would be less sensitive to changes in illumination and pose than texture-based approaches. Evidence has accumulated over the past 10 to 20 years that texture based approaches perform consistently better than geometric approaches [30], [23], [9]. The reason is that in many cases, encoding the location and shape of facial landmarks turns out to be more difficult than directly encoding the target facial expressions. Moreover, some expressions (*e.g.* AU6 and

AU7 in the Facial Action Coding System) are very similar in shape but different in appearance, *e.g.* furrow lateral to the eyes in AU6 but not in AU7 [9]. In practice current expression recognition systems are predominantly texture based or use geometry information to marginally boost the performance of a texture based system. Pure shape-based approaches have for the most part been abandoned.

The simplest texture-based approaches send the pixel values from the face images directly to the statistical classifiers. It is clear by now that such systems are consistently outperformed by architectures that have an intermediate layer of texture filters between the pixel values and the statistical classifier [10], [45]. Two of the most popular texture filters in facial expression recognition are Gabor Energy Filters (GEFs) [32], and Local Binary Patterns (LBPs) [48]. Gabor energy filters became of interest to the computer vision community in part due to the fact that they are mathematical models of the complex cells in primary visual cortex in the brain [32]. Some of the most successful expression recognition systems to date use banks of GEFs as their primary representation [24], [4], [3], [25]. LBPs originated in the computer vision literature on texture analysis [48] and since then have become popular in the facial expression recognition literature [45]. Interestingly, the basic kernels used in LBP filters are reminiscent of high spatial frequency, high bandwidth Gabor filters (see Section 2.4).

1.1 Shallow vs. Deep Architectures

In the language of neural networks, current expression recognition architectures are predominantly “shallow” or have a “single hidden layer”: The input layer is the set of raw pixel intensities, the hidden layer is a bank of non-linear texture filters, and the output layer is a single layer classifier. Shallow architectures are currently

dominant in the machine learning community due in part to the great popularity and success of two shallow machine learning algorithms: SVMs and Boosting [43], [14], [15]. However, recent years have seen a renewed interest in architectures that use more than one hidden layer (*i.e.* deep networks). This interest was rekindled partially by the development of efficient learning algorithms for training multilayer (deep) networks [19]. In computer vision deep architectures proved successful in problems such as digit recognition [22] and are now seeing a comeback [40], [7]. Interestingly the winning architecture for the 2011 Facial Expression Recognition and Analysis (FERA) challenge [41] utilized a 2-hidden layer architecture, named LGBP. This architecture applies LBP texture filters to the output of a layer of GEFs.

In this paper we present a thorough empirical analysis of the performance of single layer and dual layer texture based approaches on a difficult facial expression recognition task: FACS Coding [12]. In particular we analyze how these different approaches generalize within and between datasets and how they resist perturbations due to local changes in illumination and registration jitter. The paper is organized as follows: First we describe the system we submitted to the FERA competition, which utilized a single hidden layer architecture and placed second. In the second part of the paper we describe a post-competition analysis of texture based architectures. Particularly we focus on the comparison between LBP and Gabor filters and the comparison between single and dual hidden layer architectures.

2 COMPETITION SYSTEM

The goal of the FERA competition was to automatically FACS code a collection of videos of actors making facial expressions. The Facial Action Coding System (FACS) is a taxonomy of facial expressions based on a combination of 57 elementary components. These elementary expressions, known as Action Units (AUs) and Action Descriptors (ADs), can be seen as the “phonemes” of facial expressions: words are temporal combinations of phonemes. Facial expressions are spatial combinations of AUs. We entered the FERA competition with a previously developed FACS recognition system, named the Computer Expression Recognition Toolbox (CERT) [4], [3], [25]. Our focus was on analyzing how to best adapt an existing system, such as CERT, to generalize on a new dataset. Figure 1 describes the pipeline of the CERT system. In this section, we examine the various building blocks: their properties, common failure modes and our efforts to optimize them.

2.1 Datasets

FERA - The focus of the FERA 2011 (AU) competition was the GEMEP-FERA [2] dataset. This dataset consists of recordings of 10 actors displaying a range of expressions, while uttering a meaningless phrase, or the word “Aaah”. There are 7 subjects in the training data, and

6 subjects in the test set, 3 of which are not present in the training set. The training videos came with frame-by-frame binary labels (present or not) of 36 AUs and AD50¹ (talking). Intensity ratings, location of apex and reliability information were not available. Among the labeled AUs, only the prediction of 12AUs are evaluated on the test set (Table 2).

The authors of this paper never had access to the labels of the test set. We had access to the test images themselves but we chose not to label them in any way to estimate performance or to train our classifiers.

FFD07 - This is the original training set used to develop CERT. It is a combination of the following databases: Ekman-Hager [10], Cohn-Kanade [20], MMI [34], M3 [13], and two non-public datasets collected by the United States government which are similar in nature to M3.

2.2 Training Frame Selection

A priori it is not clear whether using every image in the data for training will result in optimal classifier performance. *E.g.*, a large sequence of nearly identical frames may result in the classifier giving excessive weight to a particular instance of a facial expression. It is also possible that training on the onset/offset points, when the intensity of the expression is very low, may be counterproductive.

A common strategy to deal with this issue is to select only the apex frame from each labeled AU event, *i.e.* the frame in which the expression intensity is maximal [26]. Practically, as the competition dataset came with only frame-by-frame binary labels instead of onset-apex-offset event information, an AU event was first localized by finding consecutive positive-labeled frames in a video, then the center frame of the sequence was selected as an approximation to the apex frame. However, this per-AU strategy ignores the fact that other co-occurring AUs may change the appearance of the target AU. For example, the brow appearances of AU1+AU4 (inner brow raiser + outer brow raise) differs greatly from the brow appearance of AU1 by itself. Therefore, instead of looking at each AU independently, we searched for unique AU-combinations in the video. When there were multiple disjoint sequences of an AU combination, only the center frame of the first event was selected.

We tried three different frame selection schemes: no frame selection, *i.e.* use all frames (NoFS); unique combinations among the 12 training AUs + AD50 (FS12) (as was done in the baseline paper [41]), and unique combinations among all 36 labeled AUs and AD50 (FS36).

Figure 2 shows the cross validation performance using the 3 different frame selection schemes. No significant performance difference were found between the three approaches. However, frame selection reduced the size of the FERA dataset from 5264 (NoFS) images into 627 (FS12) and 934 (FS36). Having less training frames

1. Action Descriptor 50.

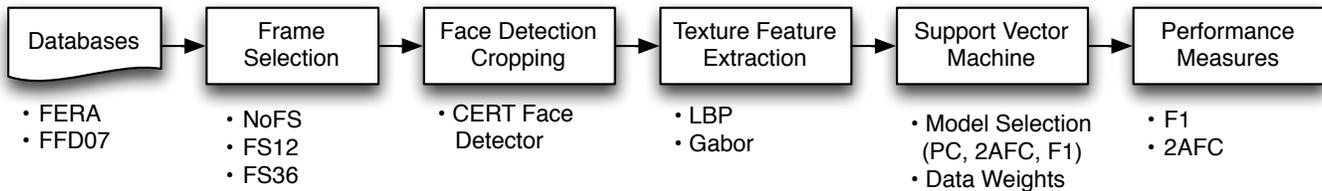


Fig. 1. The basic processing pipeline for approaches in this paper. Both the baseline method and CERT are special cases of the pipeline.

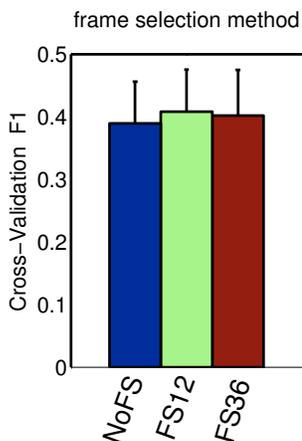


Fig. 2. Leave-one-subject-out cross validation performance on the GEMEP-FERA training set using different frame selection methods. The F1-score is averaged across cross-validation folds and AUs.

significantly sped up the training process giving up opportunities to test differences ideas. After all, we chose the FS36 approach in the rest of the paper for its slightly better performance (though not significant) and much shorter training time comparing to NoFS.

2.3 Face Detection and Registration

We used the CERT face detector and feature detectors [11] for face registration. Both detectors are of the Viola and Jones type [44] except that the face detectors discriminates face or non-face patches in a image while the feature detectors discriminates a face feature from the rest of the regions in a face. The CERT face detector detected all the faces in the competition training set with no false alarms except for the blank frames.

In the registration process, first the face detector finds a face in the image. Second, the feature detectors localize face features in the face, including eye corners, nose, and mouth corners. Third, the face is normalized using a planar least-square fitting procedure on the detected features (*i.e.* Procrustes analysis [16]). Finally, the aligned face is scaled into a 96×96 pixel matrix in which the intra-ocular distance is 48-pixels. The location of the detected facial features is discarded and the 96×96 pixel image matrix is sent to the next layer of processing.

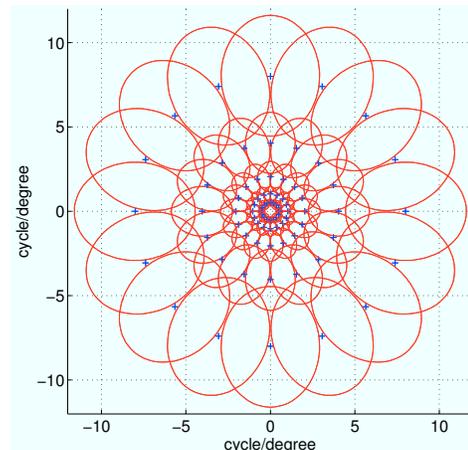


Fig. 3. Spectrum of the Gabor filter bank. The '+' centroids are the peak frequencies and the ellipsoids are the half-magnitude contour of each filter. Pairs of a opposing centroids form one GEF, so 80 centroids yield 40 GEFs. We follow the convention that a face patch spans 4° of horizontal and vertical visual angle.

2.4 Texture filters

The version of CERT we used for the competition employed a hidden layer representation with 40 Gabor energy filters (Fig. 3 plots the filter spectra), (5 frequencies, 8 orientations). This is slightly different from the standard version of CERT, which uses 72 filters [4], [3], [25]. In recent work we found this new 40 filter version to work as well as the 72 filter version.

The baseline algorithm provided by the FERA organizers used an LBP image representation consisting of histograms of the 59 neighborhood-8 radius-1 uniform local binary pattern values in non-overlapping 20×20 pixel blocks of the 200×200 cropped face, so that there were 10×10 blocks. In our implementation, the LBP operated on smaller 96×96 cropped faces.

There are interesting similarities between Gabor filters and uniform LBP filters [33]. Both approaches are spatially local oriented edge detectors with some robustness to translation and changes in illumination. Figure 4 shows visualizations of both features, which highlights their similarity. LBPs achieve robustness to translation when their histograms are pooled over local blocks. GEFs achieve it by the fact that they are invariant to local phase. LBPs detect different types of local neighborhood

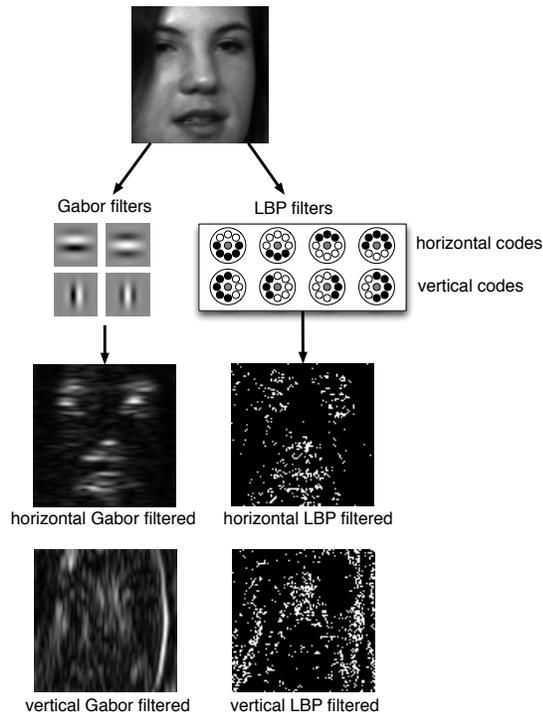


Fig. 4. Examples of Gabor and LBP filters and the filtered images. LBP encodes each pixel in the input by its neighborhood. Among the 59 possible uniform LBP_{8,1}^{u2} codes, the figure shows a subset of codes that detect horizontal (top row) or vertical (bottom row) edges passing through the center pixel. In the LBP-filtered image, the white pixels represent the pixels encoded by the corresponding horizontal or vertical edge detecting LBP subsets.

(center surround, edge, or corner) by having different binary codes. GEFs characterize local neighborhoods by combining filters of multiple frequency and orientations. Both methods can identify edges of different spatial extent by adjusting their radius/scale. However, LBPs are typically implemented with single fixed radius. From a signal processing point of view, many of the LBP kernels can be seen as high frequency, high bandwidth bandpass filters.

Empirically, we found GEFs performed better than LBPs (see Section 4.1). Therefore, we used Gabor features in our competition submission.

2.5 Data-Weighted Support Vector Machines

Support vector machines [5] were used to map the output of the texture filters, either GEFs or LBPs, into action unit categories. One of the main challenges we encountered was how to best combine FFD07, the prior data that had been used for training CERT, with the new FERA training set. The most straightforward way to adapt CERT to the FERA dataset is to combine FERA and FFD07 and retrain on the combined dataset. However, this method has some potential disadvantages: The FFD07 dataset may introduce information that is

counterproductive for the FERA challenge. For example, FFD07 has a significant proportion of individuals of Asian and African descent, while the subjects in GEMEP-FERA are strictly European. This makes CERT more applicable to a diverse population but may deteriorate performance on the FERA dataset, which is a non-representative subset. Moreover the FERA database is quite small (FS36: 934 images) after frame selection. Thus the learning algorithm (SVM) could be easily overwhelmed by the instances from FFD07 (8000+ instances).

To address this problem we developed a custom version of SVM with data weights that can be individually adjusted for each training example. Given training data $\{\mathbf{x}_i\}$, and labels $\{y_i\}$, the primal formulation of the data-weighted SVM is

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \overbrace{c_i}^{\text{data-wt.}} \xi_i \quad (1)$$

$$\begin{aligned} \text{s.t. } & y_i(\mathbf{w}^T \mathbf{x}_i - b) \geq 1 - \xi_i \\ & \xi_i \geq 0 \end{aligned} \quad (2)$$

As in the traditional SVM formulation, (\mathbf{w}, b) defines the hyperplane to be learned, ξ_i are the slack variables, and C is the master data fitness parameter; the added c_i are the data-weights, controlling the fitness to each data instance. The data-weighted SVM reduces to standard SVM when $c_i = 1$ for all data instances. Both linear kernel and non-linear radial basis function (RBF) kernel ($K(\mathbf{x}_i, \mathbf{x}_j) = e^{-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2}$) were used. For the linear kernel, the hyper-parameter is the regularization parameter (C). and for the non-linear kernel, there is an additional parameter, inverse kernel width (γ).

For each AU, a binary SVM was trained with all the selected frames using the optimal hyperparameter selected for the AU. The two hyperparameters were selected from a candidate grid of values using cross-validation accuracy [6]. However, our targeted performance measure is F1 which has very different properties from other performance measures such as 2AFC (see Sec. 2.6). We decided to try and optimize three different performance measures: percent correct (PC), two-alternative forced choice (2AFC), and F1-score.

2.6 Performance Measures

Most statistical classifiers produce real-valued scores that can be interpreted as evidence for the observed data belonging to one of two categories of interest, here referred to as the positive and negative categories. The sensitivity of the classifier depends exclusively on the statistical properties of the real valued score, not the threshold. However many applications require making binary decisions. For example, the smile shutter feature in some digital cameras needs to decide whether or not to take a picture based on the evidence that a person is smiling. These decisions are typically made based on whether the real-valued scores provided by the classifiers pass a given threshold. In such cases the

performance is a function of the sensitivity provided by the real valued scores, and of the threshold chosen for converting those scores into binary outputs. A system with good sensitivity may appear to perform poorly for a specific problem if the threshold is not chosen judiciously.

The area under the receiver operating curve (ROC) is a popular measure of performance used in the pattern recognition community [17]. An advantage of the ROC score is the invariance to the prior probabilities of the two categories, thus facilitating comparisons across datasets and categories with different priors. In addition, the ROC score is also invariant to monotonic transformations of the continuous-valued score provided by the classifiers. One problem with the ROC is that there is not a standard way to compute it. The actual value depends on the numerical integration scheme used to estimate the area under the ROC from a finite number of points. It is our experience that different methods can result in ROC scores that differ from each other by up to 5 percentage points. In psychophysics a popular task for measuring the sensitivity of an observer independently of its bias is the two-alternative forced choice task (2AFC). In this task observers are presented with two randomly chosen examples of the two categories of interest. One of the examples is from the positive category and the other from the negative category. The goal of the observer is to choose the positive example. A well known result from the theory of signal detectability shows that the expected percent correct in the 2AFC task equals the area under the ROC [17], [29]. For example if an observer has an area under the ROC of 0.9 it has a 90 % chance of correctly choosing the correct alternative in the 2AFC task. Hereafter we will use the term “2AFC score” rather than “area under the ROC”. The reasons are: (1) the 2AFC is very easy to compute (see Appendix); (2) it does not depend on numerical integration methods; (3) it has an intuitive interpretation as the probability of correctly categorizing a randomly chosen positive and negative example; (4) it can be interpreted as one of the many algorithms to estimate the area under the ROC.

The official evaluation metric for the FERA competition was the F1 score, a popular measure of performance in the document retrieval community. F1 is a function of the sensitivity of the system, the prior probabilities of the two categories, and the threshold used to make such decisions. A common misconception about F1 is that it favors low false alarm rates. This is not necessarily the case. If the system has low sensitivity (*e.g.* the random baseline in Table 1), the F1 score is maximized by having a large false alarm rate, which may be undesirable in some applications. Finally, another popular measure of performance is the percent correct (PC), *e.g.*, the proportion of expressions correctly classified. This is also a function of the sensitivity of the system, the prior probabilities of the different categories, and the chosen threshold.

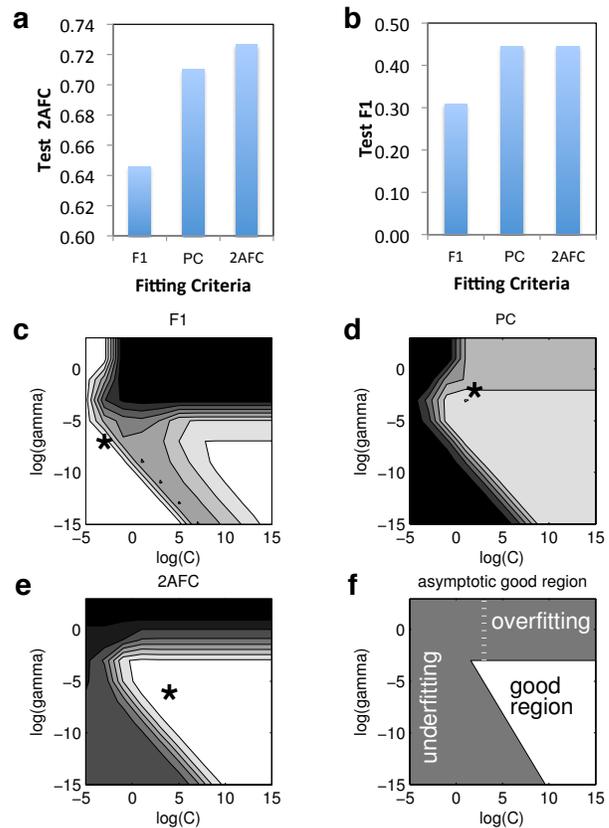


Fig. 5. (a) F1 and (b) 2AFC double cross validation scores from SVMs with parameters optimized for different performance measure (c)-(e) An example of parameter selection with different performance measure. (f) asymptotic properties on non-linear SVM model selection grid from [21]. The figures show the cross-validation performance surface for each grid search point of SVM parameters C and γ , the brighter pixels correspond to higher scores. The black “*” denotes the final chosen parameter. Unlike PC and 2AFC, the surface of F1-scores is not uni-modal, which makes parameter selection hard.

2.7 Hyperparameter Fitting Criteria

The target of the competition was to optimize the F1 score averaged across all the target AU categories in a generalization set. One question of interest to us was whether optimizing hyperparameters (*e.g.*, kernel width, regularization constant) with respect to another performance measure (*e.g.* 2AFC, PC) yields a higher F1 score on generalization sets.

To this end we compared the F1 generalization performance when the RBF-kernel SVM hyperparameters were optimized with respect to F1, 2AFC, and PC using double cross-validation. Surprisingly we found that optimizing hyperparameters with respect to 2AFC or PC resulted in better F1 score than optimizing with respect to F1 (see Fig. 5(b)). Further investigation revealed why the F1-score was not a good parameter selection criterion. As Fig. 5(c) shows, the performance landscape for F1 is multi-modal (two separate white regions) thus

TABLE 1

reported performance on the FERA blind test set. See Sec.3 text for detail.

Method	F1 Score			2AFC		
	ind	dep	all	ind	dep	all
Official random	.531	.471	.512	.500	.500	.500
Official LBP baseline	.453	.423	.451	.631	.611	.628
CERT unmodified	.569	.514	.550	.746	.692	.723
CERT retrained on FERA	.604	.539	.583	.759	.753	.758
Winner [36]	.631	.610	.628	.763	.751	.752

ind: test on new subjects not in the training set

dep: test on new videos of subjects seen in the training set

all: test on mixture of both “indep” and “dep” cases

making hyperparameter selection difficult and suggesting greater expected variability in generalization tests. In addition to the typical “good parameter region” for SVMs [21] (see Fig. 5(d-f)), the F1 score has additional peaks for small C parameters (large regularization). This leads to under-fitting of the SVM weights, and results in classification models that predict everything as positive. This may be due to the fact that for low sensitivity systems, F1 is optimized by using a threshold that encourages false alarms. The other two performance measures, 2AFC and PC, don’t seem to suffer from this issue. Therefore, we used 2AFC to optimize hyperparameters in the rest of the paper. A previous study proposed that the F1 score could be preferable because algorithms were found that performed equally well with respect to the 2AFC but differed with respect to the F1 score [38]. Our results suggest that the F1 score may just be a noisier estimator of performance than 2AFC, potentially producing spurious differences between algorithms that are unrelated to their actual sensitivity.

3 COMPETITION RESULTS

Table 1 presents the official performance results of different algorithms, including, random, baseline, two of our methods, and the winning method. The random result is from a zero-sensitivity classifier which says “yes” for all frames. The official baseline results were provided by the FERA challenge organizers. We obtained second overall performance with respect to the F1 score, the target of the competition. However we obtained the best overall performance with respect to the 2AFC score, which was not the target of the competition. Since the F1 score is sensitive to threshold while 2AFC is not, this suggests that the winning team chose thresholds better than we did.

3.1 Unmodified CERT Outputs

First we applied CERT (previously trained on M3 before competition) directly on the FERA dataset without any training or adjustment. In our original conference submission, the CERT generalization performance in all categories outperformed all of our approaches trained only on the FERA database, without having seen a single

TABLE 2

Cross database comparison of 2AFC “overall” scores.

AU	CERT on M3 CV	CERT on FERA	CERT retrained on FERA	Baseline on FERA
1	.823	.747	.805	.790
2	.812	.745	.866	.768
4	.756	.719	.776	.526
6	.955	.835	.862	.657
7	.773	.701	.707	.555
10	.731	.681	.718	.597
12	.901	.815	.869	.724
15	.831	.643	.585	.563
17	.840	.639	.761	.646
18	.780	.649	.779	.610
25	.768	.760	.720	.593
26	.801	.691	.650	.500
Avg	.814	.719	.758	.628

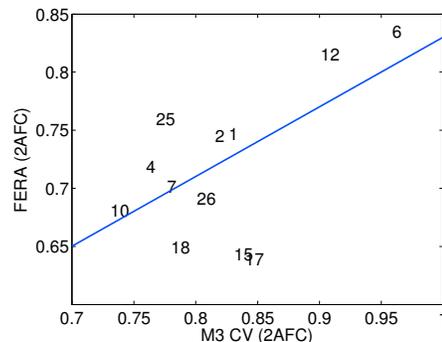


Fig. 6. CERT performance on the M3 vs the FERA dataset for different AUs. The trend is linear except for AU 15, 17, and possibly 18

FERA image [46]. Table 2 displays the per-AU 2AFC scores, which includes previously published results on the M3 dataset and the scores obtained on the entire FERA dataset. The M3 results were obtained using single cross-validation methods (labeled M3CV), and thus they represent generalization within a dataset, while the FERA results represent expected generalization to a new dataset. When applied to the FERA dataset CERT took an average performance hit of only 9.1 percentage points. This is quite remarkable considering the different nature of the FERA dataset when compared to the M3 dataset. Figure 6 shows a scatter plot of the AU by AU performance of CERT on the M3 *vs.* the FERA datasets. The performance on the two datasets is highly correlated, which suggest the relative difficulties between AUs on the two database were comparable. However there were some outliers: AU15, AU17, and arguably AU18. Performance on these AUs was worse than expected. One possible explanation is that for these AUs the criteria utilized by the human coders of FERA may have been particularly different from the criteria used by the human coders of M3.

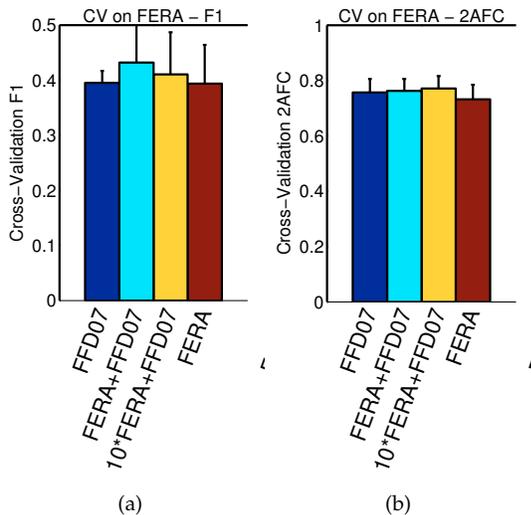


Fig. 7. Cross validation performance of SVMs trained with different data weight settings.

3.2 Retraining CERT on FERA

Table 3 describes the SVM data-weighting schemes for retraining CERT using FERA data: (1) Retrain on FERA only. (2) Retrain giving equal weight to FERA and FFD07. (3) Retrain giving 10 times more weight to FERA than FFD07.

TABLE 3
The weighting scheme used in the competition

weighting scheme	c_i for FERA	c_i for FFD07
FFD07	0	1
FERA+FFD07*	1	1
10*FERA+FFD07	10	1
FERA	1	0

Double cross validation (See Fig. 7(a-b)), indicated that best performance was obtained by retraining CERT on FERA and FFD07 with equal weights and thus our final competition system was based on that approach. Both the overall (Table 1) and per-AU performance for our best competition system are shown in Table 2 and Table 4). For subject independent tasks, the performance gains of adding FERA were found in AUs that were particularly abundant, such as AU1, AU2 and AU4. For the subject dependent case, the gain was particularly noticeable on the poor performing AUs. However the overall performance improvements were marginally better. Similar performance results could have probably been obtained by retraining CERT on FERA alone.

4 POST COMPETITION ANALYSIS

There were some interesting differences between our system and the winning system. We used exclusively texture based filters and the winning system used a combination of texture based and geometric approaches. We were particularly intrigued however by the fact

TABLE 4
Official Test F1 Score of CERT+FERA method

AU	CERT			CERT+FERA		
	indep	dep	all	indep	dep	all
1	.496	.521	.503	.765	.399	.634
2	.689	.394	.564	.736	.485	.636
4	.596	.593	.595	.608	.590	.602
6	.804	.704	.777	.788	.683	.759
7	.579	.632	.601	.563	.660	.604
10	.502	.574	.528	.545	.598	.565
12	.832	.691	.781	.857	.789	.832
15	.188	.129	.161	.160	.246	.193
17	.542	.256	.456	.570	.328	.499
18	.203	.229	.214	.353	.334	.345
25	.836	.856	.844	.809	.821	.815
26	.565	.587	.575	.499	.533	.515
Avg	.569	.514	.550	.604	.539	.583

Multi-Layer Feature

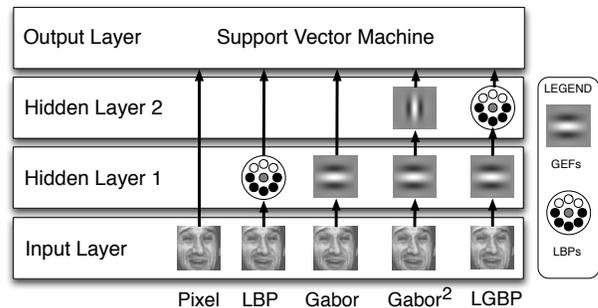


Fig. 8. A list of features from single layer to three layers.

that the winning system used an architecture, named LGBP, consisting of two hidden layers rather than the single hidden-layer architecture of our system. Post-competition we focused on comparing the performance of single layer *vs.* two layer architectures. Here we present the results of this analysis.

Figure 8 displays the different architectures we investigated: The output layer of all the systems was a SVM. The simplest architecture used no hidden layers, *i.e.*, the classifier was applied directly to the image pixels. The single hidden layer architectures used either LBPs or GEFs. The dual layer architectures used either LBP on top of GEFs (this is commonly known as LGBP) or GEFs on top of a first layer of GEFs. We refer to this architecture as Gabor² (G^2).

As these analyses were done post competition, the performance evaluation was based on leave-one-subject out double cross-validation methods on the FERA training set rather than the testing set. Therefore, the numbers are not directly comparable to those in the previous section. The performance metric was 2AFC averaged over all classified AUs. To speed up the training process, only those frames selected by FS36 were used for training.

In addition to FERA dataset, we also evaluated the different architectures on the FFD07 database using 3-fold subject independent cross validation.

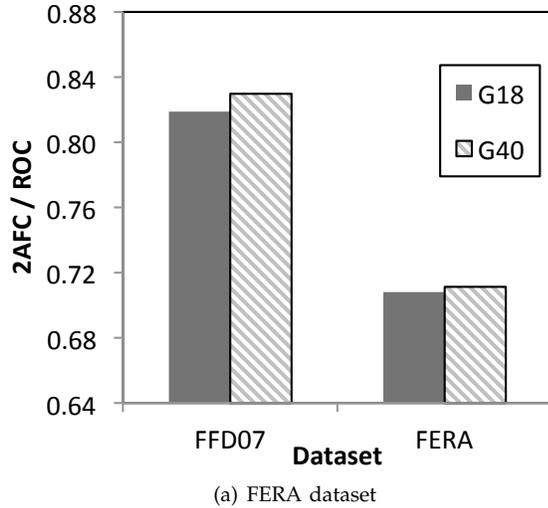


Fig. 9. Optimizing Gabor filter bank configuration: G40 (5 freq. \times 8 orientations) or G18 (3 freq. \times 6 orientations)

4.1 Optimizing Gabor, LBP and LGBP Features

To perform a fair comparison of the different architectures we first optimized their parameters. For GEFs, we evaluated the following settings²:

G40 - 5 frequency (0.5, 1, 2, 4, 8 cycles/degree) and 8 orientations ($\frac{k\pi}{8}, k = 0, 1, \dots, 7$). Figure 3 plots the spectrum of these filters.

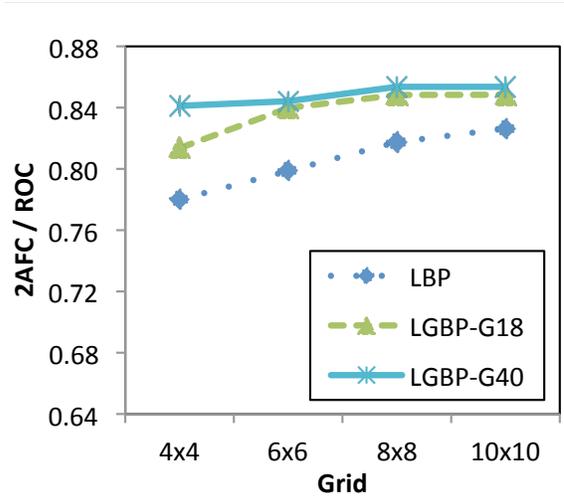
G18 - 3 frequency (2, 4, 8 cycles/degree) and 6 orientations ($0, \frac{\pi}{6}, \dots, \frac{5\pi}{6}$).

Figure 9 shows the performance of these two settings in both datasets. Consistent with our prior experience, the results suggested that G40 was better than G18.

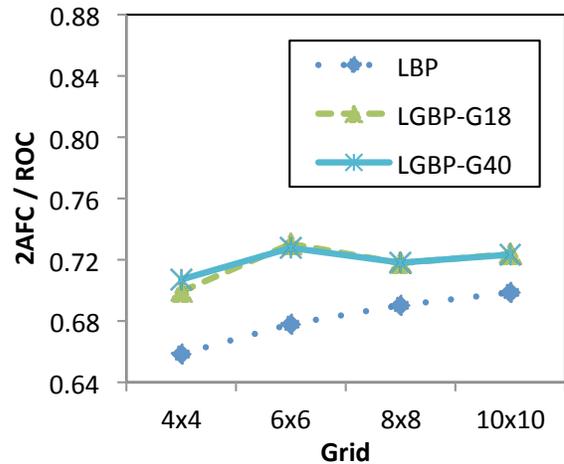
Regarding LBP/LGBP, an important parameter is the size and number of image regions over which histograms are computed. The prior literature seems to be somewhat divided on this issue. In [36] the authors partitioned the face into 16 regions organized as a 4×4 grid. In [39] the authors used a 7×6 grid and [31] used an 8×8 grid. In our analysis we tried 4×4 , 6×6 , 8×8 and 10×10 grids. Figure 10 shows the performance of LBP and LGBP on the two different datasets as a function of the grid size. In general, it appears that LBP benefited from using finer grid sizes. LGBP seemed to perform approximately the same regardless of the grid size. For further experiments, we used the 10×10 option. Next, we optimized the radius for LBP. The result is shown in Fig. 11. We found that for both LBP and LGBP, the performance deteriorated as radius increased though the gap between the smaller three radii were minimal. For further experiments, we used radius 1 pixel.

LGBP uses a first hidden layer of Gabor filters. We found that overall LGBP performed slightly better with the G40 filter bank than with the G18 filter bank. Therefore for the rest of the experiments LGBP was used with

². We are using the standard that a face patch spans 4×4 degrees of visual angle.



(a) FFD07 dataset



(b) FERA dataset

Fig. 10. Optimizing grid configuration and Gabor filters (G40 or G18) in LGBP and LBP

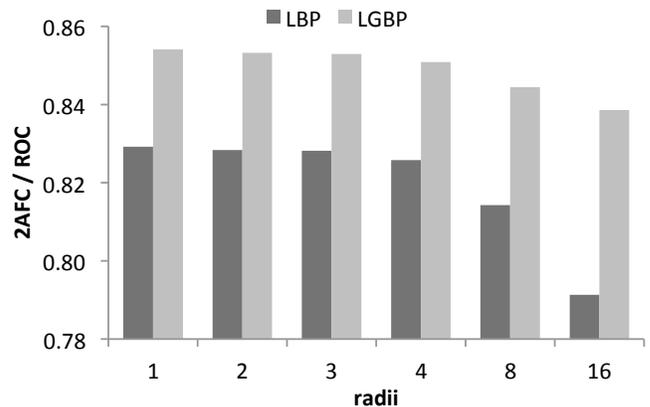


Fig. 11. Performance of LBP and LGBP for varying radii (pixels)

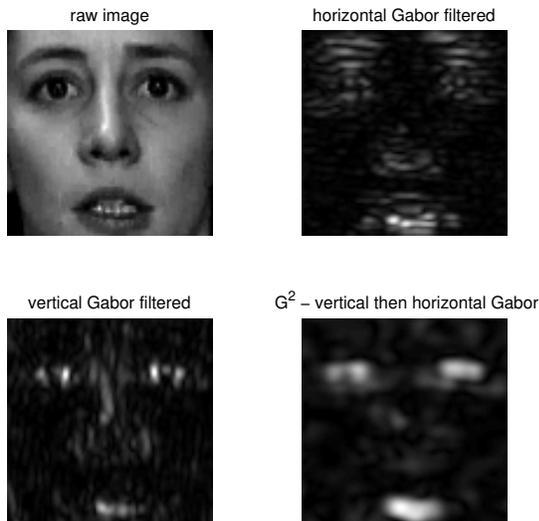


Fig. 12. Examples of Gabor energy and G^2 filtered face.

G40 filters in the first layer and 10×10 blocks in the second layer.

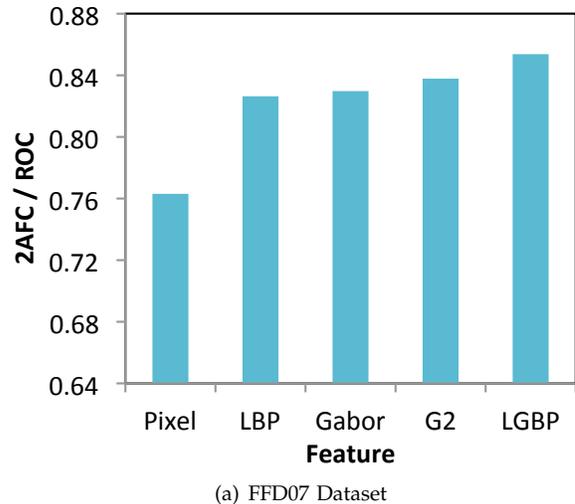
4.2 G^2 Architecture

Viola [8] pioneered the use of 2 layers of GEFs for representing image textures. The idea was that such representation could encode image textures that go beyond edges and bars. For example it could encode the fact that a horizontal structure is made of tiny vertical bars. Fig. 12 shows an example of such structure. The contrast between the sclera and the iris creates vertical bars that excite vertically oriented Gabor filters. The edges between the teeth also excite the vertically oriented Gabors. A horizontally tuned Gabor applied on top of the vertically tuned Gabor then captures the fact that the teeth align horizontally. A similar phenomenon occurs with the eyes. The final result is a representation in which the teeth and eyes clearly stand out. Note as Gabor energy filters are nonlinear, the two layers of filters cannot be combined by using a single convolution operation.

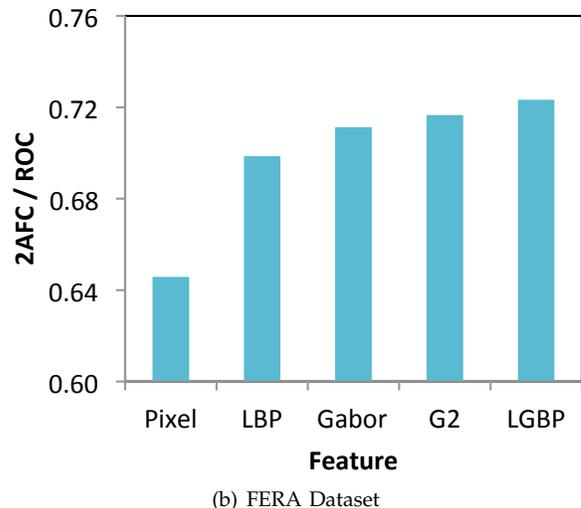
The G^2 architecture we tested used the G18 filter bank on top of the G40 filter bank. No further attempts were made to optimize this architecture.

4.3 Results

Figure 13 shows the average performance of the different architectures on the FERA and FFD07 datasets. Results are presented in terms of 2AFC score using a cross-validation procedure within each dataset. In all cases the classifier was an SVM. We selected the best known kernel for each the features. For LBP and LGBP, the histogram intersection kernel was used as in [36]. For Gabor and G^2 , since linear and RBF kernel performed comparable in our experience and literature [37], so we chose the linear kernel for simplicity. For raw pixel, we used linear kernel for comparison.



(a) FFD07 Dataset



(b) FERA Dataset

Fig. 13. Cross-Validation Performance of Features on FFD07 and FERA datasets

Figure 13 shows that the architectures with two hidden layers (LGBP and G^2) performed best, followed by one hidden layer architectures (Gabor and LBP) and lastly the architecture with no hidden layer. This trend is consistent in both datasets.

4.4 Resistance to Illumination Noise and Jitter

The original LGBP paper [47] suggested that LGBP performs better than Gabor and LBP due to the fact that it is less sensitive to local changes in illumination. Gabor energy filters are also supposed to be advantageous because of their relative insensitivity to local illumination and spatial jitter. To evaluate this hypothesis we tested how the different architectures degrade as a function of jitter noise and illumination noise. The noise was added into both training and testing sets. Our expectation was that the architectures with 2 hidden layers would be the most resistant to noise, followed by 1 hidden layer and no hidden layer architectures.



Fig. 14. Examples of simulated facebox jitter



Fig. 15. Examples of illumination noise of varying strength

To simulate jitter, we added noise to the registered face box location. The noise was sampled randomly from a 2D uniform distribution with support $[-9, 9] \times [-9, 9]$ pixels in the horizontal and vertical directions (see Figure 14).

We also tested resistance to local illumination noise by adding Gaussian “spotlights” at random positions on the face (see Fig. 15). The spotlights were generated from the weighted average between the original image and a plain-white image using a 2D isotropic Gaussian envelope with random centers and standard deviation $\text{diag}([10, 10])$. In addition to simulated illumination variation, one could also use datasets with real varied illumination, such as the MultiPIE dataset [18]). We leave this as future work.

Figure 16 presents the cross-validation results on FFD07 database. As expected, the simplest architecture

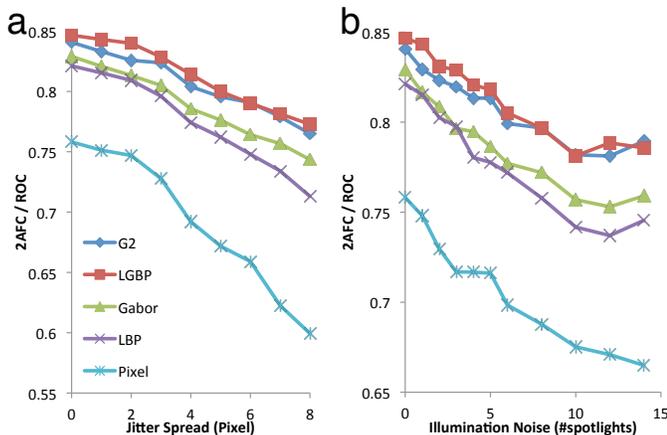


Fig. 16. Performance of features with (a) varying spreads of face registration jitter, and (b) varying intensities of illumination change

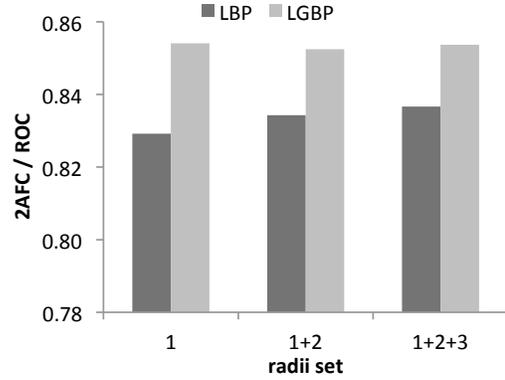


Fig. 17. Performance of LBP and LGBP for varying radii sets.

with no hidden layers was more vulnerable to jitter and illumination noise than those with hidden layers. Among the single hidden layer architectures, Gabor decayed slower than LBP in both types of noise. This suggested that Gabor features are a good choice for first hidden layer. Best overall performance was obtained by the two hidden layer architectures (LGBP and G^2). The extra layer did improve the performance substantially. The second layer seemed to marginally improve the resistance of LBP to jitter and illumination noise. This can be seen by the fact that the LBP and LGBP curves tend to diverge as noise increases. However the second layer did not seem to improve resistance to noise over the single layer Gabor architecture. This can be seen by the fact that the curves of single layer Gabor and of G^2 are parallel.

4.5 Multi-resolution LBP

We wondered why Gabor filters performed consistently better in single layer architectures than LBPs, particularly considering that Gabor and LBP are both basically oriented edge detectors. One difference is that Gabor filter banks typically utilize a range of peak spatial frequencies (Fig. 3), thus encoding multiple scales and resolutions. LBP on the other hand are typically applied with a single radius. Similarly, the advantage of LGBP over LBP could be partially attributed to the multi-scale power provided by the first Gabor layer.

To explore these issues, we tested multi-resolution LBPs with radii 1, 1+2, and 1+2+3 pixels. Our implementation of multi-resolution LBP simply concatenated all the LBP feature vectors not unlike what we do with the Gabor architecture. Figure 17 shows the results of the study. As expected, the performance of multi-resolution LBP increased as the LBP feature pool increased. Surprisingly, the last set, 1+2+3, even slightly surpassed the single layer Gabor system, although the score was still lower than the two layer architectures, G^2 and LGBP. However, as Figure 17 shows, adding multiple resolutions to the second layer of LBPs did not seem to

help. This may be due to the fact that the Gabor layer in LGBP is already providing multi-resolution capabilities and that therefore adding these to the second layer is redundant.

4.6 Generalization across datasets

Table 6 shows the performance of the 4 different architectures when tested for generalization within and across datasets. When trained on FFD07 and tested on FERA, the performance of LGBP, G^2 and single layer Gabors was similar but LBP fell a bit behind. This suggests Gabor filters provide the best representation for the first hidden layer.

When trained on FERA and tested on FFD07, LGBP was the clear winner, followed by Gabor and then G^2 . We suspect the relatively poor performance of G^2 when generalizing from a small to a larger dataset may be due to the fact that the version of G^2 we investigated has a very large dimensionality (6635520) in comparison to single layer Gabors (368640) and LGBP (236000). Table 5 describes how the number of features are calculated.

TABLE 5

Comparing the number of feature of different feature types

Feature	Configuration	#Features
LGBP	59 bins * 100 blocks * 40	236000
Gabor	96*96 pixels * 40	368640
G^2	96*96 pixels * 18 * 40	6635520

5 CONCLUSIONS

We compared single layer and dual layer architectures for texture based expression recognition systems. We showed that dual layer architectures provide small but consistent improvements in performance over single layers systems. We show that for the first layer GEFs outperformed LBPs. Experiments suggested that this may be due to the fact that Gabor filters are more resistant than LBPs to noise caused by jitter or changes in illumination. However for the second layer, LBP filters appear to perform marginally better than Gabor filters. We suspect that this may be due to the fact that the Gabor filter approaches we tried had many more parameters than the LBP approaches. Careful selection of the Gabor features in the second layer may make a difference. More work is needed to gain a better understanding of why 2 layer architectures perform better and why Gabor filters are preferable for the first layer while LBP filters are preferable for the second layer.

ACKNOWLEDGMENT

We thank the editor and anonymous reviewers. Support for this work was provided by NSF grants IIS-0905622, IIS-1002840-SoCS and NSF IIS-INT2-0808767. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do

TABLE 6

Evaluation of feature across FERA and FFD07 databases. The standard error suggests the variability of the result across different AUs.

(a) LGBP

		test on	
		FERA	FFD07
train on	FERA	.723±.032	.698±.028
	FFD07	.718±.030	.854±.017

(b) G^2

		test on	
		FERA	FFD07
train on	FERA	.717±.028	.650±.027
	FFD07	.721±.020	.838±.020

(c) Gabor

		test on	
		FERA	FFD07
train on	FERA	.711±.029	.668±.021
	FFD07	.723±.021	.830±.022

(d) LBP

		test on	
		FERA	FFD07
train on	FERA	.699±.029	.629±.041
	FFD07	.685±.028	.828±.023

(e) Pixel

		test on	
		FERA	FFD07
train on	FERA	.646±.024	.633±.029
	FFD07	.681±.027	.763±.028

not necessarily reflect the views of the National Science Foundation.

APPENDIX

Pseudo-code for computing 2AFC score

```
% x is a vector of real valued numbers
% y is a vector of 0s and 1s
function s = Calc2AFC(x, y)
    x0 = x(y==0);
    x1 = x(y==1);
    n0 = length(x0);
    n1 = length(x1);
    s = 0;
    for k = 1:n0
        n = sum(x1>x0(k)) + 0.5*sum(x1==x0(k));
        s = s + n / (n1*n0);
    end
```

REFERENCES

- [1] A. Ashraf, S. Lucey, J. Cohn, T. Chen, Z. Ambadar, K. Prkachin, and P. Solomon. The painful face-pain expression recognition using active appearance models. *Image and Vision Computing*, 27(12):1788–1796, 2009.
- [2] T. Bañziger and K. R. Scherer. Introducing the geneva multimodal emotion portrayal (gemep) corpus. K. R. Scherer, T. Bañziger, and E. B. Roesch, editors, *Blueprint for Affective Computing: A Sourcebook, Series in affective science*, pages chapter 6.1, 271–294, 2010.

- [3] M. S. Bartlett, G. Littlewort, C. Lainscsek, I. Fasel, and J. Movellan. Recognition of facial actions in spontaneous expressions. *Journal of Multimedia*, 2006.
- [4] Bartlett, M, G. Donato, Movellan, J. R., and J. Hager. Face image analysis for expression measurement and detection of deceit. In *Proceedings of the 6th Symposium on Neural Computation*, pages 8–16. California Institute of Technology, 1999.
- [5] C. Burges. A tutorial on support vector machines for pattern recognition. *Data mining and knowledge discovery*, 2(2):121–167, 1998.
- [6] C.-C. Chang and C.-J. Lin. *LIBSVM: a library for support vector machines*, 2001. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- [7] D. C. Cireřan, U. Meier, J. Masci, L. M. Gambardella, and J. Schmidhuber. High-Performance Neural Networks for Visual Object Classification. Feb. 2011.
- [8] J. De Boner and P. Viola. Structure driven image database retrieval. *Advances in Neural Information Processing Systems*, pages 866–872, 1998.
- [9] F. De la Torre and J. F. Cohn. *Guide to Visual Analysis of Humans: Looking at People*, chapter Facial Expression Analysis. Springer, 2011.
- [10] G. Donato, M. Bartlett, J. Hager, P. Ekman, and T. Sejnowski. Classifying facial actions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 21(10):974–989, 1999.
- [11] M. Eckhardt, I. Fasel, and J. Movellan. Towards practical facial feature detection. *International Journal of Pattern Recognition and Artificial Intelligence*, 23(3):379–400, 2009.
- [12] P. Ekman and W. Friesen. *Facial Action Coding System (FACS): A technique for the measurement of facial action*. Consulting, Palo Alto, CA, 1978.
- [13] M. Frank, M. Bartlett, and J. Movellan. The M3 database of spontaneous emotion expression (University of Buffalo). *In prep*, 2010.
- [14] Y. Freund and R. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. In *Computational learning theory*, pages 23–37. Springer, 1995.
- [15] J. Friedman, T. Hastie, and R. Tibshirani. Additive logistic regression: a statistical view of boosting. *Annals of Statistics*, 28(2), 2000.
- [16] C. Goodall. Procrustes methods in the statistical analysis of shape. *Journal of the Royal Statistical Society. Series B (Methodological)*, 53(2):285–339, 1991.
- [17] D. Green and J. Swets. Signal detection theory and psychophysics. 1966.
- [18] R. Gross, I. Matthews, J. Cohn, T. Kanade, and S. Baker. Multi-pie. *Image and Vision Computing*, 28(5):807–813, 2010.
- [19] G. Hinton, S. Osindero, and Y. Teh. A fast learning algorithm for deep belief nets. *Neural computation*, 18(7):1527–1554, 2006.
- [20] T. Kanade, J. Cohn, and Y. L. Tian. Comprehensive database for facial expression analysis. In *IEEE International Conference on Automatic Face & Gesture Recognition*, 2000.
- [21] S. Keerthi and C. Lin. Asymptotic behaviors of support vector machines with Gaussian kernel. *Neural computation*, 15(7):1667–1689, 2003.
- [22] Y. LeCun, B. Boser, J. Denker, D. Henderson, R. Howard, W. Hubbard, and L. Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989.
- [23] G. Littlewort, M. Bartlett, I. Fasel, J. Susskind, and J. Movellan. Dynamics of facial expression extracted automatically from video. *Image and Vision Computing*, 24(6):615–625, 2006.
- [24] G. Littlewort, M. Bartlett, I. Fasel, J. Susskind, and J. R. Movellan. An automatic system for measuring facial expression in video. *Image and Vision Computing*, in press.
- [25] G. Littlewort, J. Whitehill, T. Wu, I. Fasel, M. Frank, J. Movellan, and M. Bartlett. The computer expression recognition toolbox (CERT). In *Proceedings of Automatic Face and Gesture Recognition*, 2011.
- [26] P. Lucey, J. Cohn, T. Kanade, J. Saragih, and Z. Ambadar. The extended cohn-kanade dataset (CK+): A complete dataset for action unit and emotion-specified expression. In *CVPR Workshop on Human Communicative Behavior Analysis*, 2010.
- [27] P. Lucey, S. Lucey, and J. Cohn. Registration invariant representations for expression detection. In *2010 International Conference on Digital Image Computing: Techniques and Applications*, pages 255–261. IEEE, 2010.
- [28] T. Marks, J. R. Hershey, and J. R. Movellan. Tracking motion, deformation, and texture using conditionally gaussian processes. *Pattern Analysis and Machine Intelligence*, 32(2):348–363, February 2010.
- [29] S. Mason and N. Graham. Areas beneath the relative operating characteristics (roc) and relative operating levels (rol) curves: Statistical significance and interpretation. *Quarterly Journal of the Royal Meteorological Society*, 128(584):2145–2166, 2002.
- [30] P. Michel and R. El Kaliouby. Real time facial expression recognition in video using support vector machines. In *Proceedings of the 5th international conference on Multimodal interfaces*, pages 258–264. ACM, 2003.
- [31] S. Moore, R. Bowden, and U. Guildford. The effects of pose on facial expression recognition. In *In Proceedings of British Machine Vision Conference (BMVC2009), London, UK*. Citeseer, 2009.
- [32] J. R. Movellan. Optimal control of dynamic Bayes nets. Technical report, MPLab Tutorials, UCSD, 2005.
- [33] T. Ojala, M. Pietikäinen, and T. Mäenpää. Multiresolution gray-scale and rotation invariant texture classification with local binary patterns. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(7):971–987, July 2002.
- [34] M. Pantic, M. Valstar, R. Rademaker, and L. Maat. Web-based database for facial expression analysis. In *International Conference on Multimedia and Expo*, 2005.
- [35] A. Ryan, J. Cohn, S. Lucey, J. Saragih, P. Lucey, F. De la Torre, and A. Rossi. Automated facial expression recognition system. In *Security Technology, 2009. 43rd Annual 2009 International Carnahan Conference on*, pages 172–177. IEEE, 2009.
- [36] T. Senechal, V. Rapp, H. Salam, R. Seguier, K. Bailly, and L. Prevost. Combining aam coefficients with lgbp histograms in the multi-kernel svm framework to detect facial action units. In *Automatic Face & Gesture Recognition and Workshops (FG 2011), 2011 IEEE International Conference on*, pages 860–865. IEEE.
- [37] C. Shan, S. Gong, and P. McOwan. Facial expression recognition based on Local Binary Patterns: A comprehensive study. *Image and Vision Computing*, 27(6):803–816, 2009.
- [38] T. Simon, M. Nguyen, F. De La Torre, and J. Cohn. Action unit detection with segment-based svms. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 2737–2744. IEEE.
- [39] X. Sun, H. Xu, C. Zhao, and J. Yang. Facial expression recognition based on histogram sequence of local gabor binary patterns. In *Cybernetics and Intelligent Systems, 2008 IEEE Conference on*, pages 158–163. IEEE, 2008.
- [40] J. Susskind, G. Hinton, J. Movellan, and A. Anderson. Generating facial expressions with deep belief nets. *Affective Computing, Emotion Modelling, Synthesis and Recognition*, pages 421–440.
- [41] M. F. Valstar, B. Jiang, M. Méhu, M. Pantic, and K. Scherer. The first facial expression recognition and analysis challenge. In *Proc. IEEE Intl Conf. Automatic Face and Gesture Recognition*, in print, 2011.
- [42] M. F. Valstar and M. Pantic. Fully automatic recognition of the temporal phases of facial actions. *Part B: Cybernetics, IEEE Transactions on Systems, Man, and Cybernetics*, PP(99):1–16, 2011.
- [43] V. Vapnik. *Statistical learning theory*. 1998.
- [44] P. Viola and M. Jones. Robust real-time face detection. *International Journal of Computer Vision*, 2004.
- [45] J. Whitehill, G. Littlewort, I. Fasel, M. Bartlett, and J. R. Movellan. Toward practical smile detection. *Transactions on Pattern Analysis and Machine Intelligence*, (11):2106–2111, 2009.
- [46] T. Wu, N. Butko, P. Ruvolo, J. Whitehill, M. Bartlett, and J. Movellan. Action unit recognition transfer across datasets. In *Automatic Face & Gesture Recognition and Workshops (FG 2011), 2011 IEEE International Conference on*, pages 889–896. IEEE, 2011.
- [47] W. Zhang, S. Shan, W. Gao, X. Chen, and H. Zhang. Local gabor binary pattern histogram sequence (lgbphs): A novel non-statistical model for face representation and recognition. 2005.
- [48] G. Zhao and M. Pietikainen. Dynamic texture recognition using local binary patterns with an application to facial expressions. *IEEE transactions on pattern analysis and machine intelligence*, pages 915–928, 2007.