



## Framework for reliable, real-time facial expression recognition for low resolution images

Rizwan Ahmed Khan <sup>a,b,\*</sup>, Alexandre Meyer <sup>a,b</sup>, Hubert Konik <sup>a,c</sup>, Saïda Bouakaz <sup>a,b</sup>

<sup>a</sup> Université de Lyon, CNRS, France

<sup>b</sup> Université Lyon 1, LIRIS, UMR5205, F-69622, France

<sup>c</sup> Université Jean Monnet, Laboratoire Hubert Curien, UMR5516, 42000 Saint-Etienne, France

### ARTICLE INFO

#### Article history:

Received 30 July 2012

Available online 4 April 2013

Communicated by S. Sarkar

#### Keywords:

Facial expression recognition

Low resolution images

Local binary pattern

Image pyramid

Salient facial regions

### ABSTRACT

Automatic recognition of facial expressions is a challenging problem specially for low spatial resolution facial images. It has many potential applications in human-computer interactions, social robots, deceit detection, interactive video and behavior monitoring. In this study we present a novel framework that can recognize facial expressions very efficiently and with high accuracy even for very low resolution facial images. The proposed framework is memory and time efficient as it extracts texture features in a pyramidal fashion only from the perceptual salient regions of the face. We tested the framework on different databases, which includes Cohn-Kanade (CK+) posed facial expression database, spontaneous expressions of MMI facial expression database and FG-NET facial expressions and emotions database (FEED) and obtained very good results. Moreover, our proposed framework exceeds state-of-the-art methods for expression recognition on low resolution images.

© 2013 Elsevier B.V. All rights reserved.

### 1. Introduction

Communication in any form i.e. verbal or non-verbal is vital to complete various routine tasks and plays a significant role in daily life. Facial expression is the most effective form of non-verbal communication and it provides a clue about emotional state, mindset and intention (Ekman, 2001). Human visual system (HVS) decodes and analyzes facial expressions in real time despite having limited neural resources. As an explanation for such performance, it has been proposed that only some visual inputs are selected by considering “salient regions” (Zhaoping, 2006), where “salient” means most noticeable or most important.

For computer vision community it is a difficult task to automatically recognize facial expressions in real-time with high reliability. Variability in pose, illumination and the way people show expressions across cultures are some of the parameters that make this task difficult. Low resolution input images makes this task even harder. Smart meeting, video conferencing and visual surveillance are some of the real world applications that require facial expression recognition system that works adequately on low resolution images. Another problem that hinders the development of such system for real world application is the lack of databases with

natural displays of expressions (Valstar and Pantic, 2010). There are number of publicly available benchmark databases with posed displays of the six basic emotions (Ekman, 1971) exist but there is no equivalent of this for spontaneous basic emotions. While, it has been proved that Spontaneous facial expressions differ substantially from posed expressions (Bartlett et al., 2002). In this work, we propose a facial expression recognition system that caters for illumination changes and works equally well for low resolution as well as for good quality/high resolution images. We have tested our proposed system on spontaneous facial expressions as well and recorded encouraging results.

We propose a novel descriptor for facial features analysis, pyramid of local binary pattern (PLBP) (refer Section 3). PLBP is a spatial representation of local binary pattern (LBP) (Ojala et al., 1996) and it represents stimuli by its local texture (LBP) and the spatial layout of the texture. We combined pyramidal approach with LBP descriptor for facial feature analysis as this approach has already been proved to be very effective in a variety of image processing tasks (Hadjidemetriou et al., 2004). Thus, the proposed descriptor is a simple and computationally efficient extension of LBP image representation, and it shows significantly improved performance for facial expression recognition tasks for low resolution images. We base our framework for automatic facial expression recognition (FER) on human visual system (HVS) (refer Section 5), so it extracts PLBP features only from the salient regions of the face. To determine which facial region(s) are the most important or salient according to HVS, we conducted a psycho-visual

\* Corresponding author at: Université Lyon 1, LIRIS, UMR5205, F-69622, France. Tel.: +33 (0) 4 72 43 19 75; fax: +33 (0) 4 72 43 15 36.

E-mail addresses: [Rizwan-ahmed.khan@liris.cnrs.fr](mailto:Rizwan-ahmed.khan@liris.cnrs.fr) (R.A. Khan), [Alexandre.meyer@liris.cnrs.fr](mailto:Alexandre.meyer@liris.cnrs.fr) (A. Meyer), [Hubert.Konik@univ-st-etienne.fr](mailto:Hubert.Konik@univ-st-etienne.fr) (H. Konik), [Saida.bouakaz@liris.cnrs.fr](mailto:Saida.bouakaz@liris.cnrs.fr) (S. Bouakaz).

experiment using an eye-tracker (refer Section 4). We considered six universal facial expressions for psycho-visual experimental study as these expressions are proved to be consistent across cultures (Ekman, 1971). These six expressions are anger, disgust, fear, happiness, sadness and surprise. The novelty of the proposed framework is that, it is illumination invariant, reliable on low resolution images and works adequately for both i.e. posed and spontaneous expressions.

Generally, facial expression recognition system consists of three steps: face detection, feature extraction and expression classification. The same has been shown in Fig. 1. In our framework we tracked face/salient facial regions using Viola-Jones object detection algorithm (Viola and Jones, 2001) as it is the most cited and considered the fastest and most accurate pattern recognition method for face detection (Kolsch and Turk, 2004). The second step in the framework is feature extraction, which is the area where this study contributes. The optimal features should minimize within-class variations of expressions, while maximize between class variations. If inadequate features are used, even the best classifier could fail to achieve accurate recognition (Shan et al., 2009). Section 3 presents the novel method for facial features extraction which is based on human visual system (HVS). To study and understand HVS we performed psycho-visual experiment. Psycho-visual experimental study is briefly described in Section 4. Expression classification or recognition is the last step in the pipeline. In literature two different ways are prevalent to recognize expressions i.e. direct recognition of prototypic expressions or recognition of expressions through facial action coding system (FACS) action units (AUs) (Ekman and Friesen, 1978). In our proposed framework, which is described in Section 5 we directly classify six universal prototypic expressions (Ekman, 1971). The performance of the framework is evaluated for five different classifiers (from different families i.e. classification tree, instance based learning, SVM, etc.) and results are presented in Section 6. Next section presents the brief literature review for facial features extraction methods.

## 2. Related work

In the literature, various methods are employed to extract facial features and these methods can be categorized either as appearance-based methods or geometric feature-based methods.

*Appearance-based methods.* One of the widely studied method to extract appearance information is based on Gabor wavelets (Littlewort et al., 2006; Tian, 2004; Donato et al., 1999). Generally, the drawback of using Gabor filters is that it produces extremely large number of features and it is both time and memory intensive to convolve face images with a bank of Gabor filters to extract multi-scale and multi-orientational coefficients. Another promising approach to extract appearance information is by using Haar-like features, see Yang et al. (2010). Recently, texture descriptors and classification methods i.e. local binary pattern (LBP) (Ojala et al., 1996) and local phase quantization (LPQ) (Ojansivu and Heikkilä,

2008) are also studied to extract appearance-based facial features. Zhao and Pietikäinen (2007) proposed to model texture using volume local binary patterns (VLBP) an extension to LBP, for expression recognition.

*Geometric-based methods.* Geometric feature-based methods (Zhang and Ji, 2005; Pantic and Patras, 2006; Valstar et al., 2005; Bai et al., 2009) extracts shapes and locations of facial components information to form a feature vector. The problem with using geometric feature-based methods is that they usually require accurate and reliable facial feature detection and tracking which is difficult to achieve in many real world applications where illumination changes with time and images are recorded in very low resolution.

Generally, we have found that all the reviewed methods for automatic facial expression recognition are computationally expensive and usually requires dimensionally large feature vector to complete the task. This explains their inability for real-time applications. Secondly, in literature, very few studies exist that tackles the issue of expressions recognition from low resolution images, this adds to lack of applicability of expression recognition system for real world applications. Lastly, all of the reviewed methods, spend computational time on whole face image or divides the facial image based on some mathematical or geometrical heuristic for features extraction. We argue that the task of expression analysis and recognition could be done in more conducive manner, if only some regions are selected for further processing (i.e. salient regions) as it happens in human visual system. Thus, our contributions in this study are:

1. We propose a novel descriptor for facial expression analysis i.e. pyramid of local binary pattern (PLBP), which outperforms state-of-the-art methods for expression recognition on low resolution images (spatially degraded images). It also performs better than other state-of-the-art methods for good resolution images (with no degradation).
2. As the proposed framework is based on human visual system it algorithmically processes only salient facial regions which reduces the length of feature vector. This reduction in feature vector length makes the proposed framework suitable for real-time applications due to minimized computational complexity.

## 3. Pyramid of local binary pattern

The proposed framework creates a novel feature space by extracting proposed PLBP (pyramid of local binary pattern) features only from the visually salient facial region (see Section 4 for psycho-visual experiment). PLBP is a *pyramidal-based spatial* representation of local binary pattern (LBP) descriptor. PLBP represents stimuli by their local texture (LBP) and the spatial layout of the texture. The spatial layout is acquired by tiling the image into regions at multiple resolutions. The idea is illustrated in Fig. 2. If only the coarsest level is used, then the descriptor reduces to a global LBP histogram. Comparing to the multi-resolution LBP of

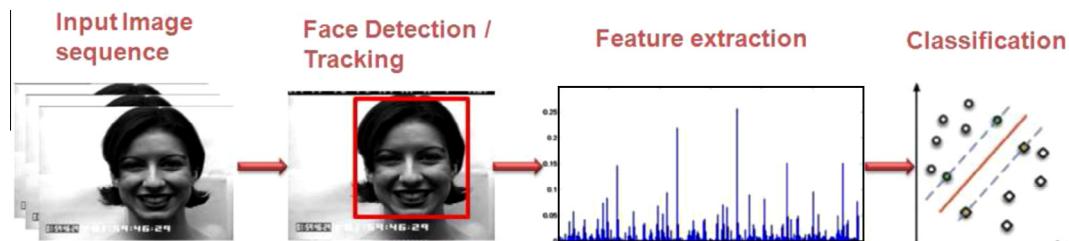
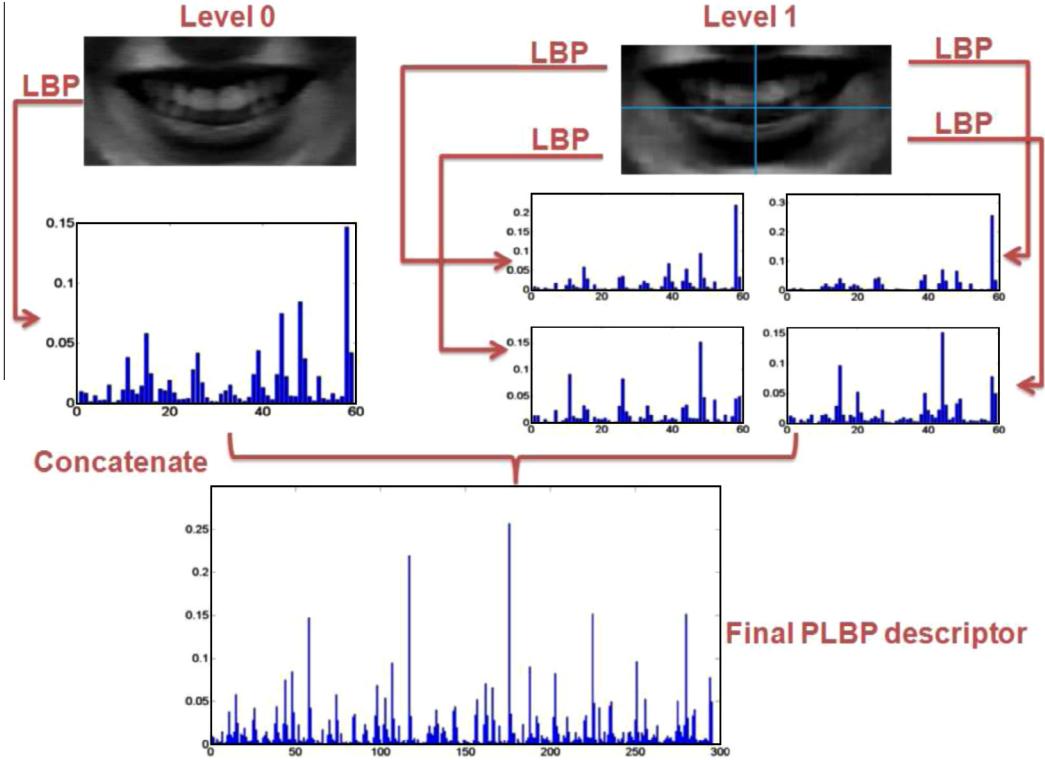


Fig. 1. Basic structure of facial expression recognition system pipeline.



**Fig. 2.** Pyramid of local binary pattern. First row: stimuli at two different pyramid levels, second row: histograms of LBP at two respective levels, and third row: final descriptor.

Ojala et al. (2002), our descriptor selects samples in a more uniformly distributed manner, whereas Ojala's LBP takes samples centered around a point leading to missing some information in the case of face (which is different than a repetitive texture).

LBP features were initially proposed for texture analysis (Ojala et al., 1996), but recently they have been successfully used for facial expression analysis (Zhao and Pietikäinen, 2007; Shan et al., 2009). The most important property of LBP features are their tolerance against illumination changes and their computational simplicity (Ojala and Pietikäinen, 1999; Ojala et al., 1996, 2002). The operator labels the pixels of an image by thresholding the  $3 \times 3$  neighborhood of each pixel with the center value and considering the result as a binary number. Then the histogram of the labels can be used as a texture descriptor. Formally, LBP operator takes the form:

$$LBP(x_c, y_c) = \sum_{n=0}^7 s(i_n - i_c) 2^n \quad (1)$$

where in this case  $n$  runs over the eight neighbors of the central pixel  $c$ ,  $i_c$  and  $i_n$  are the gray level values at  $c$  and  $n$  and  $s(u)$  is 1 if  $u \geq 0$  or 0 otherwise.

Later, the LBP operator is extended to use neighborhood of different sizes (Ojala et al., 2002) as the original operator uses  $3 \times 3$  neighborhood. Using circular neighborhoods and bilinearly interpolating the pixel values allow any radius and number of pixels in the neighborhood. The LBP operator with  $P$  sampling points on a circular neighborhood of radius  $R$  is given by:

$$LBP_{P,R} = \sum_{p=0}^{P-1} s(g_p - g_c) 2^p \quad (2)$$

where,  $g_c$  is the gray value of the central pixel,  $g_p$  is the value of its neighbors,  $P$  is the total number of involved neighbors and  $R$  is the radius of the neighborhood.

Another extension to the original operator is the definition of *uniform patterns*, which can be used to reduce the length of the feature vector and implement a simple rotation-invariant descriptor. A local binary pattern is called uniform if the binary pattern contains at most two bitwise transitions from 0 to 1 or vice versa when the bit pattern is traversed circularly. Accumulating the patterns which have more than two transitions into a single bin yields an LBP operator, denoted  $LBP_{P,R}^{u_2}$  patterns. These binary patterns can be used to represent texture primitives such as spot, flat area, edge and corner.

We extend LBP operator so that the stimuli can be represented by its local texture and the spatial layout of the texture. We call this extended LBP operator as pyramid of local binary pattern or PLBP. PLBP creates the spatial pyramid by dividing the stimuli into finer spatial sub-regions by iteratively doubling the number of divisions in each dimension. It can be observed from Fig. 2 that the pyramid at level  $l$  has  $2^l$  sub-regions along each dimension ( $R_0, \dots, R_{m-1}$ ). Histograms of LBP features at the same levels are concatenated. Then, their concatenation at different pyramid levels gives final PLBP descriptor (as shown in Fig. 2). It can be defined as:

$$H_{ij} = \sum_l \sum_{xy} I\{f_i(x,y) = i\} I\{(x,y) \in R_l\} \quad (3)$$

where  $l = 0, \dots, m-1$ ,  $i = 0, \dots, n-1$ .  $n$  is the number of different labels produced by the LBP operator and

$$I(A) = \begin{cases} 1 & \text{if } A \text{ is true} \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

While, the dimensionality of the descriptor can be calculated by:

$$N \sum_l 4^l \quad (5)$$

Where, in our experiment (see Section 6)  $l = 1$  and  $N = 59$  as we created pyramid up to level 1 and extracted 59 LBP features using

**Table 1**

Comparison of time and memory consumption.

	Shan et al. (2009)(a)	Shan et al. (2009)(b)	Bartlett et al. (2003)	PLBP
Memory (feature dimension)	2478	42,650	92,160	590
Time (feature extraction time)	0.03 s	30 s	–	0.01 s

$LBP_{8,2}^{u2}$  operator, which denotes a uniform LBP operator with eight sampling pixels in a local neighborhood region of radius 2. This pattern reduces the histogram from 256 to 59 bins. In our experiment we obtained 295 dimensional feature vector from one facial region i.e. mouth region (59 dimensions/sub-region), since we executed the experiment with the pyramid of level 1 (the same is shown in Fig. 2).

### 3.1. Novelty of the proposed descriptor

There exist some methods in literature that uses pyramid of LBP for different applications and they look similar to our proposed descriptor, i.e. Wang et al. (2011), Guo et al. (2010) and Moore and Bowden (2011). Our proposition is novel and there exist differences in the methodology that creates differences in the extracted information. Method for face recognition proposed in Wang et al. (2011) creates pyramid before applying LBP operator by down sampling original image i.e. scale-space representation, whereas we propose to create the spatial pyramid by dividing the stimuli into finer spatial sub-regions by iteratively doubling the number of divisions in each dimension. Secondly, our approach reduces memory consumption (do not requires to store same image in different resolutions) and is computationally more efficient. Guo et al. (2010) proposed approach for face and palmprint recognition based on multiscale LBP. Their proposed method seems similar to our method for expression recognition but how multiscale analysis is achieved deviates our approach. Approach proposed in Guo et al. (2010) achieves multiscale analysis using different values of  $P$  and  $R$ , where  $LBP(P, R)$  denotes a neighborhood of  $P$  equally spaced sampling points on a circle of radius  $R$  (discussed earlier). Same approach has been applied by Moore and Bowden (2011) for facial features analysis. Generally the drawback of using such approach is that it increases the size of the feature histogram and increases the computational cost. Moore and Bowden (2011) reports dimensionality of feature vector as high as 30,208 for multiscale face expression analysis as compared to our proposition which creates 590 dimensional feature vector (see Section 5) for the same task. We achieve the task of multiscale analysis much more efficiently than any other earlier proposed methods. By the virtue of efficient multiscale analysis our framework can be used for real time applications (see Table 1 for the time and memory consumption comparison) which is not the case with other methods.

As mentioned earlier, we base our framework for facial expression recognition on human visual system (HVS), which selects only few facial regions (salient) to extract information. In order to determine the saliency of facial region(s) for a particular expression, we conducted psycho-visual experiment with the help of an eye-tracker. Next section briefly explains the psycho-visual experimental study.

## 4. Psycho-visual experiment

The aim of our experiment was to record the eye movement data of human observers in free viewing conditions. The data were analyzed in order to find which components of face are salient for specific displayed expression.

### 4.1. Participants, apparatus and stimuli

Eye movements of fifteen human observers were recorded using video based eye-tracker (EyelinkII system, SR Research), as

the subjects watched the collection of 54 videos selected from the extended Cohn-Kanade (CK+) database (Lucey et al., 2010), showing one of the six universal facial expressions (Ekman, 1971). Observers include both male and female aging from 20 to 45 years with normal or corrected to normal vision. All the observers were naïve to the purpose of an experiment.

### 4.2. Eye movement recording

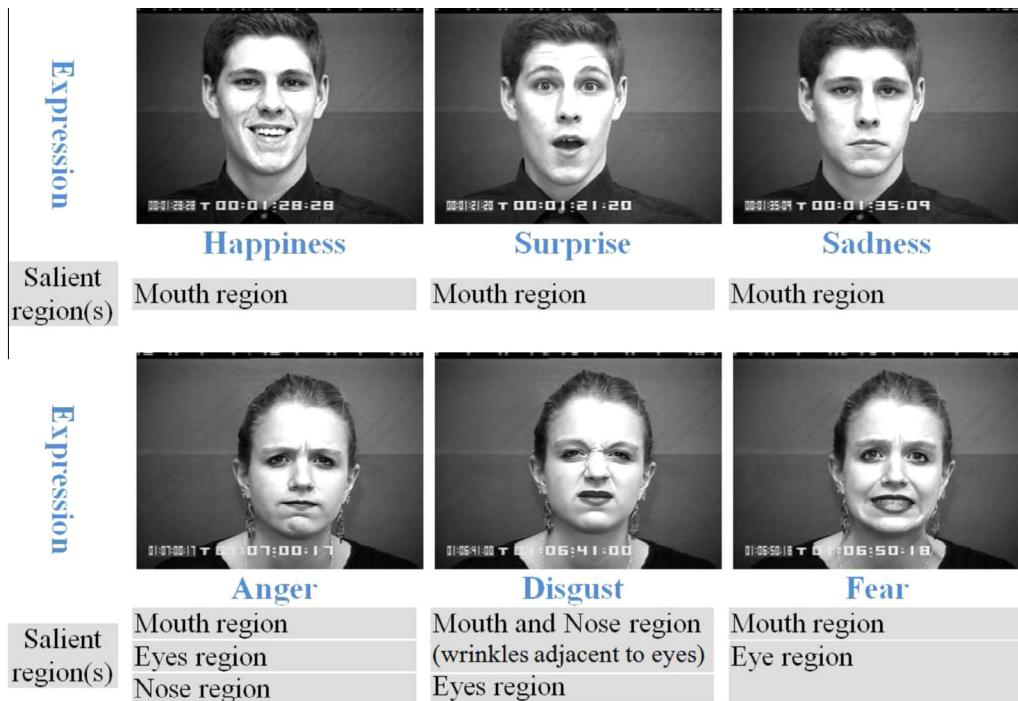
Eye position was tracked at 500 Hz with an average noise less than 0.01°. Head mounted eye-tracker allows flexibility to perform the experiment in free viewing conditions as the system is designed to compensate for small head movements.

### 4.3. Psycho-visual experiment results

In order to statistically quantify which region is perceptually more attractive for specific expression, we have calculated the average percentage of trial time observers have fixated their gazes at specific region(s) in a particular time period. As the stimuli used for the experiment is dynamic i.e. video sequences, it would have been incorrect to average all the fixations recorded during trial time (run length of the video) for the data analysis as this could lead to biased analysis of the data. To meaningfully observe and analyze the gaze trend across one video sequence we have divided each video sequence in three mutually exclusive time periods. The first time period correspond to initial frames of the video sequence i.e. neutral face. The last time period encapsulates the frames where the expression is shown with full intensity (apex frames). The second time period is a encapsulation of the frames which has a transition of facial expression i.e. transition from neutral face to the beginning of the desired expression (i.e. neutral to the onset of the expression). Then the fixations recorded for a particular time period are averaged across fifteen observers. For drawing the conclusions we considered second and third time periods as they have the most significant information in terms of specific displayed expression. Conclusions drawn are summarized in Fig. 3. Refer Khan et al. (2012a) for the detailed explanation of the psycho-visual experimental study.

Conclusions drawn from this psycho-visual experimental study suggests that for some expressions (i.e. happiness, sadness and surprise) only one facial region is salient while for other expressions two facial regions are salient or contain most discriminative information. We argue that the task of expression analysis and recognition could be done in more conducive manner, if only same perceptual salient regions are selected for further processing as it happens in human visual system. By processing only perceptual salient regions the proposed framework (refer Section 5) reduces the feature vector dimensionality. This reduction in feature vector length makes the proposed framework suitable for real-time applications due to minimized computational complexity.

Machine learning research community has also proposed a mathematical model called feature subset selection (FSS), to reduce feature vector dimensionality. The objective of FSS is to reduce the number of features used to characterize a dataset so as to improve a learning algorithm performance on a given task (Aha and Bankert, 1994), and it has recently been used for facial expression recognition application (Dornaika et al., 2011). Most FSS methods involve evaluating different feature subsets, employ



**Fig. 3.** Summary of the facial regions that emerged as salient for six universal expressions. Salient regions are mentioned according to their importance (for example facial expression of “fear” has two salient regions but mouth is the most important region according to HVS).

**Table 2**  
Comparison with the state-of-the-art methods for posed expressions.

	Sequence num.	Class num.	Performance measure	Recog. rate (%)
Littlewort et al. (2006)	313	7	Leave-one-out	93.3
Zhao and Pietikäinen (2007)	374	6	2-Fold	95.19
Zhao and Pietikäinen (2007)	374	6	10-Fold	96.26
Kotsia et al. (2008)	374	6	5-Fold	94.5
Tian (2004)	375	6	–	93.8
Yang et al. (2010)(a)	352	6	66% split	92.3
Yang et al. (2010)(b)	352	6	66% split	80
<b>Ours</b>	<b>309</b>	<b>6</b>	<b>10-Fold</b>	<b>96.7</b>
<b>Ours</b>	<b>309</b>	<b>6</b>	<b>2-Fold</b>	<b>95.2</b>

some criterion such as probability of error (Jain and Zongker, 1997). One difficulty with this approach when applied to real problems with large feature dimensionality, is the high computational complexity involved in searching the exponential space of feature subsets (Guo and Dyer, 2003). Jain and Zongker (1997) evaluated different search algorithms for FSS and found that the sequential forward floating selection (SFFS) algorithm proposed by Pudil et al. (1994) performed best. However, SFFS is very time consuming when the number of features is large. For example, Vailaya (2000) used the SFFS method to select 67 features from 600 for a two-class problem and reported that SFFS required 12 days of computation time.

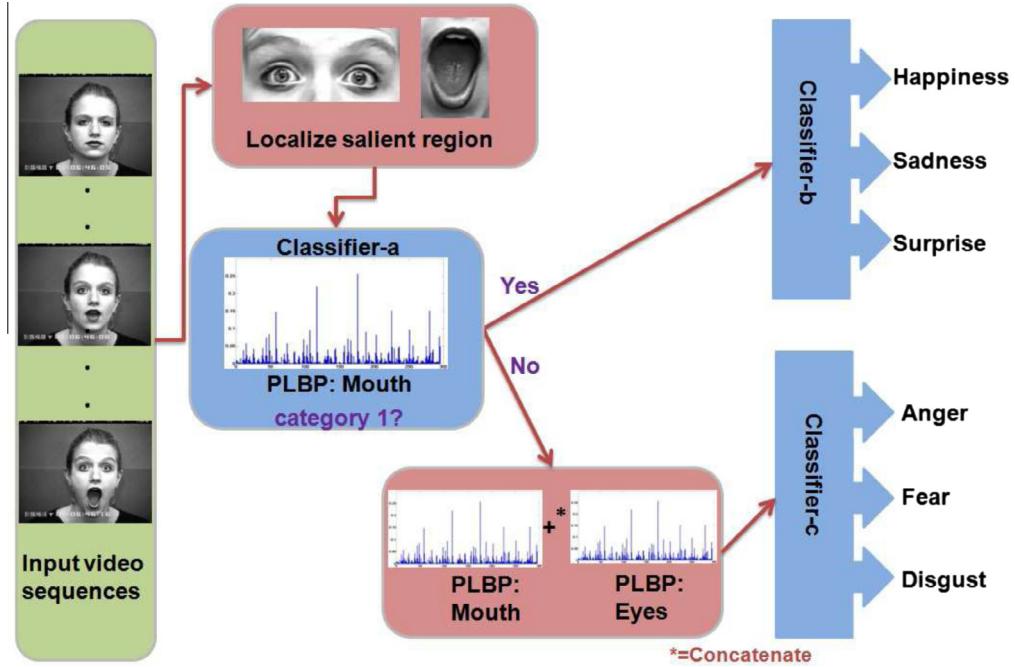
As the proposed framework is based on psycho-visual experimental study and the notion of saliency, it does not employ FSS technique to reduce feature vector dimensionality. Rather reduction in feature vector dimensionality is inherited due to the notion of saliency. Proposed framework (Section 5) extracts PLBP features only from the perceptual salient facial region(s) which contains highly discriminative information as suggested by Itti et al. (1998) and proved by the recorded results (refer Tables 1 and 2). Specifically Table 1 proves that the dimensionality of the proposed descriptor is very low as compared to other state-of-the-art descriptor. Thus, making it suitable for real-time applications. In the perspective of this study, FSS still can be utilized to reduce

the dimensionality of proposed descriptor even further (i.e. to reduce dimensionality from hundreds to tens) on the cost of increasing computational complexity.

## 5. Expression recognition framework

Feature selection along with the region(s) from where these features are going to be extracted is one of the most important step to recognize expressions. As the proposed framework draws its inspiration from the human visual system (HVS), it extracts proposed features i.e. PLBP, only from the perceptual salient facial region(s) which were determined through psycho-visual experiment. Schematic overview of the framework is presented in Fig. 4. Steps of the proposed framework are as follows:

1. The framework first localizes salient facial regions using Viola-Jones object detection algorithm (Viola and Jones, 2001). We selected this algorithm as it is the most cited and considered the fastest and most accurate pattern recognition method for face and facial region detection (Kolsch and Turk, 2004).
2. Then, the framework extracts PLBP features from the mouth region, feature vector of 295 dimensions ( $f_1, \dots, f_{295}$ ). The classification (“Classifier-a” in Fig. 4) is carried out on the basis of extracted features in order to make two groups of facial



**Fig. 4.** Schematic overview of the framework.

- expressions. First group comprises of those expressions that has one perceptual salient region i.e. happiness, sadness and surprise while the second group is composed of those expressions that have two or more perceptual salient regions i.e. anger, fear and disgust (see Section 4.3). Purpose of making two groups of expressions is to reduce feature extraction computational time.
3. If the stimuli is classified in the first group, then it is classified either as happiness, sadness or surprise by the “Classifier-b” using already extracted PLBP features from the mouth region.
  4. If the stimuli is classified in the second group, then the framework extracts PLBP features from the eyes region and concatenates them with the already extracted PLBP features from the mouth region, feature vector of 590 dimensions ( $f_1, \dots, f_{295} + f_1, \dots, f_{295}$ ). Then, the concatenated feature vector is fed to the classifier (“Classifier-c”) for the final classification. It is worth mentioning here that for the expression of “disgust” nose region emerged as one of the salient regions but the framework do not explicitly extracts features from this region. This is due to the fact that, the region of nose that emerged as salient is the upper nose (wrinkles) area which is connected and already included in the localization of the eyes region, refer Fig. 3.

## 6. Experiment and results

We performed person-independent facial expression recognition using proposed PLBP features.<sup>1</sup> We performed four experiments to test different scenarios.

1. First experiment was performed on the extended Cohn–Kanade (CK+) database (Lucey et al., 2010). This database contains 593 sequences of posed universal expressions.
2. Second experiment was performed to test the performance of the proposed framework on low resolution image sequences.

3. Third experiment tests the robustness of the proposed framework when generalizing on a new dataset.
4. Fourth experiment was performed on the MMI facial expression database (parts IV and V of the database) (Valstar and Pantic, 2010) which contains spontaneous/natural expressions.

For the first two experiments we used all the 309 sequences from the CK+ database which have FACS coded expression label (Ekman and Friesen, 1978). The experiment was carried out on the frames which covers the status of onset to apex of the expression, as done by Yang et al. (2010). Region of interest was obtained automatically by using Viola–Jones object detection algorithm (Viola and Jones, 2001) and processed to obtain PLBP feature vector. We extracted LBP features only from the salient region(s) using  $LBP_{8,2}^{\mu 2}$  operator which denotes a uniform LBP operator with eight sampling pixels in a local neighborhood region of radius 2. Only exception was in the second experiment, when we adopted  $LBP_{4,1}^{\mu 2}$  operator when the spatial facial resolution gets smaller than  $36 \times 48$ .

In our framework we created image pyramid up to level 1, so in turn got five sub-regions from one facial region i.e. mouth region (see Fig. 2). In total we obtained 295 dimensional feature vector (59 dimensions/sub-region). As mentioned earlier we adopted  $LBP_{4,1}^{\mu 2}$  operator when the spatial facial resolution was  $18 \times 24$ . In this case we obtained 75 dimensional feature vector (15 dimensions/sub-region).

We recorded correct classification accuracy in the range of 95% for image pyramid level 1. We decided not to test framework with further image pyramid levels as it would double the size of feature vector and thus increase the feature extraction time and likely would add few percents in the accuracy of the framework which will be insignificant for a framework holistically.

### 6.1. First experiment: posed expressions

This experiment measures the performance of the proposed framework on the classical database i.e. extended Cohn–Kanade (CK+) database (Lucey et al., 2010). Most of the methods in

<sup>1</sup> Video showing the result of the proposed framework on good quality image sequences is available at: [http://www.youtube.com/watch?v=RPeXBdS\\_pd8](http://www.youtube.com/watch?v=RPeXBdS_pd8).

literature report their performance on this database, so this experiment could be considered as the benchmark experiment for facial expression recognition framework.

The performance of the framework was evaluated for five different classifiers:

1. Support vector machine (SVM) with  $\chi^2$  kernel and  $\gamma = 1$
2. C4.5 decision tree (DT) with reduced-error pruning
3. Random forest (RF) of 10 trees
4. 2 Nearest neighbor (2NN) based on Euclidean distance
5. Naive Bayes (NB) classifier

Above mentioned classifiers are briefly described below.

*Support vector machine (SVM).* SVM performs an implicit mapping of data into a higher dimensional feature space, and then finds a linear separating hyperplane with the maximal margin to separate data in this higher dimensional space (Vapnik, 1995). Given a training set of labeled examples  $\{(x_i, y_i), i = 1, \dots, l\}$  where  $x_i \in \Re^n$  and  $y_i \in \{-1, 1\}$ , a new test example  $x$  is classified by the following function:

$$f(x) = \text{sgn} \left( \sum_{i=1}^l \alpha_i y_i K(x_i, x) + b \right) \quad (6)$$

where  $\alpha_i$  are Langrange multipliers of a dual optimization problem that describe the separating hyperplane,  $K(\cdot, \cdot)$  is a kernel function, and  $b$  is the threshold parameter of the hyperplane. We used Chi-Square kernel as it is best suited for histograms. It is given by:

$$K(x, y) = 1 - \sum_i \frac{2 \times (x_i - y_i)^2}{(x_i + y_i)} \quad (7)$$

*Classification Trees.* A classification tree is a classifier composed by nodes and branches which break the set of samples into a set of covering decision rules. In each node, a single test is made to obtain the partition. The starting node is called the root of the tree. In the final nodes or leaves, a decision about the classification of the case is made. In this work, we have used C4.5 paradigm (Quinlan, 1993). Random forest (RFs) are collections of decision trees (DTs) that have been constructed randomly. RFs generally performs better than DT on unseen data.

*Instance Based Learning.*  $k$ -NN classifiers are instance-based algorithms taking a conceptually straightforward approach to approximate real or discrete valued target functions. The learning process consists in simply storing the presented data. All instances correspond to points in an  $n$ -dimensional space and the nearest neighbors of a given query are defined in terms of the standard Euclidean distance. The probability of a query  $q$  belonging to a class  $c$  can be calculated as follows:

$$p(c|q) = \frac{\sum_{k \in K} W_k \cdot 1_{(kc=c)}}{\sum_{k \in K} W_k} \quad (8)$$

$$W_k = \frac{1}{d(k, q)} \quad (9)$$

$K$  is the set of nearest neighbors,  $kc$  the class of  $k$  and  $d(k, q)$  the Euclidean distance of  $k$  from  $q$ .

*Naive Bayes classifiers.* The naive-Bayes (NB) classifier uses the Bayes theorem to predict the class for each case, assuming that the predictive genes are independent given the category. To classify a new sample characterized by  $d$  genes  $X = (X_1, X_2, \dots, X_d)$ , the NB classifier applies the following rule:

$$C_N - B = \arg \max_{c_j \in C} p(c_j) \prod_{i=1}^d p(x_i | c_j) \quad (10)$$

where  $C_N - B$  denotes the class label predicted by the naive-Bayes classifier and the possible classes of the problem are grouped in  $C = \{c_1, \dots, c_l\}$ .

### 6.1.1. Results

The framework achieved average recognition rate of 96.7%, 97.9%, 96.2%, 94.7% and 90.2% for SVM, 2 nearest neighbor (2NN), random forest (RF), C4.5 decision tree (DT) and naive-Bayes (NB) respectively using 10-fold cross validation technique. One of the most interesting aspects of our approach is that it gives excellent results for a simple 2NN classifier which is a non-parametric method. This points to the fact that framework do not need computationally expensive methods such as SVM, random forests or decision trees to obtain good results. In general, the proposed framework achieved high expression recognition accuracies irrespective of the classifiers, proves the descriptive strength of the extracted features (feature minimizes within-class variations of expressions, while maximizes between class variations). For comparison and reporting results, we have used the classification results obtained by the SVM as it is the most cited method for classification in the literature.

### 6.1.2. Comparisons

We chose to compare average recognition performance of our framework with the framework proposed by Shan et al. (2009) with different SVM kernels. Our choice was based on the fact that both have common underlying descriptor i.e. local binary pattern (LBP), secondly framework proposed by Shan et al. (2009) is highly cited in the literature. Our framework obtained average recognition percentage of 93.5% for SVM linear kernel while for the same kernel Shan et al. (2009) have reported 91.5%. For SVM with polynomial kernel and SVM with RBF kernel our framework achieved recognition accuracy of 94.7% and 94.9% respectively, as compared to 91.5% and 92.6%.

In terms of time and memory cost of feature extraction process, we have measured and compared our descriptor with the frameworks proposed by Shan et al. (2009) and Bartlett et al. (2003). The results are presented in Table 1. Shan et al. (2009) have reported their results for two set of features i.e. LBP and Gabor. In Table 1 “Shan et al. (2009)(a)” corresponds to results obtained with LBP features, while “Shan et al. (2009)(b)” corresponds to Gabor feature result. Table 1 shows the effectiveness of the proposed descriptor for facial feature analysis i.e. PLBP, for real-time applications as it is memory efficient and its extraction time is much lower than other compared descriptor (see Section 5 for the dimensionality calculation). In Table 1 feature dimension reported are stored in a data type “float” and float occupies four bytes. The proposed framework is compared with the other state-of-the-art frameworks using same database (i.e. Cohn-Kanade database) and the results are presented in Table 2.

Table 2 shows the comparison of the achieved average recognition rate of the proposed framework with the state-of-the-art methods using same database (i.e. Cohn-Kanade database). Results from Yang et al. (2010) are presented for the two configurations. “Yang et al. (2010)(a)” shows the result when the method was evaluated for the last three frames from the sequence while “Yang et al. (2010)(b)” presents the reported result for the frames which encompasses the status from onset to apex of the expression. It can be observed from Table 2 that the proposed framework is comparable to any other state-of-the-art method in terms of expression recognition accuracy. The method discussed in “Yang et al. (2010)(b)” is directly comparable to our method (frames which covers the status of onset to apex of the expression). In this configuration, our framework is better in terms of average recognition accuracy.

In general, Tables 1 and 2 show that the framework is better than the state-of-the-art frameworks in terms of average expression recognition performance, time and memory costs of feature extraction processes. These results show that the system could be used with the high degree of confidence for real-time

applications as its unoptimized Matlab implementation runs at more than 30 frames/second.

### 6.2. Second experiment: low resolution image sequences

Most of the existing state-of-the-art systems for expressions recognition report their results on high resolution images without reporting results on low resolution images. As mentioned earlier there are many real world applications that require expression recognition system to work amicably on low resolution images. Smart meeting, video conferencing and visual surveillance are some examples of such applications. To compare with Tian's work (Tian, 2004), we tested our proposed framework on low resolution images of four different facial resolutions ( $144 \times 192$ ,  $72 \times 96$ ,  $36 \times 48$ ,  $18 \times 24$ ) based on Cohn-Kanade database. Tian's work can be considered as the pioneering work for low resolution image facial expression recognition. Fig. 5 shows the images at different spatial resolution along with the average recognition accuracy achieved by the different methods. Low resolution image sequences were obtained by down sampling the original sequences. All the other experimental parameters i.e. descriptor, number of sequences and region of interest, were same as mentioned earlier in the Section 6.

Fig. 5 reports the recognition results of the proposed framework with the state-of-the-art methods on four different low facial resolution images. Reported results of our proposed method i.e. are obtained using support vector machine (SVM) with  $\chi^2$  kernel and  $\gamma = 1$ . In Fig. 5 recognition curve for our proposed method is shown as PLBP-SVM, recognition curves of LBP (Shan et al., 2009) and Gabor (Shan et al., 2009) are shown as LBP[JIVC] and Gabor[JIVC] respectively, curve for Tian's work (Tian, 2004) is shown as Gabor[CVPRW] while Khan et al. (2012b) proposed system's curve is shown as PHOG[ICIP]. Results reports in LBP (Shan et al., 2009)

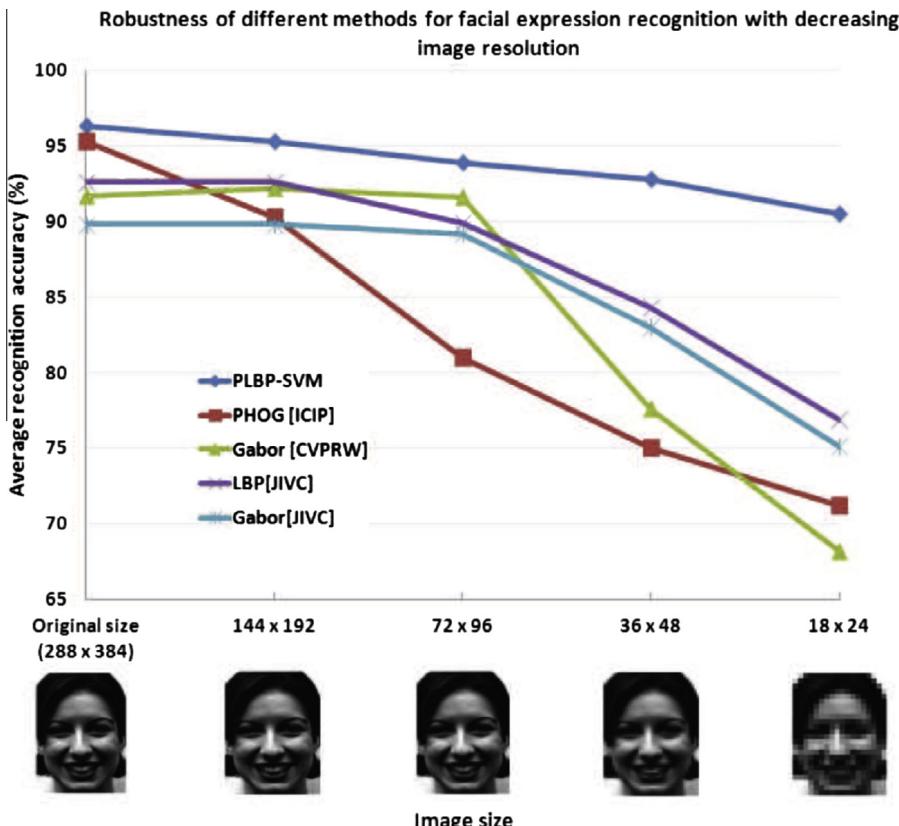
and Gabor (Shan et al., 2009), the different facial image resolution are  $110 \times 150$ ,  $55 \times 75$ ,  $27 \times 37$  and  $14 \times 19$  which are comparable to the resolutions of  $144 \times 192$ ,  $72 \times 96$ ,  $36 \times 48$ ,  $18 \times 24$  pixels in our experiment. Referenced figure shows the supremacy of the proposed framework for low resolution images. Specially for the smallest tested facial image resolution ( $18 \times 24$ ) our framework performs much better than any other compared state-of-the-art method.

Results from the first and second experiment show that the proposed framework for facial expression recognition works amicably on classical dataset (CK dataset) and its performance is not effected significantly for low resolution images. Secondly, the framework has a very low memory requirement and thus it can be utilized for real-time applications.

### 6.3. Third experiment: generalization on the new dataset

The aim of this experiment is to study how well the proposed framework generalizes on the new dataset. Valstar et al. (2005) have reported such data earlier to show generalization ability of their expression recognition system. Thus, this experiment helps to understand how the framework will behave when it will be used to classify expressions in real life videos. We used image sequences from CK+ dataset and FG-NET FEED (facial expressions and emotion database) (Wallhoff, 2006). FG-NET FEED contains 399 video sequences across 18 different individuals showing seven facial expressions i.e. six universal expression (Ekman, 1971) plus one neutral. In this dataset individuals were not asked to act rather expressions were captured while showing them video clips or still images i.e. natural expressions.

The experiment was carried out on the frames which covers the status of onset to apex of the expression as done in the previous



**Fig. 5.** Robustness of different methods for facial expression recognition with decreasing image resolution. PHOG[ICIP] corresponds to framework proposed by Khan et al. (2012b), Gabor[CVPRW] corresponds to Tian's work (Tian, 2004), LBP[JIVC] and Gabor[JIVC] corresponds to results reported by Shan et al. (2009).

**Table 3**

Average recognition accuracy (%).

	SVM	C4.5 DT	RF	2NN
<i>Training on CK+ database and testing it with FG-NET FEED</i>				
Training samples	96.7	94.7	96.2	97.9
Test samples	81.9	74.8	79.5	83.1
<i>Training on FG-NET FEED and testing it with CK+ database</i>				
Training samples	92.3	91.2	90.5	93.3
Test samples	80.5	77.3	79	84.7

experiment. This experiment was performed in two different scenarios, with the same classifier parameters as the first experiment:

- (a) In the first scenario samples from the CK+ database were used for the training of different classifiers and samples from FG-NET FEED (Wallhoff, 2006) were used for the testing. Obtained results are presented in Table 3.
- (b) In the second scenario we used samples from the FG-NET FEED for the training and testing was carried out with the CK+ database samples. Results obtained are presented in last two rows of Table 3.

This experiment simulates the real life situation when the framework would be employed to recognize facial expressions on the unseen data. Obtained results are presented in Table 3. Reported average recognition percentages for training phase were calculated using 10-fold cross validation method. It can be observed from the result that generally average recognition accuracy drops by 15% across different classifiers during testing phase. This could be due to the fact that the two databases used for training and testing phase have inherent differences i.e. played emotions differ from the natural/spontaneous ones. Obtained results can be probably further improved by training classifiers on more than one dataset before using in real life scenario.

#### 6.4. Fourth experiment: spontaneous expressions

Spontaneous/natural facial expressions differ substantially from posed expressions (Bartlett et al., 2002). The same has also been proved by psychophysical work (Ekman, 2001). To test the performance of the proposed framework on the spontaneous facial expressions we used 392 video segments from parts IV and V of the MMI facial expression database (Valstar and Pantic, 2010). Parts IV and V of the database contains spontaneous/naturalistic expressions recorded from 25 participants aged between 20 and 32 years in two different settings. Due to ethical concerns the database contains only the video recording of the expressions of happiness, surprise and disgust (Valstar and Pantic, 2010).

The framework achieved average recognition rate of 91%, 91.4%, 90.3% and 88% for SVM, 2-nearest neighbor, random forest and C4.5 decision tree respectively using 10-fold cross validation technique. Algorithm of Park and Kim (2008) for spontaneous expression recognition achieved results for three expressions in the range of 56–88% for four different configurations which is less than recognition rate of our proposed algorithm, although results cannot be compared directly as they used different database.

## 7. Conclusions and future work

We presented a novel descriptor and framework for automatic and reliable facial expression recognition. Framework is based on initial study of human vision and works adequately on posed as well as on spontaneous expressions. The key conclusions drawn from the study are:

1. Facial expressions can be analyzed automatically by mimicking human visual system i.e. extracting features only from the salient facial regions.
2. Features extracted using proposed pyramidal local binary pattern (PLBP) operator have strong discriminative ability as the recognition result for six universal expressions is not effected by the choice of classifier.
3. The proposed framework is robust for low resolution images, spontaneous expressions and generalizes well on unseen data.
4. The proposed framework can be used for real-time applications since its unoptimized Matlab implementation run at more than 30 frames/second on a Windows 64 bit machine with i7 processor running at 2.4 GHz having 6 GB of RAM.

In future we plan to investigate the effect of occlusion as this parameter could significantly impact the performance of the framework for real world applications. Secondly, the notion of movement could improve the performance of the proposed framework for real world applications as the experimental study conducted by Bassili (1979) suggested that dynamic information is important for facial expression recognition. Another parameter that needs to be investigated is the variations of camera angle as for many applications frontal facial pose is difficult to record. Lastly, in future we would also like to evaluate the proposed framework on streaming artifacts i.e. ringing, contouring, posterization, etc.

## Acknowledgment

This research work is funded by the Région Rhône-Alpes, France through ARC 6.

## References

- Aha, D., Bankert, R.L., 1994. Feature selection for case-based classification of cloud types: an empirical comparison. In: Association for the Advancement of Artificial Intelligence (AAAI) Workshop on Case-Based Reasoning. AAAI Press, pp. 106–112.
- Bai, Y., Guo, L., Jin, L., Huang, Q., 2009. A novel feature extraction method using pyramid histogram of orientation gradients for smile recognition. In: International Conference on Image Processing.
- Bartlett, M.S., Littlewort, G., Braathen, B., Sejnowski, T.J., Movellan, J.R., 2002. A prototype for automatic recognition of spontaneous facial actions. In: Advances in Neural Information Processing Systems.
- Bartlett, M.S., Littlewort, G., Fasel, I., Movellan, J.R., 2003. Real time face detection and facial expression recognition: development and applications to human computer interaction. In: Conference on Computer Vision and Pattern Recognition Workshop, 2003.
- Bassili, J.N., 1979. Emotion recognition: the role of facial movement and the relative importance of upper and lower areas of the face. Journal of Personality and Social Psychology 37, 2049–2058.
- Donato, G., Bartlett, M.S., Hager, J.C., Ekman, P., Sejnowski, T.J., 1999. Classifying facial actions. IEEE Transaction on Pattern Analysis and Machine Intelligence 21, 974–989.
- Dornaika, F., Lazcano, E., Sierra, B., 2011. Improving dynamic facial expression recognition with feature subset selection. Pattern Recognition Letters 32, 740–748.
- Ekman, P., 1971. Universals and cultural differences in facial expressions of emotion. In: Nebraska Symposium on Motivation. Lincoln University of Nebraska Press, pp. 207–283.
- Ekman, P., 2001. Telling Lies: Clues to Deceit in the Marketplace, Politics, and Marriage, third ed. W.W. Norton & Company, New York.
- Ekman, P., Friesen, W., 1978. The facial action coding system: a technique for the measurement of facial movements. Consulting Psychologist.
- Guo, G., Dyer, C.R., 2003. Simultaneous feature selection and classifier training via linear programming: a case study for face expression recognition. In: IEEE Conference on Computer Vision and Pattern Recognition. IEEE Computer Society, Washington, DC, USA, pp. 346–352.
- Guo, Z., Zhang, L., Zhang, D., Mou, X., 2010. Hierarchical multiscale lbp for face and palmprint recognition. In: IEEE International Conference on Image Processing, pp. 4521–4524.
- Hadjidemetriou, E., Grossberg, M., Nayar, S., 2004. Multiresolution histograms and their use for recognition. IEEE Transactions on Pattern Analysis and Machine Intelligence 26, 831–847.

- Itti, L., Koch, C., Niebur, E., 1998. A model of saliency-based visual attention for rapid scene analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 20, 1254–1259.
- Jain, A., Zongker, D., 1997. Feature selection: evaluation, application, and small sample performance. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 19, 153–158.
- Khan, R.A., Meyer, A., Konik, H., Bouakaz, S., 2012a. Exploring human visual system: study to aid the development of automatic facial expression recognition framework. In: Computer Vision and Pattern Recognition Workshop.
- Khan, R.A., Meyer, A., Konik, H., Bouakaz, S., 2012b. Human vision inspired framework for facial expressions recognition. In: IEEE (Ed.), *IEEE International Conference on Image Processing*.
- Kolsch, M., Turk, M., 2004. Analysis of rotational robustness of hand detection with a viola-jones detector. In: 17th International Conference on Pattern Recognition, vol. 3, pp. 107 – 110.
- Kotsia, I., Zafeiriou, S., Pitas, I., 2008. Texture and shape information fusion for facial expression and facial action unit recognition. *Pattern Recognition* 41, 833–851.
- Littlewort, G., Bartlett, M.S., Fasel, I., Susskind, J., Movellan, J., 2006. Dynamics of facial expression extracted automatically from video. *Image and Vision Computing* 24, 615–625.
- Lucey, P., Cohn, J.F., Kanade, T., Saragih, J., Ambadar, Z., Matthews, I., 2010. The extended Cohn-Kanade dataset (CK+): a complete facial expression dataset for action unit and emotion-specified expression. In: IEEE Conference on Computer Vision and Pattern Recognition Workshop.
- Moore, S., Bowden, R., 2011. Local binary patterns for multi-view facial expression recognition. *Computer Vision and Image Understanding* 115, 541–558.
- Ojala, T., Pietikäinen, M., 1999. Unsupervised texture segmentation using feature distributions. *Pattern Recognition* 32, 477–486.
- Ojala, T., Pietikäinen, M., Harwood, D., 1996. A comparative study of texture measures with classification based on featured distribution. *Pattern Recognition* 29, 51–59.
- Ojala, T., Pietikäinen, M., Mäenpää, T., 2002. Multiresolution gray-scale and rotation invariant texture classification with local binary patterns. *IEEE Transaction on Pattern Analysis and Machine Intelligence* 24, 971–987.
- Ojansivu, V., Heikkilä, J., 2008. Blur insensitive texture classification using local phase quantization. In: International conference on Image and Signal Processing.
- Pantic, M., Patras, I., 2006. Dynamics of facial expression: recognition of facial actions and their temporal segments from face profile image sequences. *IEEE Transactions on Systems, Man, and Cybernetics* 36, 433–449.
- Park, S., Kim, D., 2008. Spontaneous facial expression classification with facial motion vector. In: IEEE Conference on Automatic Face and Gesture Recognition.
- Pudil, P., Novovičová, J., Kittler, J., 1994. Floating search methods in feature selection. *Pattern Recognition Letters* 15, 1119–1125.
- Quinlan, J.R., 1993. C4.5: Programs for Machine Learning. Morgan Kaufmann.
- Shan, C., Gong, S., McOwan, P.W., 2009. Facial expression recognition based on local binary patterns: a comprehensive study. *Image and Vision Computing* 27, 803–816.
- Tian, Y., 2004. Evaluation of face resolution for expression analysis. In: Computer Vision and Pattern Recognition Workshop.
- Vailaya, A., 2000. Semantic classification in image database. Ph.D. thesis. Michigan State University.
- Valstar, M., Pantic, M., 2010. Induced disgust, happiness and surprise: an addition to the MMI facial expression database. In: International Language Resources and Evaluation Conference.
- Valstar, M., Patras, I., Pantic, M., 2005. Facial action unit detection using probabilistic actively learned support vector machines on tracked facial point data. In: IEEE Conference on Computer Vision and Pattern Recognition Workshop, pp. 76–84.
- Vapnik, V.N., 1995. *The Nature of Statistical Learning Theory*. Springer-Verlag, New York.
- Viola, P., Jones, M., 2001. Rapid object detection using a boosted cascade of simple features. In: IEEE Conference on Computer Vision and Pattern Recognition.
- Wallhoff, F., 2006. Facial expressions and emotion database. [www.mmk.ei.tum.de/~waf/fgnet/feedtum.html](http://www.mmk.ei.tum.de/~waf/fgnet/feedtum.html).
- Wang, W., Chen, W., Xu, D., 2011. Pyramid-based multi-scale lbp features for face recognition. In: International Conference on Multimedia and Signal Processing (CMSP), pp. 151–155.
- Yang, P., Liu, Q., Metaxas, D.N., 2010. Exploring facial expressions with compositional features. In: IEEE Conference on Computer Vision and Pattern Recognition.
- Zhang, Y., Ji, Q., 2005. Active and dynamic information fusion for facial expression understanding from image sequences. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 27, 699–714.
- Zhao, G., Pietikäinen, M., 2007. Dynamic texture recognition using local binary patterns with an application to facial expressions. *IEEE Transaction on Pattern Analysis and Machine Intelligence* 29, 915–928.
- Zhaoping, L., 2006. Theoretical understanding of the early visual processes by data compression and data selection. *Network: Computation in Neural Systems* 17, 301–334.