

# Histogram-based Local Descriptors for Facial Expression

## Recognition (FER): A comprehensive Study

Cigdem Turan\*, Kin-Man Lam

Centre for Signal Processing, Department of Electronic and Information Engineering, The Hong Kong Polytechnic University, Kowloon, Hong Kong

\*Corresponding author.

E-mail addresses: cigdem.turan@connect.polyu.hk (C. Turan), enkmlam@polyu.edu.hk (K.-M. Lam)

Postal address: DE503, The Hong Kong Polytechnic University, 11 Yuk Choi Rd, Hung Hom, Hong Kong

### ABSTRACT

This paper aims to present histogram-based local descriptors applied to Facial Expression Recognition (FER) from static images, and provide a systematic review and analysis of them. First, we describe the main steps in encoding binary patterns in a local patch, which are required in every histogram-based local descriptor. Then, we list the existing local descriptors, while analysing their strengths and weaknesses. Finally, we present the experimental results of all these descriptors on commonly used facial expression databases, with varying resolution, noise, occlusion, and number of sub-regions, as well as comparing them with the results obtained by the state-of-the-art deep learning methods. This paper aims to bring together different studies of the visual features for FER by evaluating their performances under the same experimental setup, and critically reviewing various classifiers making use of the local descriptors.

### 1. Introduction

Facial expressions, which are an important aspect of non-verbal communication, have been extensively studied in different fields, such as psychology [7; 8]. Ekman and Friesen [7] identified six facial expressions (i.e. anger, disgust, fear, joy, sadness, and surprise) as prototypical expressions that are universal among humans regardless of their age, race and gender. Early research on automatic FER focused on those six emotions [14; 15; 16]. In recent years, FER has shown its importance in human-

27 computer interaction (HCI), such as assistive driving [19], embodied agents [22], and in applications  
28 such as diagnosis [27; 28] and computer games [29]. Thus, the demand has been increasing for an  
29 effective FER technology that can determine one's emotional state, based on face images, regardless of  
30 one's age, gender or race. Although much progress has been made on recognizing facial expressions, it  
31 is still a difficult task, due to the complexity and variability of facial expressions, and an effective facial-  
32 representation method is a vital step to improving the recognition rate in FER.

33 Facial feature representations proposed in the literature can now be divided into three categories:  
34 geometrical, appearance-based, and deep features. Geometrical features [32] take advantage of shape  
35 and location information of facial components and salient points, i.e. the eyes, lips, nose tip, etc. FER  
36 with Action Unit (AU) recognition is a geometrical feature-based approach, which has achieved more  
37 attention recently with the advancement in deep neural-network structures [35; 36]. However,  
38 geometrical features still require an accurate and reliable reconstruction and tracking of the facial  
39 landmarks. Therefore, it is difficult to achieve in real-life situations. Furthermore, AU-based facial  
40 expression recognition may require training data whose Action Units are already labelled by experts,  
41 which is a labor-intensive and time-consuming process. Recent studies have shown that appearance-  
42 based methods can achieve similar or better performance than AU recognition-based methods [10].

43 Appearance-based features are based on texture information related to the expressions on a face,  
44 e.g. wrinkles, skin changes, etc., which can be applied to the whole face or specific facial regions.  
45 Appearance-based features do not require the accurate reconstruction of all the facial landmarks, but  
46 only eye-pupil points, since the eyes are usually used to align the faces for further facial representation,  
47 i.e. feature extraction. Furthermore, appearance-based features only need the emotion labels of the  
48 samples for the training process. These advantages make appearance-based methods more favorable in  
49 comparison to geometrical features. One of the first attempts of FER, based on texture classification, is  
50 to use Local Binary Pattern (LBP), which was proposed by Ojala et al. [24]. LBP is one of the most  
51 widely used descriptors, due to its computational simplicity, discriminative power, and insensitivity to  
52 monotonic grayscale changes.

53 The successful application of LBP on the FER problems has inspired further studies for local  
54 descriptors. These studies focus on enhancing the coding techniques, e.g. different neighbourhood sizes,

55 processing of input images, e.g. linear filtering, transformations, etc., to emphasize the expression-  
56 specific information. Numerous variants of LBP have been proposed for the problems, such as face  
57 recognition [43; 44], facial expression recognition [46], texture classification [47], spatiotemporal  
58 feature representation [50], and medical image analysis [52]. Some comprehensive studies of LBP  
59 variants can be found in [53; 54; 55].

60 Recently, local binary feature learning methods have been proposed for efficient and data-adaptive  
61 face representation, because LBP and other hand-crafted features require strong prior knowledge of the  
62 problem in order to engineer them by hand [56; 57; 58]. The objective behind the feature learning  
63 methods is to learn a feature mapping using raw pixels to project each local pixel difference (PDV) into  
64 a low-dimensional binary vector that can efficiently represent the face data. Therefore, a codebook  
65 constructed using the learned binary codes can be used to obtain a histogram feature for each image  
66 [58]. To the best of our knowledge, local binary feature learning methods have not been applied to the  
67 FER problem, but only to age estimation [62] and face recognition [63; 64; 65]. As LBP has been  
68 successfully applied to the tasks for facial image analysis, it is worthwhile evaluating the recently  
69 proposed local binary feature learning methods on FER.

70 Recently, deep neural networks have been studied widely for many pattern-recognition tasks, such  
71 as human pose estimation [66], face recognition [68], gender recognition [69], image recognition [71;  
72 72], which require learning from a large amount of data. The increasing popularity and the success of  
73 deep features are also rooted in the FER problems [1; 3; 4; 12; 17; 20; 25]. Although the increase in  
74 recognition rate for FER is undeniable, the debate between hand-crafted features and deep features is  
75 still active. Benitez-Garcia et al. [23] proposed a local descriptor, i.e. a handcrafted feature, which can  
76 achieve a higher recognition rate than any deep neural-network structure until now. This suggests that  
77 the domain-specific knowledge and the handcrafted features are still effective and favourable for visual  
78 classification. In this paper, we present a comprehensive study of appearance-based facial features, i.e.  
79 handcrafted features, and then we compare their best results with those methods based on recently  
80 proposed local binary feature learning methods and deep features for FER.

81 The steps of a basic FER framework with the use of appearance-based features can be listed as  
82 follows: 1) detecting and aligning the face images, 2) dividing each face image into several overlapping

83 or non-overlapping regions, 3) extracting local features from these regions based on the local  
84 descriptors, 4) concatenating the respective local features to form a single feature vector, followed by  
85 unsupervised or supervised dimensionality reduction, 5) training a classifier based on the feature vectors  
86 from training samples, and 6) predicting the class label of a new query based on the trained classifier.  
87 The classification results depend on almost every step listed above. However, most of the recent studies  
88 have focused only on developing more robust local features [23; 31; 73; 75]. A robust feature should  
89 be highly discriminative, easily computed, of low dimensionality, insensitive to noise, such as  
90 illumination changes, and have low intra-class variations.

91 It is difficult to balance these properties for a local descriptor. For example, LBP is computationally  
92 simple and discriminative, but sensitive to random noise. Similarly, although Gabor-based local  
93 descriptors have shown their achievements, especially in face recognition [38; 76; 77; 78], the features  
94 suffer from the expensive computational requirement and high dimensionality. Thus, developing a  
95 robust local descriptor is still an open issue for many fields of image representation and classification,  
96 such as texture representation [75; 79; 80; 81] and face representation [44; 75].

97 In the field of computer vision, popular local descriptors are often applied to different problems or  
98 applications. For instance, although LBP was originally devised for texture classification, it has been  
99 applied to face recognition [82], image retrieval [83], facial expression recognition [16], etc. However,  
100 it might not always be true for a new descriptor. Facial expression recognition is a problem different  
101 from face recognition or other types of recognition. “A good face-recognition local descriptor” should  
102 represent discriminative identity information about face images, while “a good facial-expression local  
103 descriptor” should discard the subject’s identity information and highlight the expression-specific  
104 information of a face. Therefore, it is important to be attentive to the nature of a problem in choosing  
105 an appropriate local descriptor.

106 The local descriptors, proposed in the literature, often benchmark their results against previously  
107 reported ones. However, the reliability of the benchmarking may not be high, due to the following  
108 reasons:

- 109 - A few benchmark databases were used, and the descriptors were evaluated with different  
110 databases.

- 111 - Each of the databases may have a different set of expression categories.
- 112 - Different image preprocessing techniques, e.g. face alignment, illumination, different  
113 normalization, etc., are used in experiments.
- 114 - The evaluation procedures/testing protocols, e.g. the choice of the classifier, the cross-  
115 validation scheme used, etc., are different.
- 116 - The overall experiments cannot be reproduced because not all the experimental setup is known.

117 In the literature, there have been several attempts to compare the performances of LBP-like  
118 descriptors using the same experimental settings. One of the most recent experimental studies on the  
119 LBP-like descriptors was conducted by Liu et al. [84], which evaluated thirty-two LBP variants for  
120 texture classification. However, there are still many other texture descriptors for facial expression  
121 recognition, which should be compared.

122 Kristensen et al. [55] presented an overview of “binary flavored features” for FER. Although a set  
123 of commonly used terms was defined so as to encourage consistency in terminology and to explain the  
124 current challenges, the depth of the survey in terms of performance comparison is limited. Another aim  
125 of this paper is to fill this gap by providing a comprehensive performance analysis on those recent local  
126 descriptors used for FER.

127 In this paper, we compare the performances of 27 local descriptors on four popular databases with  
128 the same experimental setup, including the use of two classifiers, different image resolutions, and  
129 different numbers of sub-regions. In addition to their accuracy, other important aspects, such as face  
130 resolutions for best performances, are also studied. Moreover, we compare the results achieved by  
131 handcrafted features, e.g. histogram-based local features, with the results obtained by the “Compact  
132 Binary Face Descriptor (CBFD) [57]” and the state-of-the-art deep features. We also evaluate the  
133 robustness of the respective local descriptors in the scenario of a cross-dataset facial expression  
134 recognition problem. In our evaluation, we found that the best overall performances are obtained by  
135 Local Phase Quantization (LPQ) and Local Gabor Binary Pattern Histogram Sequence (LGBPHS), with  
136 consistency across most of the databases used in our experiments.

137 The rest of the paper is organized as follows: Section 2 introduces a taxonomy for histogram-based  
138 local descriptors and highlights the representative examples of the specific steps. In Section 3, the

139 experimental setup is first described, then comprehensive experimental results are presented. Section 4  
140 concludes the paper.

## 141 **2. Construction of the histogram-based local descriptors**

142 Histogram-based local descriptors compute local statistical information at key points, and describe the  
143 features in a region using a histogram representation. Almost all the local statistical feature (LSF)  
144 methods, as described in [43], have two main parts: statistical histogram feature extraction and statistical  
145 feature combination. Unlike [43] which divides the statistical histogram feature extraction further into  
146 three steps, we divide it into five steps in this paper, in order to describe different local descriptors in  
147 more detail. In the rest of this section, each step is explained while the corresponding representative  
148 descriptors are highlighted with their strengths and weaknesses.

### 149 **2.1. Local variation coding**

150 Histogram-based local-feature descriptors represent the centre pixel of a local region as a decimal  
151 number, according to its values compared to its neighbouring pixels. Regardless of the input image,  
152 local variation coding is a general method used to encode the pattern features in a local patch. For each  
153 local patch, with a given neighbourhood, a typical local variation coding has five steps, including linear  
154 filtering, quantization, binarization, encoding and binary to decimal conversion. In the following sub-  
155 sections, these five steps will be explained in detail.

#### 156 **2.1.1. Linear filtering**

157 The first step of local variation coding is to convolve a patch with a predefined set of linear filters. The  
158 most commonly used linear filters in histogram-based local descriptors are Kirsch [33; 34], Prewitt [5;  
159 6; 33], Sobel [5; 6; 13; 33; 40], and Derivative-Gaussian [33].

160 From the computational point of view, Sobel operators are more efficient than the Kirsch operators,  
161 as less pixels and multiplications are involved. These linear filters operate on a local patch with a  $3 \times 3$   
162 mask, and custom linear filters, which consider higher-order derivatives, have also been proposed. For  
163 example, Local Arc Pattern (LAP) [21] and Local Monotonic Pattern (LMP) [45] encode the first and  
164 the second-order derivatives of a local patch in different orientations, using a set of custom filters.  
165 Although LAP and LMP can represent a bigger micro pattern with multiple radii, they use intensity

166 values, as LBP, and are therefore sensitive to non-monotonic changes. Local Transitional Pattern (LTrP)  
167 [59] and Local Monotonic Pattern (LMP) [45] encode the transition of intensity change in different  
168 directions over a local patch. Local Derivative Pattern (LDP) [85] encodes the second and higher-order  
169 derivatives of a local patch. Although the higher-order derivatives can represent local variations with  
170 more details, the dimensionality of the resulting feature vector will become higher, as well as the  
171 computational cost.

### 172 **2.1.2. Quantization**

173 The second step of the local variation coding is the quantization of the linear-filter responses. The most  
174 common way of quantization used in the different descriptors is the unit step function. The local  
175 descriptors, such as LBP, Median Binary Pattern (MBP) [51], etc., quantize their filter responses using  
176 the unit-step function. However, this will generate inconsistent binary codes in uniform and near-  
177 uniform face regions, because the filter responses may vary slightly around the threshold value, usually  
178 zero. Local Ternary Pattern (LTeP) [51], Median Ternary Pattern (MTP) [51], Gradient Directional  
179 Pattern (GDP) [5; 6], and Gradient Local Ternary Pattern (GLTeP) [10; 13] add an extra level of  
180 thresholding, which facilitates the generation of more consistent codes for local patterns in smooth  
181 facial regions, as well as highly textured regions.

182 Quantization of the filter responses does not necessarily result in binary values. A common way of  
183 non-binary quantization is the  $k$ -bin method. Histogram of Oriented Gradients (HOG) [86] and  
184 Pyramids of Histogram of Oriented Gradients (PHOG) [70] are two examples, which quantize the  
185 gradient angles to  $k$  intervals, and then count the gradient magnitudes of those pixels whose gradient  
186 orientations are within a specific interval. Another method of non-binary quantization of the filter  
187 responses, such as the angle or phase information, is to use the quadrant information [37; 43; 76], i.e.  
188 the 2-D Cartesian coordinate system, where four quadrants are defined by the  $x$ - and  $y$ -axes.

### 189 **2.1.3. Binarization**

190 After quantization, the filter responses of some descriptors, such as LBP [24], GDP [5; 6], have already  
191 been in binary form, i.e. 0 and 1. However, the other descriptors need a binarization process. The filter  
192 responses can be binarized in two ways:

193 *Binarization by splitting into different levels of binary codes*: One example of this method is LTeP [51],  
194 which has three levels after thresholding. A common way of encoding these three-level responses is to  
195 split the responses into two binary codes: “1” and “0” form one binary code, while “0” and “-1” form  
196 the other one. Therefore, two histograms are formed, and this results in a higher dimensional feature  
197 vector.

198 *Binarization by logical operators*: This method can be utilized in two different circumstances: when  
199 the quantized values are in binary form [45; 59], or not in binary form [37]. The common logical  
200 operators are “AND” [45] and “XOR” [18; 59]. These two logical operators have their unique  
201 advantages in information encoding. “AND” encodes the likeness/sameness of the values, while  
202 “XOR” encodes the opposition between the values.

#### 203 **2.1.4. Encoding**

204 The bits in a binary codeword correspond to the binarized responses of the different abovementioned  
205 filters. A basic way of creating a codeword is to use all the resultant binary codes to form a string. In  
206 the case of  $3 \times 3$  neighborhood, i.e. 8 neighbours, each code string will be 8-bit long, which forms a  
207 decimal value between 0 and 255. LBP, LTeP, MBP, MTP and GDP utilize this basic code. Local  
208 Directional Pattern (LDiP) [26] computes the eight directional edge responses, by using the Kirsch  
209 masks. However, as the response values are not equally important in all the directions, LDiP encodes  
210 the  $k$  most prominent directions, i.e. a customized codeword. LDiP can provide more stable codes, in  
211 the presence of gray-level distortion, such as noise and non-monotonic illumination changes. High-  
212 frequency regions in a face carry more information about texture information, such as the human eye  
213 regions. Therefore, to achieve a more competent face representation, textural regions with high  
214 contrast/frequency should influence the LDiP code more. However, LDiP considers both low and high-  
215 frequency regions equally. To incorporate this importance into the LDiP codes, an extension of LDiP,  
216 named Local Directional Pattern Variance (LDiPv) [30], was proposed, which introduces the variance  
217 of the codes as weights in constructing the histograms. However, both LDiP and LDiPv consider the  
218 filter responses in absolute value, which lose the important direction information, e.g. different  
219 transitions in a region. Furthermore, they are sensitive to rotation variations, because a fixed start



220 position has to be defined for encoding a binary string, and they are profoundly dependent on the  
221 number of the most prominent directions considered. Local Directional Number Pattern (LDN) [33]  
222 also encodes the principal directions, i.e. the most positive and negative directions, so a more  
223 discriminative representation of directions can be achieved. Local Directional Texture Pattern (LDTP)  
224 [34] also encodes the principal directions, which discards the insignificant details that may vary on the  
225 samples belonging to the same class. However, different from the other descriptors, LDTP encodes both  
226 the principal directions and the intensity information (the intensity difference of the two principal  
227 directions). Therefore, LDTP is robust against both rotation and illumination changes.

228 Recently in the fields of texture classification, image retrieval, and facial feature representation, an  
229 extensive amount of customized coding schemes has been proposed [44; 75; 80; 81]. All these coding  
230 schemes aim at producing robust features, which are important for the image-classification problem.

#### 231 **2.1.5. Binary to decimal conversion**

232 The last step of local variation coding is to convert a binary codeword into a decimal value, which  
233 represents the local pattern of the pixel under consideration. After computing the feature values for all  
234 the pixels in a patch, the statistics of these numbers, in the form of a histogram, can be used to represent  
235 the patch.

#### 236 **2.1.6. Local Binary Pattern and other local variation coding schemes**

237 LBP, as a local variation coding method, has four steps as discussed previously: linear filtering,  
238 quantization with the unit step function, encoding the binary codeword, and binary to decimal  
239 conversion. LBP has also been extended to use different neighbourhood sizes, as well as uniform LBP  
240 codewords, i.e. those codewords have no more than two transitions from 1 to 0 or 0 to 1. A codeword  
241 is non-uniform if it has more than two transitions. This idea was inspired by the fact that the uniform  
242 codewords occur more frequently than those non-uniform codewords in images.

243 LBP encodes the relationship between the central pixel and its neighbours. Some local descriptors  
244 extract high-order local information. A high-order descriptor can capture more detailed discriminative  
245 information. Other local descriptors also encode different distinctive spatial relationships in a local  
246 region. More information about LBP variants can be found in [84].

## 247 **2.2. Local feature representation**

248 LBP and other histogram-based local descriptors encode the distribution of local variation codes within  
249 a region. A frequency-based or weighted-vote-based histogram constructed for a whole face image will  
250 lose the spatial information about the patterns encoded by a local descriptor. To represent the facial  
251 features more effectively, face images are divided into a number of overlapping or non-overlapping  
252 small sub-regions. Local features extracted from the sub-regions can achieve better recognition rates  
253 than those using holistic features, such as Eigenfaces and Fisherfaces [87].

254 Different regions in a face carry different amount of information about an expression. To eliminate  
255 the excessive and non-informative features for face or expression recognition, weighted histogram  
256 representation has been adopted. In this representation, weights are often set according to the  
257 discriminability of the regions [16], e.g. a small weight near the image's borders, and a higher weight  
258 around the eye and mouth regions.

259 Another local-feature representation uses only those regions that carry salient information about  
260 facial expressions. Benitez-Garcia et al. [23] developed an algorithm to detect salient regions based on  
261 fiducial points for feature extraction. In [15], we observed that the features extracted from the eye and  
262 mouth regions can achieve higher recognition rates than the features extracted from the sub-regions  
263 divided from a whole face.

## 264 **2.3. Inputs to local variation coding**

265 Most of the early descriptors extract local features from intensity information, using a local variation  
266 coding method. However, the intensity information is sensitive to noise and illumination variations.  
267 Therefore, other types of input have been considered for local variation coding. Since gradients are  
268 more stable than intensity under the presence of illumination variations, several descriptors utilize  
269 gradient information to encode local variations. For example, GDP encodes gradient angles, while  
270 GLTeP encodes gradient magnitudes.

271 After the successful applications of LBP, several descriptors, which are based on Gabor filtering  
272 with a predefined number of scales and orientations, have been proposed. Examples of these descriptors  
273 include Local Gabor Binary Pattern Histogram Sequence (LGBPHS) [38], Local Gabor Directional

274 Pattern (LGD<sub>i</sub>P) [39], and Local Gabor Transitional Pattern (LGTrP) [42]. These descriptors often  
275 encode the magnitude information of the transform, i.e. the Gabor Magnitude Image, because the  
276 magnitude information is robust to misalignment. Gabor features are robust to image variations in terms  
277 of illumination and noise, but extracting the features is computationally expensive and the resulting  
278 feature vector has a high dimensionality.

279 Binary Pattern of Phase Congruency (BPPC) [2] applies wavelet transform to the logarithmic Gabor  
280 features, followed by computing the phase congruency (PC). PC is a dimensionless quantity, and can  
281 be considered as the gradient where high energy values of PC occur on edges, corners, etc. Monogenic  
282 signal analysis [88], which is a 2-D generalization of the 1-D analytic signal, is an alternative method  
283 to Gabor filtering. Monogenic signal analysis can estimate the multi-resolution amplitude, orientation,  
284 and phase components of a signal, which represent the signal energetic, structural, and geometric  
285 information, respectively. One advantage of monogenic signal analysis over Gabor transformations is  
286 that it has a lower time and space complexity.

287 In 2010, two local descriptors, which use monogenic signal analysis, were proposed for texture  
288 classification [89] and face recognition [90], where only the monogenic phase information and both the  
289 amplitude and orientation information, respectively, are encoded. Several other this kind of local  
290 descriptors exist in the literature [44; 91; 92]. Monogenic signal analysis has also been used for  
291 spatiotemporal facial expression recognition, with the local descriptor named “Spatiotemporal Local  
292 Monogenic Binary Patterns (STLMBP)” [93]. However, to the best of our knowledge, Monogenic  
293 Binary Coding (MBC) [43] is the only descriptor that applied monogenic signal analysis to static facial  
294 expression images [61]. MBC encodes the amplitude (MBC<sub>A</sub>), phase (MBC<sub>P</sub>), and orientation  
295 (MBC<sub>O</sub>) information separately.

296 Local Phase Quantization (LPQ) [48] is a local descriptor, which extracts features from the discrete  
297 Fourier transform (DFT) over an image. LPQ is robust against blur and low resolution because it  
298 quantizes the phase information in local neighbourhoods. However, LPQ requires the point spread  
299 function (PSF) to be positive and valued in the low-frequency domain. Local Frequency Descriptor  
300 (LFD) [37], which also extracts information from DFT, encodes both the magnitude and phase

301 information using LBP and Local XNOR Pattern (LXNORP). LFD does not require PSF to be positive,  
302 and can carry more information than LPQ, but the dimension of the feature vector is doubled.

303 Weber Local Descriptor (WLD) [69, 70] was inspired by the Weber’s Law, which states that the  
304 significance of a change in the stimuli depends on the initial value of the stimuli. WLD, which computes  
305 the differential excitation and the orientation of an image, forms a joint histogram for the differential  
306 excitation and the orientation. WLD has been applied to several problems successfully, including facial  
307 expression recognition [74]. However, WLD discards the orientation information of the differential  
308 excitation and neighbouring pixel pairs.

309 Recently, Jang et al. [18] proposed an extension of WLD, named Improved Weber Binary Coding  
310 (IWBC), to solve the drawbacks of WLD. IWBC generates two images, which are called the Novel  
311 Weber Magnitude Image and the Novel Weber Orientation Image, which are then encoded using Local  
312 XOR Pattern and LBP, respectively. Although IWBC can represent a face more accurately than WLD  
313 by including the orientation information about the neighbouring pixels, it suffers from the problem of  
314 high dimensionality. To the best of our knowledge, IWBC has never been applied to the FER problem.  
315 Since IWBC has been shown to outperform WLD on the face recognition problem, so we include IWBC  
316 in our experiments to evaluate its performance as a local descriptor for FER.

### 317 **3. Experiments**

318 In this section, a number of histogram-based local descriptors are evaluated for facial-expression  
319 recognition, with the same experiment settings. We will first describe the experimental setup, including  
320 the benchmark databases, pre-processing, feature extraction, and classification schemes, and then  
321 analyse the experimental results.

#### 322 **3.1. Experimental setup**

##### 323 **3.1.1. Databases and the corresponding numbers of expression classes**

324 The performances of the local descriptors are compared on commonly-used, acted databases, as well as  
325 spontaneous databases. The facial-expression databases used in our experiments are BAUM-2 [94],  
326 CK+ [95], JAFFE [96], and TFEID [97].

327 The CK+ database, which is one of the acted facial-expression databases mostly used, contains a  
 328 total of 593 posed sequences across 123 subjects. 327 of the sequences were labelled with one of the  
 329 seven discrete expressions — anger, contempt, disgust, fear, happiness, sadness, and surprise. The last

A list of the descriptors, and the corresponding feature dimensions, used in our experiments.

Abbreviation	Descriptor Name	Dimension
1 BPPC [2]	Binary Pattern of Phase Congruency	1062
2 GDP [5; 6]	Gradient Directional Pattern	256
3 GDP2 [9]	Gradient Direction Pattern	8
4 GLTeP [10; 13]	Gradient Local Ternary Pattern	512
5 IWBC [18]	Improved Weber Binary Coding	2048
6 LAP [21]	Local Arc Pattern	272
7 LBP [24]	Local Binary Pattern	59
8 LDiP [26]	Local Directional Pattern	56
9 LDiPv [30]	Local Directional Pattern Variance	56
10 LDN [33]	Local Directional Number Pattern	56
11 LDTP [34]	Local Directional Texture Pattern	72
12 LFD [37]	Local Frequency Descriptor	512
13 LGBPHS [38]	Local Gabor Binary Pattern Histogram Sequence	256
14 LGDiP [39]	Local Gabor Directional Pattern	280 *
15 LGIP [40]	Local Gradient Increasing Pattern	37
16 LGP [41]	Local Gradient Pattern	7
17 LGTrP [42]	Local Gabor Transitional Pattern	256
18 LMP [45]	Local Monotonic Pattern	256
19 LPQ [48; 49]	Local Phase Quantization	256
20 LTeP [51]	Local Ternary Pattern	512
21 LTrP [59; 60]	Local Transitional Pattern	256
22 MBC [43; 61]	Monogenic Binary Coding	3072 *
23 MBP [51]	Median Binary Pattern	256
24 MRELBP [67]	Median Robust Extended Local Binary Pattern	800
25 MTP [51]	Median Ternary Pattern	512
26 PHOG [70]	Pyramid of Histogram of Oriented Gradients	168 *
27 WLD [73; 74]	Weber Local Descriptor	32 *

\* the feature dimension used in our experiments

330 three frames of each sequence and their landmarks provided are used for experiments. JAFFE and  
 331 TFEID are two acted face databases with six prototypical expressions and the neutral expression, which  
 332 contain 213 images from 10 Japanese females and 268 images from 40 Taiwanese subjects,  
 333 respectively. The BAUM-2 database consists of expression videos extracted from movies. The  
 334 expressions in the videos are in the close-to-real-life conditions, i.e. with pose, age, and illumination  
 335 variations. In our experiments, the image dataset, namely BAUM-2i, consisting of images with peak  
 336 expressions extracted from the videos in BAUM-2 are considered. There are 1,057 face images from  
 337 250 subjects, which have seven discrete expressions and the neutral expression in BAUM-2i.

338 The abovementioned databases have their own characteristics in terms of where the expression  
 339 images were taken, the expression classes, the race, age and gender of the participants, etc.

340 **3.1.2. Descriptors**

341 From the list of descriptors in [55], we select those descriptors based on spatial features, because this

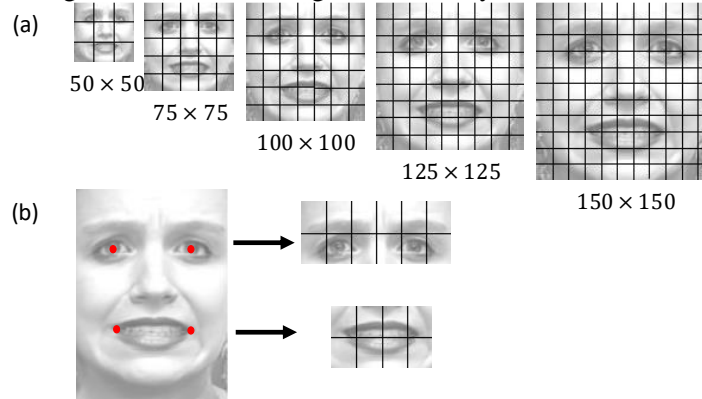
Table 1. A list of selected descriptors for our experiments and a comparison of the types of input data used and the local variation coding methods.

Descriptors	Input for local variation coding	Local variation coding
IWBC	Weber magnitude, Weber orientation	Local Xor Pattern (LXP) and Local Binary Pattern (LBP)
LAP	Intensity	first- and second-order derivatives using a set of custom filters
LBP	Intensity	-
LGBPHS	Gabor image	Local Binary Pattern (LBP)
LGIP	Intensity	Horizontal and vertical responses of sobel masks
LMP	Intensity	Local And Pattern with sign information of two level intensity differences
LPQ	Phase from Fourier Transform	Quantization
LTeP	Intensity	Two-level LBP
MBC_P	Phase, orientation or amplitude from Riesz transform	Local Nand (not and) Pattern
WLD	Differential excitations and orientations	Quantization

342 paper considers FER on static images only. The descriptors described in this paper and used in FER are  
 343 also included in the experiments. Because of numerous LBP variants, only the basic LBP variants and  
 344 MRELBP [67], which achieved the best performances in a recent comparative study for texture  
 345 classifications [84], are chosen for our comparative analysis. The other local descriptors, which are not  
 346 based on LBP, but inspired by LBP, for facial expression recognition are also included. Most of the  
 347 descriptors presented and evaluated in this paper belong to the sixth category defined in [84], which is  
 348 called “other methods inspired by LBP”. A feature learning method, named “Compact Binary Face  
 349 Descriptor (CBFD) [57]”, is also used in our experiments to evaluate its performance on FER, in  
 350 comparison to other state-of-the-art methods.

351 The descriptors (represented by their abbreviations), evaluated in our experiments, are listed in  
 352 **Error! Reference source not found..** To conduct a more detailed performance analysis, the best ten  
 353 descriptors, along with the corresponding input information used and coding methods, are listed in  
 354 Table 2. It is worth noting that MRELBP, IWBC, and CBFD are applied for the first time to the FER  
 355 problem.

Figure 1. Examples of the sub-regions used in our experiments: (a) regular sub-regions in an image, and (b) the sub-regions for the eye window and mouth window.



### 3.1.3. Pre-processing and feature extraction

For the first set of experiments, face images from the different databases are scaled to different resolutions, including  $50 \times 50$ ,  $75 \times 75$ ,  $100 \times 100$ ,  $125 \times 125$ , and  $150 \times 150$ . Then, features are extracted from the images with different numbers of sub-regions.

Table 2. The recognition rates for different resolutions, different numbers of sub-regions, on the CK+ database. “-” means that the corresponding results are unavailable because the dimensionality of the feature vectors are too high for experiments.

Database	Resolution	CK+ – LOSO – 6-class														
		50x50			75x75			100x100			125x125			150x150		
		3x3	3x3	5x5	5x5	7x7	5x5	7x7	9x9	5x5	7x7	9x9	11x11			
1	BPPC [2]	85.33	85.76	90.40	89.75	<b>90.83</b>	90.40	90.40	89.86	87.06	90.40	89.21	89.86			
2	GDP [5; 6]	74.54	75.73	83.82	86.30	86.62	86.30	85.98	<b>86.84</b>	85.65	85.76	86.08	86.08			
3	GDP2 [9; 10; 11]	57.71	57.39	83.06	81.98	90.51	82.20	89.97	92.45	83.06	89.43	94.28	<b>94.82</b>			
4	GLTP [13]	82.85	85.98	91.15	92.66	92.66	92.34	91.69	91.05	92.13	<b>93.64</b>	91.59	92.99			
5	IWBC [18]	88.67	90.72	91.69	90.51	93.42	91.15	93.10	<b>94.82</b>	90.83	92.13	93.31	92.99			
6	LAP [21]	83.17	80.26	89.75	89.21	91.69	90.29	91.26	93.42	91.05	91.05	93.42	<b>94.17</b>			
7	LBP [24]	84.68	84.03	92.23	91.48	91.69	91.91	93.53	93.85	91.80	93.20	93.74	<b>95.25</b>			
8	LDiP [26]	68.72	71.52	86.73	85.44	89.00	86.08	89.54	89.54	84.68	89.64	89.32	<b>89.75</b>			
9	LDiPv [30; 31]	68.93	71.20	83.17	82.85	86.95	83.50	87.70	89.00	85.11	85.98	88.67	<b>89.21</b>			
10	LDN [33]	80.91	82.96	88.46	88.24	90.29	90.40	91.15	92.66	90.83	90.40	<b>92.66</b>	91.91			
11	LDTF [34]	82.74	80.69	85.87	85.65	90.08	86.19	89.75	93.10	83.60	87.06	<b>93.53</b>	89.75			
12	LFD [37]	86.62	82.09	<b>90.61</b>	88.78	90.51	87.38	89.43	89.21	86.62	88.57	88.78	87.49			
13	LGBPHS [38]	86.19	87.27	92.02	92.88	92.99	91.26	90.72	91.48	90.29	89.75	91.48	<b>95.25</b>			
14	LGDIP [39]	71.09	69.15	75.19	80.15	79.72	77.35	78.86	79.07	80.04	<b>83.39</b>	80.80	78.64			
15	LGIP [40]	83.28	84.14	93.20	91.59	92.88	91.69	92.66	93.96	91.69	92.34	93.31	<b>95.15</b>			
16	LGP [41]	50.70	51.13	79.50	77.13	87.38	76.27	85.33	92.45	76.27	86.41	93.10	<b>93.31</b>			
17	LGTfP [42]	48.76	50.16	62.46	64.51	68.72	65.26	64.40	66.67	62.03	<b>69.26</b>	64.40	68.82			
18	LMP [45]	86.30	87.38	90.83	92.23	92.34	92.99	95.04	<b>95.25</b>	91.59	94.50	93.85	93.96			
19	LPQ [48; 49]	90.08	92.45	93.96	<b>94.39</b>	93.31	93.31	94.28	94.17	92.77	93.74	93.74	93.64			
20	LTfP [51]	88.35	89.10	91.80	92.45	93.31	92.99	93.96	<b>95.69</b>	92.56	94.50	95.04	94.93			
21	LTrP [59; 60]	74.76	75.73	85.65	85.44	88.13	84.36	87.70	88.24	87.38	<b>89.54</b>	89.43	87.27			
22	MBC_A [43; 61]	<b>92.56</b>	89.54	89.97	90.51	88.35	89.43	89.43	-	90.08	89.32	-	-			
23	MBC_P [43; 61]	88.89	89.32	92.88	<b>94.28</b>	90.51	91.80	93.42	-	92.56	92.45	-	-			
24	MBC_O [43; 61]	88.89	87.81	92.02	91.80	91.37	90.94	<b>92.56</b>	-	<b>92.56</b>	92.02	-	-			
25	MBP [51]	83.71	82.85	90.08	90.61	90.94	91.05	91.48	93.53	90.40	91.69	93.42	<b>94.07</b>			
26	MRELBP [67]	87.70	88.13	92.13	90.29	92.45	90.72	92.02	<b>93.53</b>	91.05	92.88	92.88	93.31			
27	MTP [51]	90.72	87.92	90.51	90.08	89.97	89.64	89.97	<b>92.77</b>	89.43	89.00	91.59	90.94			
28	PHOG [70]	87.59	89.54	89.54	90.29	89.32	89.00	<b>91.80</b>	90.72	89.21	90.83	90.51	90.40			
29	WLD [73; 74]	81.45	79.61	91.37	90.94	92.23	90.83	93.10	<b>95.90</b>	92.23	93.31	95.47	95.47			

In the second set of experiments, face images from the different databases are all scaled to the size of  $126 \times 189$  pixels, with a distance of 64 pixels between the two eyes. To locate the eye and mouth windows, the facial landmarks, i.e. the eye and mouth corners, are used. If facial landmarks are not

363 provided for a database, the required facial-feature points are marked manually. The eye window and  
 364 the mouth window are further divided into 12 and 8 sub-regions, respectively. Figure 1 shows examples  
 365 of selected sub-regions in both the first and the second set of experiments.

Table 3. The recognition rates for different resolutions and different numbers of sub-regions, on the BAUM-2i database. “-” means that the corresponding results are unavailable because the dimensionality of the feature vectors are too high for experiments.

Database Resolution # of sub-regions		BAUM-2i – 10-fold – 6-class												
		50x50		75x75			100x100		125x125			150x150		
		3x3	3x3	5x5	5x5	7x7	5x5	7x7	9x9	5x5	7x7	9x9	11x11	
1	BPPC [2]	51.18	52.48	53.90	56.26	58.98	54.37	<b>59.10</b>	57.45	56.15	55.67	57.21	54.85	
2	GDP [5; 6]	40.78	45.39	50.71	46.10	<b>52.84</b>	45.98	50.83	51.89	46.22	50.24	51.89	49.76	
3	GDP2 [9; 10; 11]	23.88	26.12	29.91	27.90	37.35	28.01	40.54	48.94	28.84	40.66	48.11	<b>53.31</b>	
4	GLTP [13]	50.00	54.14	57.57	55.44	59.57	53.90	59.10	58.27	54.73	59.22	57.57	<b>60.05</b>	
5	IWBC [18]	55.67	57.33	<b>59.22</b>	58.39	58.16	56.97	57.92	57.57	57.09	57.21	56.86	56.26	
6	LAP [21]	46.22	44.80	53.66	48.70	51.77	48.35	50.24	56.03	49.05	50.24	56.03	<b>56.38</b>	
7	LBP [24]	48.35	48.23	54.26	54.02	56.62	53.07	56.86	58.63	53.31	54.61	56.62	<b>59.46</b>	
8	LDiP [26]	27.30	30.61	43.62	45.39	53.07	48.70	52.25	53.66	45.15	52.60	<b>56.03</b>	<b>56.03</b>	
9	LDiPv [30; 31]	24.11	28.25	40.19	36.05	49.88	40.54	48.94	52.25	40.78	48.82	52.84	<b>53.90</b>	
10	LDN [33]	37.23	34.16	48.11	47.52	51.06	47.28	54.14	55.67	47.16	54.61	55.91	<b>60.99</b>	
11	LDTF [34]	30.26	34.04	48.11	43.38	48.82	42.79	46.81	49.41	44.33	47.64	46.81	<b>51.06</b>	
12	LFD [37]	46.69	44.56	53.90	50.59	57.21	51.77	56.74	57.57	51.06	54.73	56.50	<b>57.92</b>	
13	LGBPHS [38]	49.41	50.00	56.62	57.57	59.46	59.22	60.28	60.76	59.81	61.11	<b>62.41</b>	57.92	
14	LGDIP [39]	30.97	32.62	39.60	42.67	44.09	41.25	46.10	<b>47.52</b>	39.83	42.55	44.44	43.85	
15	LGIP [40]	30.50	30.97	49.88	47.04	54.61	49.17	54.02	56.74	48.11	52.96	55.44	<b>58.39</b>	
16	LGP [41]	23.40	21.75	23.52	26.36	31.32	25.41	32.98	42.43	25.30	30.02	41.84	<b>46.22</b>	
17	LGTfP [42]	24.47	23.05	31.44	30.38	31.09	32.03	34.04	35.11	31.68	36.29	<b>40.07</b>	36.29	
18	LMP [45]	49.76	50.35	55.91	56.86	58.75	56.50	58.16	<b>60.64</b>	52.36	55.32	58.27	60.17	
19	LPQ [48; 49]	56.38	56.03	61.35	61.35	60.28	59.46	<b>61.47</b>	61.23	57.68	59.57	60.28	<b>61.47</b>	
20	LTeP [51]	52.36	50.59	55.32	52.96	58.87	52.96	59.46	59.22	54.02	59.57	<b>60.28</b>	<b>60.28</b>	
21	LTrP [59; 60]	35.34	38.89	42.79	46.22	51.06	45.15	49.17	52.01	46.57	50.47	51.77	<b>53.78</b>	
22	MBC_A [43; 61]	56.62	56.62	<b>59.81</b>	57.57	59.34	56.38	58.63	55.08	56.03	55.67	55.20	-	
23	MBC_P [43; 61]	56.03	54.96	59.93	59.81	61.58	59.57	61.47	61.94	59.46	60.87	<b>62.06</b>	-	
24	MBC_O [43; 61]	57.68	55.79	61.35	<b>61.82</b>	60.99	60.05	60.17	61.23	58.27	60.99	60.40	-	
25	MBP [3; 4; 51]	43.62	47.04	54.37	54.37	54.49	53.55	55.32	<b>59.46</b>	52.60	53.43	55.08	57.33	
26	MRELBP [67]	46.34	48.70	55.56	56.86	57.68	57.80	57.92	<b>59.34</b>	57.45	57.57	58.98	<b>59.34</b>	
27	MTP [51]	43.97	41.96	51.54	50.24	<b>54.02</b>	47.87	53.78	51.65	42.91	51.65	52.13	52.60	
28	PHOG [70]	47.52	50.35	51.42	53.43	53.90	51.77	54.26	52.96	54.14	54.02	<b>54.61</b>	53.43	
29	WLD [73; 74]	30.26	24.82	51.77	46.10	57.80	47.52	55.44	56.15	46.93	54.73	55.79	<b>58.75</b>	

### 366 3.1.4. Dimensionality reduction and classification

367 In the first two sets of experiments, the local descriptors listed in **Error! Reference source not found.**  
 368 were first extracted. Then, the subspace-learning method, Soft Locality Preserving Projection (SLPM)  
 369 [98], is applied for manifold learning and dimensionality reduction. SLPM is a graph-based subspace-  
 370 learning method, which uses the  $k$ -neighborhood information and the class information. The key feature  
 371 of SLPM is that it aims to control the level of spread of the different classes, because the spread of the  
 372 classes in the underlying manifold is closely connected to the generalizability of the learned subspace.  
 373 In our experiments, we employ SLPM for dimensionality reduction and for increasing the  
 374 discriminative ability of the extracted features. Finally, the nearest neighbour (NN) classifier is used for  
 375 classification. The third set of experiments were conducted, with the best setting for each of the  
 376 databases, using the Support Vector Machine (SVM) classifier, with the linear kernel. The results are



377 then compared to those based on the nearest neighbour classifier. Two different cross-validation  
 378 schemes are adopted in our experiments: Leave-One-Subject-Out (LOSO) to encourage the  
 379 reproducibility of the experiments, and 10-fold cross-validation, which is used when there are sufficient  
 380 number of images for each subject in the database, i.e. BAUM-2i. Furthermore, both the 10-fold and  
 381 LOSO cross-validation schemes are used for comparison on the JAFFE and TFEID databases.

### 382 3.2. Experimental results

Table 4. The comparison of recognition rates obtained by the selected local descriptors on the BAUM-2i database (the best of sub-regions) using 10-fold cross validation. 6-class: AN, DI, FE, HA, SA, and SU. 7-class: AN, CO, DI, FE, HA, SA, and SU. 8-class: AN, CO, DI, FE, HA, NE, SA, and SU.

	BAUM-2i		
	6-class	7-class	8-class
IWBC	59.22	55.53	52.53
LAP	56.38	54.97	49.00
LBP	59.46	58.32	52.44
LGBPHS	<b>62.41</b>	57.99	54.15
LGIP	58.39	54.75	49.86
LMP	60.64	57.54	52.53
LPQ	61.47	<b>58.99</b>	<b>54.73</b>
LT <sub>e</sub> P	60.28	57.21	52.63
MBC_P	62.06	58.10	54.25
WLD	58.75	54.41	50.53

383 In this section, the experiment results on the four facial-expression databases (BAUM-2i, CK+, JAFFE,  
 384 TFEID) under different experimental settings are presented and discussed. The experiments are  
 385 designed to measure the performances of the respective descriptors, for face images at different  
 386 resolutions and divided into different sub-regions, and with different classifiers.

#### 387 3.2.1. Performance analysis for varying resolution and number of sub-regions

388 All the face images are first aligned based on the positions of the two eye pupils, and cropped to the  
 389 different resolutions. For each resolution, face images are divided into different numbers of sub-regions,  
 390 say  $l \times l$ , where  $l$  varies from 3 to 11.

391 Table 3 and Table 4 present the results on CK+ and BAUM-2i for all the descriptors. As observed  
 392 from the results shown in Tables 3 and 4, in general, the classification performances improve when the  
 393 image resolution and the number of sub-regions increase. Therefore, higher resolution and more sub-  
 394 regions lead to better classification performances. However, with more sub-regions, the feature  
 395 dimension will become very high. In other words, the better performance is at the expenses of higher  
 396 computational requirements.

397 For more detailed performance analysis, the best ten descriptors, which have achieved promising  
398 results, were chosen to repeat the first set of experiments on the four databases separately with different  
399 numbers of expression classes, as well as the two different classification schemes. Table 5 shows the  
400 best classification rates on BAUM-2i with different numbers of expression classes. In Tables 6 to 8, the  
401 columns named “best of sub-regions” show the best classification rates for the number of sub-regions  
402 being used. We only show the best results, otherwise there are too many data to be shown.

### 403 **3.2.2. Performance analysis of the eye and mouth regions**

404 The second set of experiments was conducted with the features extracted from the eye and mouth  
405 windows of face images. The CK+, JAFFE and TFEID databases are used to test the performances of  
406 the respective features extracted from the eye and the mouth regions. The BAUM-2i database is not  
407 used because it consists of images in the wild. Labelling the facial landmarks is a complicated task. In  
408 Tables 6 to 8, the two columns under “eye and mouth windows” show the classification accuracies of  
409 the selected features, using the LOSO and 10-fold cross-validation schemes.

410 As observed from the tables, using the features extracted from the eye and the mouth windows  
411 achieves lower classification accuracies than that using features extracted from the sub-regions of whole  
412 face images. However, for the results based on sub-regions, we show the best classification accuracies  
413 achieved for the five different resolutions and the five different numbers of sub-regions. Furthermore,  
414 each descriptor achieves the best performance on a different resolution and a different number of sub-

415 regions. Experiment results show that there are not a particular resolution and a particular number of

Table 5. The recognition rates of selected local descriptors on the CK+ database, with 6 classes (AN, DI, FE, HA, SA, and SU) and 7 classes (AN, CO, DI, FE, HA, SA, and SU), using LOSO.

	CK+			
	Eye and mouth windows		Best of sub-regions	
	6-class	7-class	6-class	7-class
IWBC	94.61	93.68	94.82	94.50
LAP	91.37	91.44	94.17	92.86
LBP	93.31	92.56	95.25	93.99
LGBPHS	92.23	90.72	95.25	93.99
LGIP	91.26	92.35	95.15	94.50
LMP	<b>94.71</b>	94.19	95.25	<b>94.90</b>
LPQ	94.61	<b>94.90</b>	94.39	94.19
LTeP	93.53	93.17	95.69	94.80
MBC_P	91.69	89.40	94.28	92.46
WLD	93.31	91.44	<b>95.90</b>	94.80

Table 6. The recognition rates of the selected best local descriptors on the JAFFE database.

	JAFFE			
	Eye and mouth windows		Best of sub-regions	
	LOSO	10-fold	LOSO	10-fold
IWBC	58.69	88.73	65.73	90.61
LAP	<b>68.08</b>	90.61	68.54	94.84
LBP	61.50	86.38	65.73	93.43
LGBPHS	63.38	<b>93.90</b>	<b>71.83</b>	93.90
LGIP	62.91	87.32	66.20	93.90
LMP	60.09	85.92	67.14	93.43
LPQ	67.61	92.02	69.95	93.43
LTeP	61.03	89.20	62.44	94.37
MBC_P	63.38	92.96	66.67	93.90
WLD	63.38	86.85	69.01	<b>96.24</b>
CBFD	66.20	89.67	-	-

Table 7. The comparison of recognition rates obtained by the selected local descriptors on the TFEID database.

	TFEID			
	Eye and mouth windows		Best of sub-regions	
	LOSO	10-fold	LOSO	10-fold
IWBC	89.55	90.67	92.91	91.79
LAP	91.04	91.04	94.40	<b>95.15</b>
LBP	91.79	92.54	93.66	94.78
LGBPHS	<b>94.40</b>	91.04	<b>95.15</b>	93.66
LGIP	89.18	86.19	94.78	93.28
LMP	91.42	92.16	94.03	94.03
LPQ	<b>94.40</b>	<b>93.28</b>	94.03	94.40
LTeP	90.30	92.16	94.40	<b>95.15</b>
MBC_P	94.30	91.79	94.40	93.66
WLD	92.16	91.42	94.78	94.40
CBFD	93.66	92.16	-	-

416 sub-regions that can work the best for all the descriptors.

### 417 3.2.3. Performance analysis of the classifiers

418 Table 9 presents the experiment results obtained with the NN and the SVM classifiers. We can observe  
 419 that both LGBPHS and LPQ achieve similar performances in the use of NN and SVM. However, the

420 NN classifier can achieve equal or higher performance than the SVM classifier if a supervised  
 421 dimensionality reduction method is employed. In our experiments, we utilize SLPM for dimensionality

Table 8. The comparison of the recognition rates obtained with features extracted from the eye and mouth regions by the nearest neighbor classifier (NN) and SVM classifier using LOSO.

	CK+		JAFPE		TFEID	
	SLPM + NN	SVM	SLPM + NN	SVM	SLPM + NN	SVM
LGBPMS	92.23	91.91	63.38	61.50	94.40	94.40
LPQ	94.61	<b>94.93</b>	<b>67.61</b>	67.14	94.40	94.40

422 reduction.

### 423 3.2.4. Performance analysis of cross-dataset facial expression recognition

424 In real-life applications, query samples are often different from the training samples in terms of  
 425 uncontrolled variations such as illumination. Therefore, it is important for a local descriptor to have a  
 426 good generalization power, and the descriptor can still achieve a good performance when the training  
 427 and test sets are from different databases. In this paper, we also conduct experiments to test the  
 428 robustness and accuracy of the best selected descriptors in the scenario of cross-dataset FER.

429 Table 10 shows the experiment results when the training and the testing sets are two different datasets,  
 430 which have different acquisition conditions. As you can observe in Table 10, the recognition rates for

Table 9. The comparison of the recognition rates of the ten selected descriptors on cross-dataset facial expression recognition, with features extracted from the eye and mouth windows.

Trained on	CK+		JAFPE		TFEID	
	JAFPE	TFEID	CK+	TFEID	CK+	JAFPE
IWBC	25.00	33.77	34.52	42.98	38.30	30.98
LAP	29.89	32.46	24.16	42.98	<b>45.85</b>	26.63
LBP	21.74	34.65	29.02	44.74	35.81	28.26
LGBPMS	18.48	33.33	37.22	60.09	39.48	44.02
LGIP	<b>30.43</b>	31.14	31.18	41.23	42.61	31.52
LMP	29.35	32.46	<b>37.32</b>	48.25	37.32	23.37
LPQ	19.57	<b>38.16</b>	32.58	50.44	42.07	35.33
LTep	19.57	35.96	25.03	25.88	26.86	35.87
MBC_P	25.54	31.14	37.00	<b>63.60</b>	38.83	<b>47.28</b>
WLD	15.76	35.09	27.18	32.89	42.61	24.46

431 the 6 basic emotions decrease significantly, because cross-dataset FER is a challenging task. Although  
 432 no local descriptor can perform consistently better than the others, MBC\_P achieves the highest  
 433 recognition rates when the model is trained using JAFPE while tested on TFEID, and vice versa.  
 434 MBC\_P uses monogenic signal analysis to estimate the phase component of the images, which  
 435 represents the images' geometric information. Since the JAFPE database consists of images of Japanese  
 436 women, while TFEID consists of images of Taiwanese men and women, we can observe that the phase  
 437 information of the monogenic signal analysis is insensitive to cross-cultural face representation for FER.

### 438 3.2.5. Comparison with deep features

439 Recently, convolutional deep neural networks have been applied to FER [1; 3; 4; 12; 17; 20; 25]. Table  
 440 11 presents the performances of deep learning methods applied on the CK+ database. 3DCNN-DAP [1]  
 441 adapts a deformable parts learning component to detect discriminative facial action parts for  
 442 spatiotemporal FER, where a Boosted Deep Belief Network [3] (BDBN) is used to learn and select the  
 443 expression-related facial features to develop a strong classifier in a unified loopy framework iteratively.  
 444 Iterative learning of the BDBN framework strengthens the discriminative capabilities of the features.  
 445 STM-ExpLet [4] learns a spatiotemporal manifold (STM) from low-level features from each expression  
 446 video clip, followed by learning a universal manifold model that statistically unify all the STMs. With  
 447 this method, expression videos are also aligned. Different from these methods, DTAGN [12] trains two  
 448 models, with temporal geometry features and temporal appearance features, respectively, from image

Table 10. The comparison of recognition rates of deep learning methods and the best recognition rate obtained with handcrafted features

Method	Feature Type	Accuracy (%)
3DCNN-DAP [1]	Deep features	92.4
BDBN [3]	Deep features	96.7
STM-ExpLet [4]	Deep features	94.2
DTAGN [12]	Deep features	97.3
Inception [17]	Deep features	93.2
PPDN [20]	Deep features	97.3
LFC + FFD [23]	Handcrafted features	97.9
FN2EN [25]	Deep features	<b>98.6</b>
<b>LPQ-SLPM-NN</b>	Handcrafted features	<b>95.9</b>

449 sequences, and these two features are complementary to each other. In [17], a network, which consists  
 450 of two convolutional layers with max pooling and four inception layers, was proposed. The network  
 451 was evaluated for its generalizability by experiments, with cross-database classification. To boost the  
 452 generalizability of learning, [20] presented a peak-piloted deep network (PPDN), which uses the  
 453 samples with high-intensity expressions to supervise the samples with low-intensity expression that are  
 454 hard to classify. Until now, FN2EN [25], which uses a two-stage training algorithm, achieved the best  
 455 performance on the CK+ dataset. FN2EN, in the first stage, trains the convolutional layers, whose  
 456 outputs from the last pooling layer are used to supervise the expression net in the second stage.

457 As observed in Table 11, LPQ with NN outperforms several deep learning methods. LFC + FFD  
 458 [23] is also a histogram-based feature extraction method, which achieves higher classification accuracy

459 than all the listed methods, except FN2EN [25]. To the best of our knowledge, FN2EN achieves the  
460 highest classification accuracy on the CK+ database. However, expensive computational cost is a  
461 drawback of most of the methods based on deep convolutional neural networks. Furthermore, the CK+  
462 database consists of images taken under controlled environments, i.e. posed expressions and the number  
463 of expression samples are limited in the CK+ database. These factors direct us the need of a large-scale  
464 facial-expression database in the wild. There have been several attempts to collect facial-expression  
465 images in the wild [99; 100; 101]. [102] and [103] are two recently published databases, which contain  
466 large-scale face images with varying expressions. These databases will be very useful for FER based  
467 on deep learning.

#### 468 **4. Conclusion**

469 This paper provides a systematic review and analysis of current histogram-based local feature  
470 descriptors, which have been applied for facial-expression recognition. The weaknesses and strengths  
471 of the existing descriptors, as well as their underlying connections, have also been discussed and  
472 analysed. Then, a comprehensive evaluation of the performances of different descriptors for facial-  
473 expression recognition is conducted and presented. In total, 27 local descriptors have been applied on  
474 four facial-expression databases, under the same experimental settings. The robustness of the respective  
475 local descriptors is tested under different conditions, such as varying image resolutions and number of  
476 sub-regions, and the classifiers. Moreover, a brief performance comparison with seven recent deep  
477 features and two handcrafted features has been conducted.

478 Several remarks from the experiment results are listed as follows:

- 479 • The databases have different characteristics, which affect the choice of the ideal descriptor for  
480 a particular database. Even the number of expression classes can also affect the performances  
481 of the descriptors.
- 482 • The results show a trade-off between the number of sub-regions and the overall classification  
483 accuracy. The use of the eye and the mouth windows decreases the number of sub-regions and  
484 the dimensionality of the resulting feature vectors, with a slight loss in terms of accuracy.

- 485       • The resolution of face images and the number of sub-regions are the two most important factors  
486       that affect the overall classification accuracies.
- 487       • The highest classification accuracies are obtained mostly by LGBPHS and LPQ. This shows  
488       that Gabor wavelets and phase information are important features for representing expression-  
489       specific information. However, we should keep in mind that Gabor features suffer from high  
490       computational cost.
- 491       • According to the comprehensive analysis shown in this paper, the best local descriptors for  
492       FER, by considering the feature length, computational cost, and the classification accuracy  
493       simultaneously, is LPQ.
- 494       • Deep neural-network-based methods indeed can achieve excellent classification accuracies on  
495       FER. However, these methods also suffer from time and space complexities as LGBPHS.

496 In conclusion, our comprehensive experiment results show that the trade-off between the computational  
497 cost and the classification accuracy still exists today.

## 498 **5. Acknowledgements**

499 The work described in this paper was supported by a research grant from The Hong Kong Polytechnic  
500 University (project code: G-YBKF).

## 501 **References**

- 502 [1] M. Liu, S. Li, S. Shan, R. Wang, and X. Chen, Deeply learning deformable facial action parts model  
503       for dynamic expression analysis, Asian Conference on Computer Vision, Springer, 2014, pp.  
504       143-157.
- 505 [2] S. Shojaeilangari, W.-Y. Yau, J. Li, and E.-K. Teoh, Feature extraction through binary pattern of  
506       phase congruency for facial expression recognition, Control Automation Robotics & Vision  
507       (ICARCV), 2012 12th International Conference on, IEEE, 2012, pp. 166-170.
- 508 [3] P. Liu, S. Han, Z. Meng, and Y. Tong, Facial expression recognition via a boosted deep belief  
509       network, Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition,  
510       2014, pp. 1805-1812.
- 511 [4] M. Liu, S. Shan, R. Wang, and X. Chen, Learning expressionlets on spatio-temporal manifold for  
512       dynamic facial expression recognition, Proceedings of the IEEE Conference on Computer  
513       Vision and Pattern Recognition, 2014, pp. 1749-1756.
- 514 [5] F. Ahmed, Gradient directional pattern: a robust feature descriptor for facial expression recognition.  
515       Electronics letters 48 (2012) 1203-1204.
- 516 [6] W. Chu, Facial expression recognition based on local binary pattern and gradient directional pattern,  
517       Green Computing and Communications (GreenCom), 2013 IEEE and Internet of Things  
518       (iThings/CPSCoM), IEEE International Conference on and IEEE Cyber, Physical and Social  
519       Computing, IEEE, 2013, pp. 1458-1462.
- 520 [7] P. Ekman, and W.V. Friesen, Constants across cultures in the face and emotion. Journal of  
521       personality and social psychology 17 (1971) 124.

- 522 [8] P. Ekman, Darwin, deception, and facial expression. *Annals of the New York Academy of Sciences*  
523 1000 (2003) 205-221.
- 524 [9] M.S. Islam, Gender Classification using Gradient Direction Pattern. *Science International* 25 (2013).
- 525 [10] M. Valstar, and M. Pantic, Fully automatic facial action unit detection and temporal analysis,  
526 *Computer Vision and Pattern Recognition Workshop, 2006. CVPRW'06. Conference on, IEEE,*  
527 2006, pp. 149-149.
- 528 [11] R. Walecki, V. Pavlovic, B. Schuller, and M. Pantic, Deep Structured Learning for Facial Action  
529 Unit Intensity Estimation. *arXiv preprint arXiv:1704.04481* (2017).
- 530 [12] H. Jung, S. Lee, S. Park, I. Lee, C. Ahn, and J. Kim, Deep temporal appearance-geometry network  
531 for facial expression recognition. *arXiv preprint arXiv:1503.01532* (2015).
- 532 [13] F. Ahmed, and E. Hossain, Automated facial expression recognition using gradient-based ternary  
533 texture patterns. *Chinese Journal of Engineering* 2013 (2013).
- 534 [14] C. Turan, K.-M. Lam, and X. He, Facial expression recognition with emotion-based feature fusion,  
535 *Signal and Information Processing Association Annual Summit and Conference (APSIPA),*  
536 2015 Asia-Pacific, IEEE, 2015, pp. 1-6.
- 537 [15] C. Turan, and K.-M. Lam, Region-based feature fusion for facial-expression recognition, *Image*  
538 *Processing (ICIP), 2014 IEEE International Conference on, IEEE, 2014, pp. 5966-5970.*
- 539 [16] C. Shan, S. Gong, and P.W. McOwan, Facial expression recognition based on local binary patterns:  
540 A comprehensive study. *Image and Vision Computing* 27 (2009) 803-816.
- 541 [17] A. Mollahosseini, D. Chan, and M.H. Mahoor, Going deeper in facial expression recognition using  
542 deep neural networks, *Applications of Computer Vision (WACV), 2016 IEEE Winter*  
543 *Conference on, IEEE, 2016, pp. 1-10.*
- 544 [18] B.-Q. Yang, T. Zhang, C.-C. Gu, K.-J. Wu, and X.-P. Guan, A novel face recognition method  
545 based on iwld and iwbc. *Multimedia Tools and Applications* 75 (2016) 6979.
- 546 [19] E. Vural, M. Cetin, A. Ercil, G. Littlewort, M. Bartlett, and J. Movellan, Drowsy driver detection  
547 through facial movement analysis. *Human-computer interaction* (2007) 6-18.
- 548 [20] X. Zhao, X. Liang, L. Liu, T. Li, Y. Han, N. Vasconcelos, and S. Yan, Peak-piloted deep network  
549 for facial expression recognition, *European Conference on Computer Vision, Springer, 2016,*  
550 pp. 425-442.
- 551 [21] M.S. Islam, and S. Auwatanamongkol, Facial Expression Recognition using Local Arc Pattern.  
552 *Trends in Applied Sciences Research* 9 (2014) 113.
- 553 [22] T. Fong, I. Nourbakhsh, and K. Dautenhahn, A survey of socially interactive robots. *Robotics and*  
554 *autonomous systems* 42 (2003) 143-166.
- 555 [23] G. Benitez-Garcia, T. Nakamura, and M. Kaneko, Facial Expression Recognition Based on Local  
556 Fourier Coefficients and Facial Fourier Descriptors. *Journal of Signal and Information*  
557 *Processing* 8 (2017) 132.
- 558 [24] T. Ojala, M. Pietikainen, and T. Maenpaa, Multiresolution gray-scale and rotation invariant texture  
559 classification with local binary patterns. *IEEE Transactions on pattern analysis and machine*  
560 *intelligence* 24 (2002) 971-987.
- 561 [25] H. Ding, S.K. Zhou, and R. Chellappa, Facenet2expnet: Regularizing a deep face recognition net  
562 for expression recognition, *Automatic Face & Gesture Recognition (FG 2017), 2017 12th IEEE*  
563 *International Conference on, IEEE, 2017, pp. 118-126.*
- 564 [26] T. Jabid, M.H. Kabir, and O. Chae, Local directional pattern (LDP)—A robust image descriptor for  
565 object recognition, *Advanced Video and Signal Based Surveillance (AVSS), 2010 Seventh*  
566 *IEEE International Conference on, IEEE, 2010, pp. 482-487.*
- 567 [27] J.F. Cohn, T.S. Kruez, I. Matthews, Y. Yang, M.H. Nguyen, M.T. Padilla, F. Zhou, and F. De la  
568 Torre, Detecting depression from facial actions and vocal prosody, *Affective Computing and*  
569 *Intelligent Interaction and Workshops, 2009. ACII 2009. 3rd International Conference on,*  
570 *IEEE, 2009, pp. 1-7.*
- 571 [28] S. Jaiswal, M.F. Valstar, A. Gillott, and D. Daley, Automatic detection of ADHD and ASD from  
572 expressive behaviour in RGBD data, *Automatic Face & Gesture Recognition (FG 2017), 2017*  
573 *12th IEEE International Conference on, IEEE, 2017, pp. 762-769.*
- 574 [29] S.J. Kirsh, and J.R. Mounts, Violent video game play impacts facial emotion recognition.  
575 *Aggressive behavior* 33 (2007) 353-358.



- 576 [30] M.H. Kabir, T. Jabid, and O. Chae, A local directional pattern variance (LDPv) based face  
577 descriptor for human facial expression recognition, *Advanced Video and Signal Based*  
578 *Surveillance (AVSS)*, 2010 Seventh IEEE International Conference on, IEEE, 2010, pp. 526-  
579 532.
- 580 [31] K.-C. Fan, and T.-Y. Hung, A novel local pattern descriptor—local vector pattern in high-order  
581 derivative space for face recognition. *IEEE transactions on image processing* 23 (2014) 2877-  
582 2891.
- 583 [32] B. Martinez, M.F. Valstar, B. Jiang, and M. Pantic, Automatic analysis of facial actions: A survey.  
584 *IEEE Transactions on Affective Computing* (2017).
- 585 [33] A.R. Rivera, J.R. Castillo, and O.O. Chae, Local directional number pattern for face analysis: Face  
586 and expression recognition. *IEEE transactions on image processing* 22 (2013) 1740-1752.
- 587 [34] A.R. Rivera, J.R. Castillo, and O. Chae, Local directional texture pattern image descriptor. *Pattern*  
588 *Recognition Letters* 51 (2015) 94-100.
- 589 [35] S. Jaiswal, and M. Valstar, Deep learning the dynamic appearance and shape of facial action units,  
590 *Applications of Computer Vision (WACV)*, 2016 IEEE Winter Conference on, IEEE, 2016,  
591 pp. 1-8.
- 592 [36] Z. Tóssér, L.A. Jeni, A. Lőrincz, and J.F. Cohn, Deep learning for facial action unit detection under  
593 large head poses, *Computer Vision—ECCV 2016 Workshops*, Springer, 2016, pp. 359-371.
- 594 [37] Z. Lei, T. Ahonen, M. Pietikäinen, and S.Z. Li, Local frequency descriptor for low-resolution face  
595 recognition, *Automatic Face & Gesture Recognition and Workshops (FG 2011)*, 2011 IEEE  
596 *International Conference on*, IEEE, 2011, pp. 161-166.
- 597 [38] W. Zhang, S. Shan, W. Gao, X. Chen, and H. Zhang, Local gabor binary pattern histogram  
598 sequence (lgbphs): A novel non-statistical model for face representation and recognition,  
599 *Computer Vision, 2005. ICCV 2005. Tenth IEEE International Conference on*, IEEE, 2005, pp.  
600 786-791.
- 601 [39] S.Z. Ishraque, A.H. Banna, and O. Chae, Local Gabor directional pattern for facial expression  
602 recognition, *Computer and Information Technology (ICCIT)*, 2012 15th International  
603 *Conference on*, IEEE, 2012, pp. 164-167.
- 604 [40] Z. Lubing, and W. Han, Local gradient increasing pattern for facial expression recognition, *Image*  
605 *Processing (ICIP)*, 2012 19th IEEE International Conference on, IEEE, 2012, pp. 2601-2604.
- 606 [41] M.S. Islam, Local gradient pattern—A novel feature representation for facial expression recognition.  
607 *Journal of AI and Data Mining* 2 (2014) 33-38.
- 608 [42] T. Ahsan, T. Jabid, and U.-P. Chong, Facial expression recognition using local transitional pattern  
609 on Gabor filtered facial images. *IETE Technical Review* 30 (2013) 47-52.
- 610 [43] M. Yang, L. Zhang, S.C.-K. Shiu, and D. Zhang, Monogenic binary coding: An efficient local  
611 feature extraction approach to face recognition. *IEEE Transactions on Information Forensics*  
612 *and Security* 7 (2012) 1738-1751.
- 613 [44] J. Li, N. Sang, and C. Gao, Face recognition with Riesz binary pattern. *Digital Signal Processing*  
614 51 (2016) 196-201.
- 615 [45] T. Mohammad, and M.L. Ali, Robust facial expression recognition based on local monotonic  
616 pattern (LMP), *Computer and Information Technology (ICCIT)*, 2011 14th International  
617 *Conference on*, IEEE, 2011, pp. 572-576.
- 618 [46] R.A. Khan, A. Meyer, H. Konik, and S. Bouakaz, Framework for reliable, real-time facial  
619 expression recognition for low resolution images. *Pattern Recognition Letters* 34 (2013) 1159-  
620 1168.
- 621 [47] T. Song, H. Li, F. Meng, Q. Wu, and J. Cai, LETRIST: Locally Encoded Transform Feature  
622 Histogram for Rotation-Invariant Texture Classification. *IEEE Transactions on Circuits and*  
623 *Systems for Video Technology* (2017).
- 624 [48] V. Ojansivu, and J. Heikkilä, Blur insensitive texture classification using local phase quantization,  
625 *International conference on image and signal processing*, Springer, 2008, pp. 236-243.
- 626 [49] A. Dhalla, A. Asthana, R. Goecke, and T. Gedeon, Emotion recognition using PHOG and LPQ  
627 features, *Automatic Face & Gesture Recognition and Workshops (FG 2011)*, 2011 IEEE  
628 *International Conference on*, IEEE, 2011, pp. 878-883.
- 629 [50] G. Zhao, and M. Pietikäinen, Dynamic texture recognition using volume local binary patterns,  
630 *Dynamical Vision*, Springer, 2007, pp. 165-177.

- 631 [51] F. Bashar, A. Khan, F. Ahmed, and M.H. Kabir, Robust facial expression recognition based on  
632 median ternary pattern (MTP), Electrical Information and Communication Technology (EICT),  
633 2013 International Conference on, IEEE, 2014, pp. 1-5.
- 634 [52] L. Nanni, A. Lumini, and S. Brahnam, Local binary patterns variants as texture descriptors for  
635 medical image analysis. *Artificial intelligence in medicine* 49 (2010) 117-125.
- 636 [53] N.P. Doshi, and G. Schaefer, A comprehensive benchmark of local binary pattern algorithms for  
637 texture retrieval, *Pattern Recognition (ICPR)*, 2012 21st International Conference on, IEEE,  
638 2012, pp. 2760-2763.
- 639 [54] L. Nanni, A. Lumini, and S. Brahnam, Survey on LBP based texture descriptors for image  
640 classification. *Expert Systems with Applications* 39 (2012) 3634-3641.
- 641 [55] R.L. Kristensen, Z.-H. Tan, Z. Ma, and J. Guo, Binary pattern flavored feature extractors for Facial  
642 Expression Recognition: An overview, *Information and Communication Technology,  
643 Electronics and Microelectronics (MIPRO)*, 2015 38th International Convention on, IEEE,  
644 2015, pp. 1131-1137.
- 645 [56] V. Balntas, L. Tang, and K. Mikolajczyk, Binary online learned descriptors. *IEEE transactions on  
646 pattern analysis and machine intelligence* 40 (2018) 555-567.
- 647 [57] J. Lu, V.E. Liong, X. Zhou, and J. Zhou, Learning compact binary face descriptor for face  
648 recognition. *IEEE transactions on pattern analysis and machine intelligence* 37 (2015) 2041-  
649 2056.
- 650 [58] Y. Duan, J. Lu, J. Feng, and J. Zhou, Learning rotation-invariant local binary descriptor. *IEEE  
651 Transactions on Image Processing* 26 (2017) 3636-3651.
- 652 [59] T. Jabid, and O. Chae, Local Transitional Pattern: A Robust Facial Image Descriptor for Automatic  
653 Facial Expression Recognition, *Proc. International Conference on Computer Convergence  
654 Technology*, Seoul, Korea, 2011, pp. 333-44.
- 655 [60] T. Jabid, and O. Chae, Facial Expression Recognition Based on Local Transitional Pattern.  
656 *International Information Institute (Tokyo). Information* 15 (2012) 2007.
- 657 [61] X.X. Xia, Z.L. Ying, and W.J. Chu, Facial Expression Recognition Based on Monogenic Binary  
658 Coding, *Applied Mechanics and Materials*, Trans Tech Publ, 2014, pp. 437-440.
- 659 [62] J. Lu, V.E. Liong, and J. Zhou, Cost-sensitive local binary feature learning for facial age estimation.  
660 *IEEE Transactions on Image Processing* 24 (2015) 5356-5368.
- 661 [63] Y. Duan, J. Lu, J. Feng, and J. Zhou, Context-aware local binary feature learning for face  
662 recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2017).
- 663 [64] J. Lu, V. Erin Liong, and J. Zhou, Simultaneous local binary feature learning and encoding for face  
664 recognition, *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp.  
665 3721-3729.
- 666 [65] J. Lu, V.E. Liong, and J. Zhou, Simultaneous local binary feature learning and encoding for  
667 homogeneous and heterogeneous face recognition. *IEEE transactions on pattern analysis and  
668 machine intelligence* (2017).
- 669 [66] A. Toshev, and C. Szegedy, Deeppose: Human pose estimation via deep neural networks,  
670 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp.  
671 1653-1660.
- 672 [67] L. Liu, S. Lao, P.W. Fieguth, Y. Guo, X. Wang, and M. Pietikäinen, Median robust extended local  
673 binary pattern for texture classification. *IEEE Transactions on Image Processing* 25 (2016)  
674 1368-1381.
- 675 [68] F. Schroff, D. Kalenichenko, and J. Philbin, Facenet: A unified embedding for face recognition  
676 and clustering, *Proceedings of the IEEE Conference on Computer Vision and Pattern  
677 Recognition*, 2015, pp. 815-823.
- 678 [69] J. Mansanet, A. Albiol, and R. Paredes, Local deep neural networks for gender recognition. *Pattern  
679 Recognition Letters* 70 (2016) 80-86.
- 680 [70] A. Bosch, A. Zisserman, and X. Munoz, Representing shape with a spatial pyramid kernel,  
681 *Proceedings of the 6th ACM international conference on Image and video retrieval*, ACM,  
682 2007, pp. 401-408.
- 683 [71] K. He, X. Zhang, S. Ren, and J. Sun, Deep residual learning for image recognition, *Proceedings of  
684 the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770-778.

- 685 [72] K. Simonyan, and A. Zisserman, Very deep convolutional networks for large-scale image  
686 recognition. arXiv preprint arXiv:1409.1556 (2014).
- 687 [73] S. Li, D. Gong, and Y. Yuan, Face recognition using Weber local descriptors. *Neurocomputing*  
688 122 (2013) 272-283.
- 689 [74] S. Liu, Y. Zhang, and K. Liu, Facial expression recognition under partial occlusion based on Weber  
690 Local Descriptor histogram and decision fusion, *Control Conference (CCC), 2014 33rd*  
691 *Chinese, IEEE, 2014*, pp. 4664-4668.
- 692 [75] S. Chakraborty, S. Singh, and P. Chakraborty, Local Gradient Hexa Pattern: A Descriptor for Face  
693 Recognition and Retrieval. *IEEE Transactions on Circuits and Systems for Video Technology*  
694 (2016).
- 695 [76] B. Zhang, S. Shan, X. Chen, and W. Gao, Histogram of gabor phase patterns (hgpp): A novel object  
696 representation approach for face recognition. *IEEE Transactions on Image Processing* 16  
697 (2007) 57-68.
- 698 [77] W.-P. Choi, S.-H. Tse, K.-W. Wong, and K.-M. Lam, Simplified Gabor wavelets for human face  
699 recognition. *Pattern Recognition* 41 (2008) 1186-1199.
- 700 [78] K.-H. Pong, and K.-M. Lam, Multi-resolution feature fusion for face recognition. *Pattern*  
701 *Recognition* 47 (2014) 556-567.
- 702 [79] S. Du, Y. Yan, and Y. Ma, Local spiking pattern and its application to rotation-and illumination-  
703 invariant texture classification. *Optik-International Journal for Light and Electron Optics* 127  
704 (2016) 6583-6589.
- 705 [80] M. Verma, and B. Raman, Local tri-directional patterns: A new texture feature descriptor for image  
706 retrieval. *Digital Signal Processing* 51 (2016) 62-72.
- 707 [81] S. Fadaei, R. Amirfattahi, and M.R. Ahmadzadeh, Local derivative radial patterns: A new texture  
708 descriptor for content-based image retrieval. *Signal Processing* 137 (2017) 274-286.
- 709 [82] T. Ahonen, A. Hadid, and M. Pietikainen, Face description with local binary patterns: Application  
710 to face recognition. *IEEE transactions on pattern analysis and machine intelligence* 28 (2006)  
711 2037-2041.
- 712 [83] V. Takala, T. Ahonen, and M. Pietikäinen, Block-based methods for image retrieval using local  
713 binary patterns. *Image analysis* (2005) 13-181.
- 714 [84] L. Liu, P. Fieguth, Y. Guo, X. Wang, and M. Pietikäinen, Local binary features for texture  
715 classification: Taxonomy and experimental study. *Pattern Recognition* 62 (2017) 135-160.
- 716 [85] B. Zhang, Y. Gao, S. Zhao, and J. Liu, Local derivative pattern versus local binary pattern: face  
717 recognition with high-order local pattern descriptor. *IEEE transactions on image processing* 19  
718 (2010) 533-544.
- 719 [86] N. Dalal, and B. Triggs, Histograms of oriented gradients for human detection, *Computer Vision*  
720 *and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on, IEEE,*  
721 *2005*, pp. 886-893.
- 722 [87] P.N. Belhumeur, J.P. Hespanha, and D.J. Kriegman, Eigenfaces vs. fisherfaces: Recognition using  
723 class specific linear projection. *IEEE Transactions on pattern analysis and machine intelligence*  
724 19 (1997) 711-720.
- 725 [88] M. Felsberg, and G. Sommer, The monogenic signal. *IEEE Transactions on Signal Processing* 49  
726 (2001) 3136-3144.
- 727 [89] L. Zhang, L. Zhang, Z. Guo, and D. Zhang, Monogenic-LBP: A new approach for rotation invariant  
728 texture classification, *Image Processing (ICIP), 2010 17th IEEE International Conference on,*  
729 *IEEE, 2010*, pp. 2677-2680.
- 730 [90] M. Yang, L. Zhang, L. Zhang, and D. Zhang, Monogenic binary pattern (MBP): A novel feature  
731 extraction and representation model for face recognition, *Pattern Recognition (ICPR), 2010*  
732 *20th International Conference on, IEEE, 2010*, pp. 2680-2683.
- 733 [91] Y.-H. Oh, A.C. Le Ngo, J. See, S.-T. Liong, R.C.-W. Phan, and H.-C. Ling, Monogenic Riesz  
734 wavelet representation for micro-expression recognition, *Digital Signal Processing (DSP),*  
735 *2015 IEEE International Conference on, IEEE, 2015*, pp. 1237-1241.
- 736 [92] Z. Zeng, L. Song, Q. Zheng, and Y. Chi, A new image retrieval model based on monogenic signal  
737 representation. *Journal of Visual Communication and Image Representation* 33 (2015) 85-93.
- 738 [93] X. Huang, G. Zhao, W. Zheng, and M. Pietikainen, Spatiotemporal local monogenic binary patterns  
739 for facial expression recognition. *IEEE Signal Processing Letters* 19 (2012) 243-246.

- 740 [94] C.E. Erdem, C. Turan, and Z. Aydin, BAUM-2: a multilingual audio-visual affective face database.  
741 Multimedia Tools and Applications 74 (2015) 7429-7459.
- 742 [95] T. Kanade, J.F. Cohn, and Y. Tian, Comprehensive database for facial expression analysis,  
743 Automatic Face and Gesture Recognition, 2000. Proceedings. Fourth IEEE International  
744 Conference on, IEEE, 2000, pp. 46-53.
- 745 [96] M. Lyons, S. Akamatsu, M. Kamachi, and J. Gyoba, Coding facial expressions with gabor  
746 wavelets, Automatic Face and Gesture Recognition, 1998. Proceedings. Third IEEE  
747 International Conference on, IEEE, 1998, pp. 200-205.
- 748 [97] L.-F. Chen, and Y.-S. Yen, Taiwanese facial expression image database. Brain Mapping  
749 Laboratory, Institute of Brain Science, National Yang-Ming University, Taipei, Taiwan (2007).
- 750 [98] C. Turan, K.-M. Lam, and X. He, Soft Locality Preserving Map (SLPM) for Facial Expression  
751 Recognition. arXiv preprint arXiv:1801.03754 (2018).
- 752 [99] D. McDuff, R. El Kaliouby, and R.W. Picard, Crowdsourcing facial responses to online videos,  
753 Affective Computing and Intelligent Interaction (ACII), 2015 International Conference on,  
754 IEEE, 2015, pp. 512-518.
- 755 [100] E. Douglas-Cowie, R. Cowie, I. Sneddon, C. Cox, O. Lowry, M. Mcrorie, J.-C. Martin, L.  
756 Devillers, S. Abrilian, and A. Batliner, The HUMAINE database: addressing the collection and  
757 annotation of naturalistic and induced emotional data. Affective computing and intelligent  
758 interaction (2007) 488-500.
- 759 [101] A. Dhall, R. Goecke, S. Lucey, and T. Gedeon, Static facial expression analysis in tough  
760 conditions: Data, evaluation protocol and benchmark, Computer Vision Workshops (ICCV  
761 Workshops), 2011 IEEE International Conference on, IEEE, 2011, pp. 2106-2112.
- 762 [102] X. Peng, Z. Xia, L. Li, and X. Feng, Towards facial expression recognition in the wild: a new  
763 database and deep recognition system, Proceedings of the IEEE conference on computer vision  
764 and pattern recognition workshops, 2016, pp. 93-99.
- 765 [103] A. Mollahosseini, B. Hasani, and M.H. Mahoor, AffectNet: A Database for Facial Expression,  
766 Valence, and Arousal Computing in the Wild. arXiv preprint arXiv:1708.03985 (2017).
- 767