



# Automatic facial expression analysis: a survey

B. Fasel<sup>a,\*</sup>, Juergen Luetttin<sup>b</sup>

<sup>a</sup>IDIAP—Dalle Molle Institute for Perceptual Artificial Intelligence, Rue du Simplon 4, CH-1920 Martigny, Switzerland

<sup>b</sup>Ascom Systec AG, Applicable Research and Technology, Gewerbepark CH-5506, Maegenwil, Switzerland

Received 15 May 2001; accepted 15 February 2002

---

## Abstract

Over the last decade, automatic facial expression analysis has become an active research area that finds potential applications in areas such as more engaging human–computer interfaces, talking heads, image retrieval and human emotion analysis. Facial expressions reflect not only emotions, but other mental activities, social interaction and physiological signals. In this survey, we introduce the most prominent automatic facial expression analysis methods and systems presented in the literature. Facial motion and deformation extraction approaches as well as classification methods are discussed with respect to issues such as face normalization, facial expression dynamics and facial expression intensity, but also with regard to their robustness towards environmental changes. © 2002 Pattern Recognition Society. Published by Elsevier Science Ltd. All rights reserved.

**Keywords:** Facial expression recognition; Facial expression interpretation; Emotion recognition; Affect recognition; FACS

---

## 1. Introduction

Facial expression analysis goes well back into the nineteenth century. Darwin [1] demonstrated already in 1872 the universality of facial expressions and their continuity in man and animals and claimed among other things, that there are specific inborn emotions, which originated in serviceable associated habits. In 1971, Ekman and Friesen [2] postulated six primary emotions that possess each a distinctive content together with a unique facial expression. These prototypic emotional displays are also referred to as so called *basic emotions*. They seem to be universal across human ethnicities and cultures and comprise happiness, sadness, fear, disgust, surprise and anger. In the past, facial expression analysis was primarily a research subject for psychologists, but already in 1978, Suwa et al. [3] presented a preliminary investigation on automatic facial expression analysis from an image sequence. In the nineties, automatic facial expression analysis research gained much inertia

starting with the pioneering work of Mase and Pentland [4]. The reasons for this renewed interest in facial expressions are multiple, but mainly due to advancements accomplished in related research areas such as face detection, face tracking and face recognition as well as the recent availability of relatively cheap computational power. Various applications using automatic facial expression analysis can be envisaged in the near future, fostering further interest in doing research in different areas, including image understanding, psychological studies, facial nerve grading in medicine [5], face image compression and synthetic face animation [6], video-indexing, robotics as well as virtual reality. Facial expression recognition should not be confused with human emotion recognition as is often done in the computer vision community. While *facial expression recognition* deals with the classification of facial motion and facial feature deformation into abstract classes that are purely based on visual information, human emotions are a result of many different factors and their state might or might not be revealed through a number of channels such as emotional voice, pose, gestures, gaze direction and facial expressions. Furthermore, emotions are not the only source of facial expressions, see Fig. 1. In contrast to facial expression recognition, emotion recognition is an *interpretation attempt* and

---

\*Corresponding author.

E-mail addresses: [beat.fasel@idiap.ch](mailto:beat.fasel@idiap.ch) (B. Fasel),  
[juergen.luetttin@ascom.ch](mailto:juergen.luetttin@ascom.ch) (J. Luetttin).

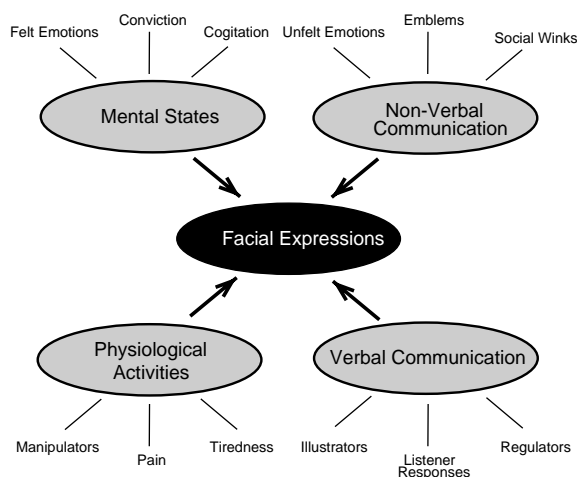


Fig. 1. Sources of facial expressions.

often demands understanding of a given situation, together with the availability of full contextual information.

## 2. Facial expression measurement

Facial expressions are generated by contractions of facial muscles, which results in temporally deformed facial features such as eye lids, eye brows, nose, lips and skin texture, often revealed by wrinkles and bulges. Typical changes of muscular activities are brief, lasting for a few seconds, but rarely more than 5 s or less than 250 ms. We would like to accurately measure facial expressions and therefore need a useful terminology for their description. Of importance is the *location* of facial actions, their *intensity* as well as their *dynamics*. Facial expression intensities may be measured by determining either the geometric deformation of facial features or the density of wrinkles appearing in certain face regions. For example the degree of a smiling is communicated by the magnitude of cheek and lip corner raising as well as wrinkle displays. Since there are inter-personal variations with regard to the amplitudes of facial actions, it is difficult to determine *absolute* facial expression intensities, without referring to the neutral face of a given subject. Note that the intensity measurement of spontaneous facial expressions is more difficult in comparison to posed facial expressions, which are usually displayed with an exaggerated intensity and can thus be identified more easily. Not only the nature of the deformation of facial features conveys meaning, but also the relative timing of facial actions as well as their temporal evolution. Static images do not clearly reveal subtle changes in faces and it is therefore essential to measure also the dynamics of facial expressions. Although the importance of correct timing is widely accepted, only a few studies have investigated this aspect systematically, mostly for smiles [7]. Facial expressions can be described

with the aid of three temporal parameters: *onset* (attack), *apex* (sustain), *offset* (relaxation). These can be obtained from human coders, but often lack precision. Few studies relate to the problem of automatically computing the onset and offset of facial expressions, especially when not relying on intruding approaches such as Facial EMG [8]. There are two main methodological approaches of how to measure the aforementioned three characteristics of facial expressions, namely *message judgment* based and *sign vehicle*-based approaches [9]. The former directly associate specific facial patterns with mental activities, while the latter represent facial actions in a coded way, prior to eventual interpretation attempts.

### 2.1. Judgment-based approaches

Judgment-based approaches are centered around the messages conveyed by facial expressions. When classifying facial expressions into a predefined number of emotion or mental activity categories, an agreement of a group of coders is taken as ground truth, usually by computing the average of the responses of either experts or non-experts. Most automatic facial expression analysis approaches found in the literature attempt to directly map facial expressions into one of the basic emotion classes introduced by Ekman and Friesen [2].

### 2.2. Sign-based approaches

With sign vehicle-based approaches, facial motion and deformation are coded into visual classes. Facial actions are hereby abstracted and described by their location and intensity. Hence, a complete description framework would ideally contain all possible perceptible changes that may occur on a face. This is the goal of facial action coding system (FACS), which was developed by Ekman and Friesen [10] and has been considered as a foundation for describing facial expressions. It is appearance-based and thus does not convey any information about e.g. mental activities associated with expressions. FACS uses 44 action units (AUs) for the description of facial actions with regard to their location as well as their intensity, the latter either with three or five levels of magnitude. Individual expressions may be modeled by single action units or action unit combinations. Similar coding schemes are EMFACS [11], MAX [12] and AFFEX [13]. However, they are only directed towards emotions. Finally, the MPEG-4-SNHG [6] is a standard that encompasses analysis, coding [14] and animation of faces (talking heads) [15]. Instead of describing facial actions only with the aid of purely descriptive AUs, scores of sign-based approaches may be interpreted by employing *facial expression dictionaries*. Friesen and Ekman introduced such a dictionary for the FACS framework [16]. Ekman et al. [17] presented also a database called facial action coding system affect interpretation database (FACS-AID), which allows to translate emotion related FACS scores into affective meanings. Emotion

interpretations were provided by several experts, but only agreed affects were included in the database.

### 2.3. Reliability of ground truth coding

The labeling of employed databases determines not only, whether a given system attempts to recognize or interpret facial expressions, but may also influence the achievable recognition accuracy, especially when it comes to facial expression timing and intensity estimations. Furthermore, chosen classification schemes affect the design of facial expression classifiers, e.g. they have an influence on the number and nature of facial action categories that have to be treated. According to Ekman [18], there are several points that need to be addressed when measuring facial expressions: (a) a separate agreement index about the scoring of specific facial actions, as typically some actions are easier to recognize than others, (b) spontaneous rather than posed facial actions, (c) various subjects including infants, children, adults and aged populations, (d) limiting the disagreement in the judgment of facial actions by providing a minimal intensity threshold of facial actions, (e) inclusion of both expert and beginners for the measurement of facial actions and (f) the reliability should be reported not only for the type, but also the intensity and dynamics of facial actions. These points can probably be easier fulfilled with sign than with judgment-based approaches as the latter can only provide a limited labeling accuracy. For example, within a single basic emotion category, there is too much room for interpretation. Cross-cultural studies have shown furthermore, that the judgment of facial expression is also culturally dependent and partially influenced by learned display rules [19]. Even though the aforementioned basic emotions are universal across cultures, the assessment is hampered, if the encoder and decoder are of different cultures [20]. Sign-based coding schemes on the other hand increase objectivity, as coders are only required to record specific concerted facial components instead of performing facial expression interpretation. An advantage of sign-based methods is also the possibility of decomposing facial expression recognition and facial expression interpretation. Hence, the performance of the employed analysis methods may be evaluated directly with regard to their visual performance.

## 3. Automatic facial expression analysis

Automatic facial expression analysis is a complex task as physiognomies of faces vary from one individual to another quite considerably due to different age, ethnicity, gender, facial hair, cosmetic products and occluding objects such as glasses and hair. Furthermore, faces appear disparate because of pose and lighting changes. Variations such as these have to be addressed at different stages of an automatic facial expression analysis system, see Fig. 2. We have a

closer look at the individual processing stages in the remainder of this chapter.

### 3.1. Face acquisition

Ideally, a face acquisition stage features an automatic face detector that allows to locate faces in complex scenes with cluttered backgrounds. Certain face analysis methods need the exact position of the face in order to extract facial features of interest while others work, if only the coarse location of the face is available. This is the case with e.g. active appearance models [21]. Hong et al. [22] used the PersonSpotter system by Steffens et al. [23] in order to perform realtime tracking of faces. The exact face dimensions were then obtained by fitting a labeled graph onto the bounding box containing the face previously detected by the PersonSpotter system. Essa and Pentland [24] located faces by using the view-based and modular eigenspace method of Pentland et al. [25]. Face analysis is complicated by face *appearance changes* caused by pose variations and illumination changes. It might therefore be a good idea to normalize acquired faces prior to their analysis:

- **Pose:** The appearance of facial expressions depends on the angle and distance at which a given face is being observed. Pose variations occur due to *scale changes* as well as *in-plane* and *out-of-plane rotations* of faces. Especially out-of-plane rotated faces are difficult to handle, as perceived facial expression are distorted in comparison to frontal face displays or may even become partly invisible. Limited out-of-plane rotations can be addressed by warping techniques, where the center positions of distinctive facial features such as the eyes, nose and mouth serve as reference points in order to normalize test faces according to some generic face models e.g. see Ref. [24]. Scale changes of faces may be tackled by scanning images at several resolutions in order to determine the size of present faces, which can then be normalized accordingly [24,26].
- **Illumination:** A common approach for reducing lighting variations is to filter the input image with Gabor wavelets, e.g. see Ref. [27]. The problem of partly lightened faces is difficult to solve. This has been addressed for the task of face recognition by Belhumeur et al. [28], but not yet sufficiently for facial expression analysis. Finally, specular reflections on eyes, teeth and wet skin may be encountered by using brightness models [29].

Note that even though *face normalization* may be a reasonable approach in conjunction with some face analysis approaches, it is *not mandatory*, as long as extracted feature parameters are normalized prior to their classification. Indeed, appearance-based model [21] and local motion model [30] approaches have dealt with significant out-of-plane rotations without relying on face normalization.

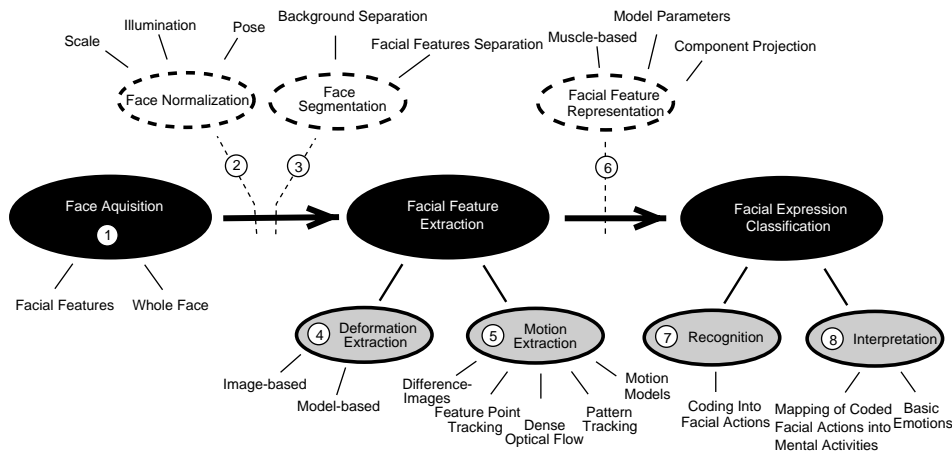


Fig. 2. Generic facial expression analysis framework. The encircled numbers are used in the system diagrams presented further below and indicate relevant processing stages. Note that the face normalization, the face segmentation as well as the facial feature representation stages are only necessary in conjunction with some specific facial feature extraction and classification methods.

Table 1

Facial feature extraction methods: Overview of prominent deformation and motion extraction methods, which were used for the task of facial expression analysis

Deformation extraction	Holistic methods	Local methods
Image-based	Neural network [27,31,32] Gabor wavelets [27,31]	Intensity profiles [33] High gradient components [34] PCA + Neural networks [35,36] Geometric face model [39] Two view point-based models [41]
Model-based	Active appearance model [21,37,38] Point distribution model [40] Labeled graphs [22,42,43]	
Motion extraction	Holistic methods	Local methods
Dense optical flow	Dense flow fields [33,34]	Region-based flow [4,44,45]
Motion models	3D motion models [24,46] 3D deformable models [48]	Parametric motion models [30,47] 3D motion models [49] Feature tracking [34,50–53]
Feature point tracking		Region-based difference-images [56]
Difference-images	Holistic diff.-imgs [27,33,54,55]	Highlighted facial features [57]
Marker-based		Dot markers [3,58]

### 3.2. Feature extraction and representation

Feature extraction methods can be categorized according to whether they focus on motion or deformation of faces and facial features, respectively, whether they act locally or holistically. Table 1 gives an overview of methods that were employed by the computer vision community for the task of facial expression analysis.

#### 3.2.1. Local versus holistic approaches

Facial feature processing may take place either holistically, where the face is processed as a whole, or locally, by focusing on facial features or areas that are prone to change

with facial expressions. We can distinguish two types of facial features:

- *Intransient facial features* are always present in the face, but may be deformed due to facial expressions. Among these, the eyes, eyebrows and the mouth are mainly involved in facial expression displays. Tissue texture, facial hair as well as permanent furrows constitute other types of intransient facial features that influence the appearance of facial expressions.
- *Transient facial features* encompass different kind of wrinkles and bulges that occur with facial expressions. Especially, the forefront and the regions surrounding the mouth and the eyes are prone to contain transient facial

features. Opening and closing of eyes and the mouth may furthermore lead to iconic changes [29], local changes of texture that cannot be predicted from antecedent frames.

*Face segmentation* allows to isolate transient and intransient features within faces or can be used to separate faces of interest from the background. Segmentation boundaries were often determined heuristically, where guidance was provided by a priori knowledge of human observers. Note that holistic feature extraction methods are good at determining prevalent facial expressions, whereas local methods are able to detect subtle changes in small areas. The latter are suitable especially with rule-based interpretation attempts, e.g. see Ref. [41].

### 3.2.2. Deformation versus motion-based approaches

Motion extraction approaches directly focus on facial changes occurring due to facial expressions, whereas deformation-based methods do have to rely on neutral face images or face models in order to extract facial features that are relevant to facial actions and not caused by e.g. intransient wrinkles due to old age. In contrast to motion-based approaches, deformation-based methods can be applied to both single images and image sequences, in the latter case by processing frames independently from each other. However, deformation-based feature extractors miss low-level directional flow information, i.e. they cannot reconstruct pixel motion. Nonetheless, high-level motion of intransient facial features may be inferred by using e.g. face and facial feature models that allow to estimate possible flow directions.

### 3.2.3. Image versus model-based approaches

Image-based methods extract features from images without relying on extensive knowledge about the object of interest. They have the advantage of being typically fast and simple. However, image-based approaches can become unreliable and unwieldy, when there are many different views of the same object that must be considered. The facial structure can also be described with the aid of 2D or 3D face models. The former allows to model facial features and faces based on their appearance, without attempting to recover the volumetric geometry of the scene, e.g. see Ref. [21]. There are two types of 3D models, namely muscle and motion models [24,59]. The latter can significantly improve the precision of motion estimations, since only physically possible motion is considered. However, 3D models often require complex mapping procedures that generate heavy computational requirements. In addition, accurate head and face models have to be constructed manually, which is a tedious undertaking.

### 3.2.4. Appearance versus muscle-based approaches

In contrast to appearance-based image and 2D model approaches, where processing focuses on the effects of facial muscle activities, muscle-based based frameworks attempt

to interfere muscle activities from visual information. This may be achieved e.g. by using 3D muscle models that allow to map extracted optical flow into muscle actions [59,60]. Modeled facial motion can hereby be restricted to muscle activations that are allowed by the muscle framework, giving control over possible muscle contractions, relaxation and orientation properties. However, the musculature of the face is complex, 3D information is not readily present and muscle motion is not directly observable. For example, there are at least 13 groups of muscles involved in the lip movements alone [61]. Mase and Pentland [4] did not use complex 3D models to determine muscle activities. Instead they translated 2D motion in predefined windows directly into a coarse estimate of muscle activity. Muscle-based approaches are not only well suited for the recognition of facial expression, but are also used to animate synthetic faces.

### 3.3. Deformation extraction

Deformation of facial features are characterized by shape and texture changes and lead to high spatial gradients that are good indicators for facial actions and may be analyzed either in the image or the spatial frequency domain. The latter can be computed by high-pass gradient or Gabor wavelet-based filters, which closely model the receptive field properties of cells in the primary visual cortex [62,63]. They allow to detect line endings and edge borders over multiple scales and with different orientations. These features reveal much about facial expressions, as both transient and intransient facial features often give raise to a contrast change with regard to the ambient facial tissue. As we have mentioned before, Gabor filters remove also most of the variability in images that occur due to lighting changes. They have shown to perform well for the task of facial expression analysis and were used in image-based approaches [27,31,33] as well as in combination with labeled graphs [22,42,43]. For an illustration of Gabor filters applied to face images see Fig. 3. We can distinguish local or holistic image-based deformation extraction approaches:

- *Holistic image-based approaches:* Several authors have taken either whole faces as features [27,31,32] or Gabor wavelet filtered whole-faces [27,31]. The main emphasis is hereby put on the classifier, which has to deal not only with face physiognomies, but in the case of image-domain-based face processing also with lighting variations. Common for most holistic face analysis approaches is the need of a thorough face-background separation in order to prevent disturbance caused by clutter.
- *Local image-based approaches:* Padgett and Cottrell [35] as well as Cottrell and Metcalfe [36] extracted facial expressions from windows placed around intransient facial feature regions (both eyes and mouth) and employed local principal component analysis (PCA) for representation purposes. Local transient facial features such as wrinkles

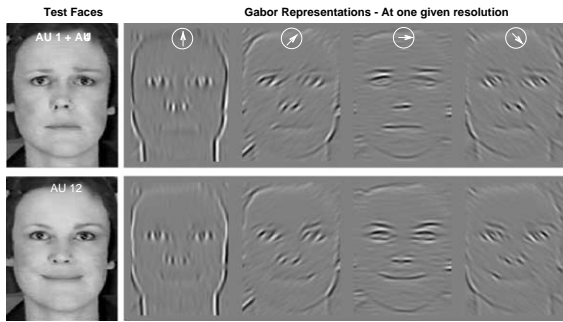


Fig. 3. Facial feature extraction using Gabor wavelets: Shown are two distinct facial expression displays on the left-hand side with the corresponding Gabor representations on the right-hand side. The latter were obtained by convoluting the face images on the left-hand side with four differently oriented wavelet kernels (at one given resolution).

can be measured by using image intensity profiles along segments [33] or by determining the density of high gradient components over windows of interest [34].

*Model-based approaches* constitute an alternative to image-based deformation extraction. *Appearance-based model* approaches allow to separate fairly well different information sources such as facial illumination and deformation changes. Lanitis et al. [21] interpreted face images by employing active appearance models (AAM) [37,38]. Faces were analyzed by a dual approach, using both shape and texture models. Active shape models (ASM) allow to simultaneously determine shape, scale and pose by fitting an appropriate point distribution model (PDM) to the object of interest, see Fig. 4. A drawback of appearance-based models is the manual labor necessary for the construction of the shape models. The latter are based on landmark points that need to be precisely placed around intransient facial features during the training of the models. Huang and Huang [40] used a point distribution model to represent the shape of a face, where shape parameters were estimated by employing a gradient-based method. Another type of holistic face models constitute the so-called *labeled graphs*, which are comprised of sparsely distributed fiducial feature points [22,42,43]. The nodes of these feature graphs consist of *Gabor jets*, where each component of a jet is a filter response of a specific Gabor wavelet extracted at a given image point. A labeled graph is matched to a test face by varying its scale and position. The obtained graph can then be compared to reference graphs in order to determine the facial expression display at hand. Kobayashi and Hara [39] used a geometric face model consisting of 30 facial characteristic points (FCP). They measured the intensity distribution along 13 vertical FCPs crossing facial lines with the aid of a neural network. Finally, Pantic and Rothkrantz [41] used a 2D point-based model composed

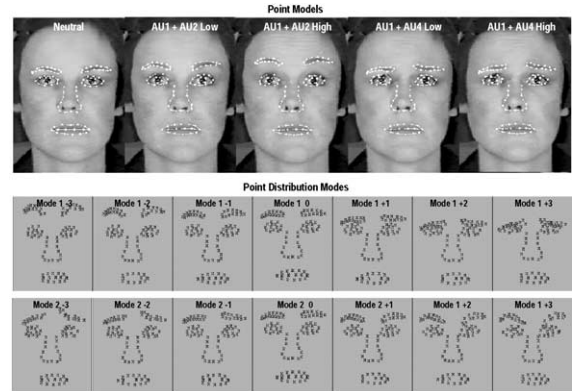


Fig. 4. Facial feature representation using active shape models (ASM): The first row shows manually placed point models (PM) that were employed to create a point distribution model (PDM), represented by a few discrete instances of two point distribution modes shown in row two (mode 1) and three (mode 2), with intensities ranging from  $-3$  to  $+3$ . The point distribution modes were computed using the active shape model toolbox implemented by Matthews [64].

of both frontal and a side views. Multiple feature detectors were applied redundantly in order to localize contours of prominent facial features prior to their modeling.

### 3.4. Motion extraction

Among the motion extraction methods that have been used for the task of facial expression analysis we find *dense optical flow*, *feature point tracking* and *difference-images*. Dense optical flow has been applied both locally and holistically:

- *Holistic dense optical flow* approaches allow for whole-faces analysis and were employed e.g. in Refs. [33,34]. Lien [34] analyzed holistic face motion with the aid of wavelet-based, multi-resolution dense optical flow. For a compacter representation of the resulting flow fields they computed PCA-based eigenflows both in horizontal and vertical directions. Fig. 5 shows sample dense optical flow fields, computed from two facial expression sequences.
- *Local dense optical flow*: Region-based dense optical flow was used by Mase and Pentland [4] in order to estimate the activity of 12 of the totally 44 facial muscles. For each muscle, a window in the face image was defined as well as an axis along which each muscle expands and contracts. Dense optical flow motion was quantified into eight directions and allowed for a coarse estimation of muscle activity. Otsuka and Ohya [44] estimated facial motion in local regions surrounding the eyes and the mouth. Feature vectors were obtained by taking 2D Fourier transforms of the vertical and horizontal optical flow fields. Yoneyama et al. [45]

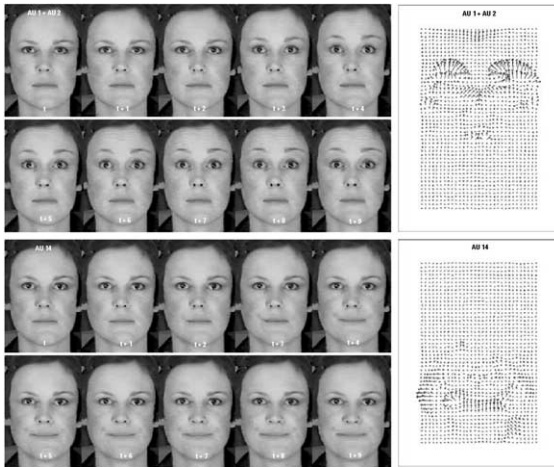


Fig. 5. Facial motion extraction using dense optical flow: Shown are two sample facial expression sequences on the left-hand side and the corresponding optical flow images on the right-hand side, which were computed with Nagel's algorithm [65]. Note the asymmetric facial action display in the lower facial expression sequence.

divided normalized test faces into  $8 \times 10$  regions, where local dense optical flow was computed and quantified region-wise into ternary feature vectors (+1/0/−1), indicating upwards, none and downwards movements, while neglecting horizontal facial movements.

Different optical flow algorithms have been applied to facial motion analysis. For instance, Lien et al. [34] employed Wu's [66] approach of optical flow to estimate facial motion by using scaling functions and wavelets from Cai and Wang [67] to capture both local and global facial characteristics. Essa and Pentland [24] used Simoncelli's [68] coarse-to-fine optical flow, while Yacoob and Davis [47] as well as Rosenblum et al. [53] employed Abdek-Mottaleb et al. [69] optical flow. Apart from a certain vulnerability to image noise and non-uniform lighting, holistic dense optical flow methods often result in prodigious computational requirements and tend to be sensitive to motion discontinuities (iconic changes) as well as non-rigid motion. Optical flow analysis can also be done in conjunction with motion models that allow for increased stability and better interpretation of extracted facial motion, e.g. muscle activations:

- *Holistic motion models:* Terzopoulos and Waters [70] have used 11 principal deformable contours (also known as "snakes") to track lip and facial features throughout image sequences with the aid of a force field, which is computed from gradients found in the images. Only frontal faces were allowed and some facial make-up was used to enhance contrast. Essa and Pentland [24] employed sophisticated 3D motion and muscle models for facial expression recognition and increased tracking stability by

Kalman filtering. DeCarlo and Metaxas [48] presented a formal methodology for the integration of optical flow and 3D deformable models and applied it to human face shapes and facial motion estimation. A relatively small number of parameters were used to describe a rich variety of face shapes and facial expressions. Eisert and Girod [46] employed 3D face models to specify shape, texture and motion. These models were also used to describe facial expressions caused by speech and were parameterized by FAPs of the MPEG-4 coding scheme.

- *Local motion models:* Black and Yacoob [30] as well as Yacoob and Davis [47] introduced *local parametric motion models* that allow, within local regions in space and time, to not only accurately model non-rigid facial motions, but to provide also a concise description of the motion associated with the edges of the mouth, nose, eyelids and eyebrows in terms of a small number of parameters. However, the employed motion models focus on the main intransient facial features involved with facial expressions (eyes, eye-brows and mouth) and the analysis of transient facial features, occurring in residual facial areas, was not considered. Last but not least, Basu et al. [49] presented a convincing approach of how to track human lip motions by using 3D models.

In contrast to low-level dense optical flow, there are also higher level variants that focus either on the movement on generic features points, patterns or markers:

- *Feature point tracking:* Here, motion estimates are obtained only for a selected set of prominent features such as intransient facial features [34,50,51]. In order to reduce the risk of tracking loss, feature points are placed into areas of high contrast, preferably around intransient facial features as is illustrated on the right-hand side of Fig. 6. Hence, the movement and deformation of the latter can be measured by tracking the displacement of the corresponding feature points. Motion analysis is directed towards objects of interest and therefore does not have to be computed for extraneous background patterns. However, as facial motion is extracted only at selected feature point locations, other facial activities are ignored altogether. The automatic initialization of feature points is difficult and was often done manually. Otsuka and Ohya [52] presented a feature point tracking approach, where feature points are not selected by human expertise, but chosen automatically in the first frame of a given facial expression sequence. This is achieved by acquiring potential facial feature points from local extrema or saddle points of luminance distributions. Tian et al. [50] used different component models for the lips, eyes, brows as well as cheeks and employed feature point tracking to adapt the contours of these models according to the deformation of the underlying facial features. Finally, Rosenblum et al. [53] tracked rectangular, facial feature enclosing regions of interest with the aid of feature points.

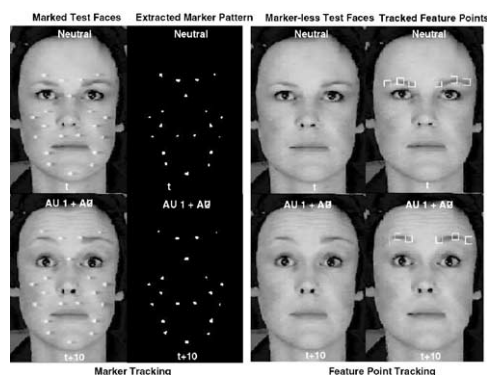


Fig. 6. Marker versus feature point tracking: On the left-hand side are shown two faces with affixed markers. The corresponding extracted marker patterns is depicted in the next column. Scale-normalized distances between marker points allow to determine underlying muscle activities [58]. A marker-less feature point tracking approach is shown on the right-hand side. Here, six feature points  $10 \times 10$  pixel windows were used to determine the displacement of the eyebrows.

- *Marker tracking:* It is possible to determine facial actions with more reliability than with previously discussed methods, namely by measuring deformation in areas, where underlying muscles interact. Unfortunately, these are mostly skin regions with relatively poor texture. Highlighting is necessary and can be done by either applying color to salient facial features and skin [57] or by affixing colored plastic dots to predefined locations on the subject's face, see the illustration on left-hand side of Fig. 6. Markers allow to render tissue motion visible and were employed in Refs. [3,58].

Note that even though the tracking of feature points or markers allows to extract motion, often only relative feature point locations, i.e. deformation information was used for the analysis of facial expressions, e.g. in Ref. [50] or [58]. Yet another way of how to extract image motion are *difference-images*: Specifically for facial expression analysis, difference-images are mostly created by subtracting a given facial image from a previously registered reference image, containing a neutral face of the same subject. However, in comparison to optical flow approaches, no flow direction can be extracted, but only differences of image intensities. In addition, accurate face normalization procedures are necessary in order to align reference faces onto the test faces. Holistic difference-image-based motion extraction was employed in Refs. [27,33,54,55]. Choudhury and Pentland [56] used motion field histograms for the modeling of eye and eye brow actions. Motion was also extracted by difference-images, but taken from consecutive image frames and further processed by using local receptive field histograms [71] in order to increase robustness with regard to rotation, translation and scale changes.

### 3.5. Classification

Feature classification is performed in the last stage of an automatic facial expression analysis system. This can be achieved by either attempting facial expression recognition using sign-based facial action coding schemes or interpretation in combination with judgment or sign/dictionary-based frameworks. We can distinguish spatial and spatio-temporal classifier approaches:

- *Spatio-temporal approaches: Hidden Markov models (HMM)* are commonly used in the field of speech recognition, but are also useful for facial expression analysis as they allow to model the dynamics of facial actions. Several HMM-based classification approaches can be found in the literature [44,72,73] and were mostly employed in conjunction with image motion extraction methods. *Recurrent neural networks* constitute an alternative to HMMs and were also used for the task of facial expression classification [53,74]. Another way of taking temporal evolution of facial expression into account are so-called spatio-temporal motion-energy templates. Here, facial motion is represented in terms of 2D motion fields. The Euclidean distance between two templates can then be used to estimate the prevalent facial expression [24].
- *Spatial approaches: Neural networks* were often used for facial expression classification [32,33,35,39,42,45,75]. They were either applied directly on face images [27,32] or combined with facial features extraction and representation methods such as PCA independent component analysis (ICA) or Gabor wavelet filters [27,31]. The former are unsupervised statistical analysis methods that allow for a considerable dimensionality reduction, which both simplifies and enhances subsequent classification. These methods have been employed both in a holistic manner [33,54,55] or locally, using mosaic-like patches extracted from small facial regions [31,33,35,54]. Dailey and Cottrell [31] applied both local PCA and Gabor jets for the task of facial expression recognition and obtained quantitatively indistinguishable results for both representations. Fig. 7 shows an illustration of PCA and ICA components obtained from facial expression images. Unfortunately, neural networks are difficult to train if used for the classification of not only basic emotions, but unconstrained facial expressions. A problem is the great number of possible facial action combinations, about 7000 AU combinations have been identified within the FACS framework [18]. An alternative to classically trained neural networks constitute compiled, rule-based neural networks that were employed e.g. in Ref. [58].

#### 3.5.1. Facial expression recognition

Traditional approaches for modeling characteristics of facial motion and deformation have relied on hand-crafted rules and symbolic mid-level representations for emotional states, which have been introduced by computer scientists



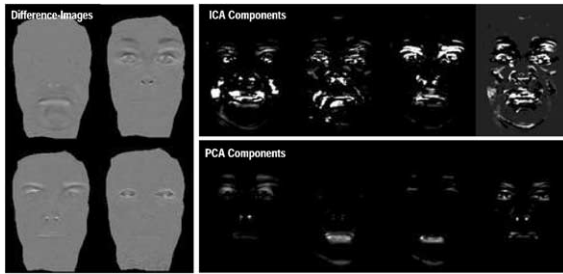


Fig. 7. Facial feature representation using data-driven methods: Sample difference-images are shown in the first row, corresponding holistic ICA components in the second and PCA components in the third row. The difference-images were computed by subtracting a neutral face reference image from face images displaying facial actions.

in the course of their investigations on facial expressions [4,30,47]. Human expertise is necessary to map these symbolic representations into e.g. emotions. However, facial signals consist of numerous distinct expressions, each with specific facial action intensity evolutions. In addition, individual realizations of facial expressions often differ only in subtle ways. This makes the task of manually creating facial expression classes rather difficult. Therefore, another group of researchers have relied on facial expression coding schemes such as MPEG-4 [46,60] or FACS [33,41,50,52,54,55,72,76]. Essa and Pentland [24] proposed an extension to FACS called FACS+, which consists of a set of control parameters using vision-based observations. In contrast to FACS, FACS+ describes also the dynamics of facial expressions.

### 3.5.2. Facial expression interpretation

Many automatic facial expression analysis systems found in the literature attempt to directly interpret observed facial expressions and mostly in terms of basic emotions [21,22,24,39,40,42,43,51,75,77,78]. Only a few systems use rules or facial expression dictionaries in order to translate coded facial actions into emotion categories [30,41]. The latter approaches have not only the advantage of accurately describing facial expressions without resorting to interpretation, but allow also to animate synthetic faces, e.g. within the FACS coding framework [15]. This is of interest, as animated synthetic faces make a direct inspection of automatically recognized facial expressions possible. See also Ref. [79] for an introduction to automatic facial expression interpretation.

## 4. Representative facial expression recognition systems

In this section, we have a closer look at a few representative facial expression analysis systems. First, we discuss *deformation* and *motion-based* feature extraction

systems. Then we introduce *hybrid* facial expression analysis systems, which employ several image analysis methods that complete each other and thus allow for a better overall performance. *Multi-modal frameworks* on the other hand integrate other non-verbal communication channels for improved facial expression interpretation results. Finally, *unified frameworks* focus on multiple facial characteristics, allowing for synergy effects between different modalities.

### 4.1. Deformation extraction-based systems

Padgett et al. [77] presented an automatic facial expression interpretation system that was capable of identifying six basic emotions. Facial data was extracted from  $32 \times 32$  pixel blocks that were placed on the eyes as well as the mouth and projected onto the top 15 PCA eigenvectors of 900 random patches, which were extracted from training images. For classification, the normalized projections were fed into an ensemble of 11 neural networks. Their output was summed and normalized again by dividing the average outputs for each possible emotion across all networks by their respective deviation over the entire training set. The largest score for a particular input was considered to be the emotion found by the ensemble of networks. Altogether 97 images of six emotions from 6 males and 6 females were analyzed and a 86% generalization performance was measured on novel face images. Lyons et al. [43] presented a Gabor wavelet-based facial expression analysis framework, featuring a node grid of Gabor jets, similar to what was used by the Von der Malsburg group for the task of face recognition [80]. Hereby, each test image was convolved with a set of Gabor filters, whose responses are highly correlated and redundant at neighboring pixels. Therefore, it was only necessary to acquire samples at specific points on a sparse grid covering the face. The projections of the filter responses along discriminant vectors, calculated from the training set, were compared at corresponding spatial frequency, orientation and locations of two face images, where the normalized dot product was used to measure the similarity of two Gabor response vectors. Lyons et al. placed graphs manually onto the faces in order to obtain a better precision for the task of facial expression recognition. Experiments were carried out on subsets of totally six different posed expressions and neutral faces of 9 Japanese female undergraduates. A generalization rate of 92% was obtained for the recognition of new expressions of known subjects and 75% for the recognition of facial expressions of novel expressers.

### 4.2. Motion extraction-based systems

Black and Yacoob [30] analyzed facial expressions with parameterized models for the mouth, the eyes and the eye brows and represented image flow with low-order polynomials [81]. A concise description of facial motion was achieved with the aid of a small number of parameters from which they derived mid- and high-level description of facial

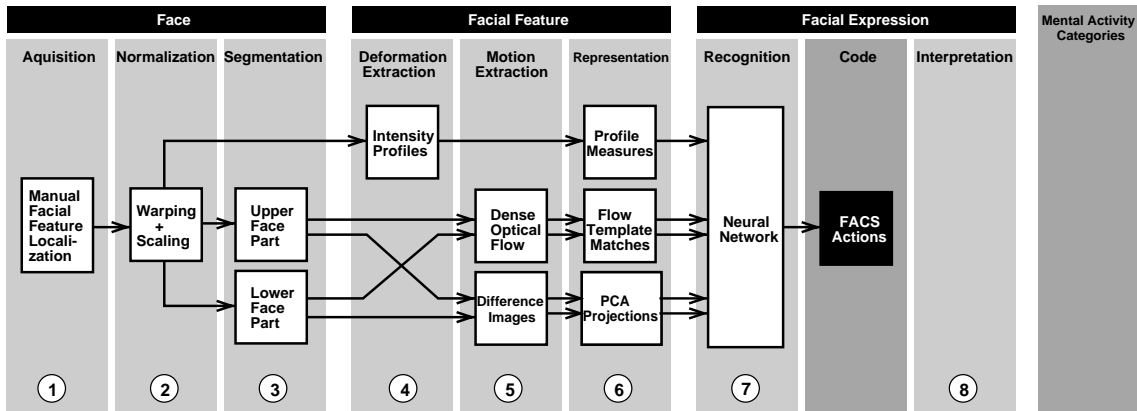


Fig. 8. Hybrid facial expression analysis system proposed by Bartlett et al. [82]. The three analysis methods employed showed to produce different error patterns and thus allow for an improved recognition performance if combined.

actions. The latter considered also temporal consistency of the mid-level predicates in order to minimize the effects of noise and inaccuracies with regard to the motion and deformation of the models. Hence, each facial expression was modeled by registering the intensities of the mid-level parameters within temporal segments (beginning, apex, ending). Extensive experiments were carried out on 40 subjects in the laboratory with a 95–100% correct recognition rate and also with television and movie sequences resulting in a 60–100% correct recognition rate. Black and Yacoob proved that a recognition of basic emotions was possible also in presence of significant pose variations and head motion. Essa and Pentland [24] presented a rather complete computer vision system featuring both automatic face detection and face analysis. Facial motion was extracted with the aid of holistic dense optical flow and coupled with 3D motion and muscle-based face models. The latter allowed to describe the facial structure including facial tissue as well as muscle actuators and their force-based deformation. Essa and Pentland located test faces automatically by using a view-based and modular eigenspace method and determined also the position of facial features. The latter were then employed in order to warp face images to match canonical face meshes, which in turn allowed to extract additional feature points corresponding to fixed nodes on the meshes. After the initial model to image registration, Simoncelli's [68] coarse-to-fine optical flow was used to compute image motion. In addition, a Kalman filter-based control framework was applied in order to prevent chaotic responses of the physical system. The employed dynamic face model allowed not only to extract muscle actuations of observed facial expressions, but it was also possible to produce noise corrected 2D motion fields via the control-theoretic approach. The latter were then classified with motion energy templates in order to extract facial actions. Experiments were carried out on 52 frontal view image sequences with a correct recognition rate of 98% for both the muscle and the 2D motion energy models.

#### 4.3. Hybrid systems

Hybrid facial expression analysis systems combine several facial expression analysis methods. This is most beneficial, if the individual estimators produce very different error patterns. Bartlett et al. [82] proposed a system that integrates holistic difference-images motion extraction coupled with PCA, feature measurements along predefined intensity profiles for the estimation of wrinkles and holistic dense optical flow for whole-face motion extraction, see Fig. 8. These three methods were compared with regard to their contribution to the facial expressions recognition task. Bartlett et al. estimated that without feature measurement, there would have been a 40% decrease of the improvement gained by all methods combined. Faces were normalized by alignment through scaling, rotation and warping of aspect ratios. However, eye and mouth centers were located manually in the neutral face frame, each test sequence had to start with. Facial expression recognition was achieved with the aid of a feed-forward neural network, made up of 10 hidden and six output units. The input of the neural network consisted of 50 PCA component projections, five feature density measurements and six optical flow-based template matches. A winner takes it all (WTA) judgment approach was chosen to select the final AU candidates. Initially, Bartlett et al.'s hybrid facial expression analysis system was able to classify six upper FACS action units on a database containing 20 subjects, correctly recognizing 92% of the AU activations, but no AU intensities. Later it was extended to allow also for the classification of lower FACS action units and achieved a 96% accuracy for 12 lower and upper face actions [33,54]. Cohn et al. [72] and Lien et al. [76] introduced systems that were based on holistic image flow analysis, feature point tracking and high-gradient component analysis methods, which were integrated into a spatio-temporal framework and combined with an HMM recognition stage, see Fig. 9. Local face motion was estimated by feature point tracking. Hereby, a

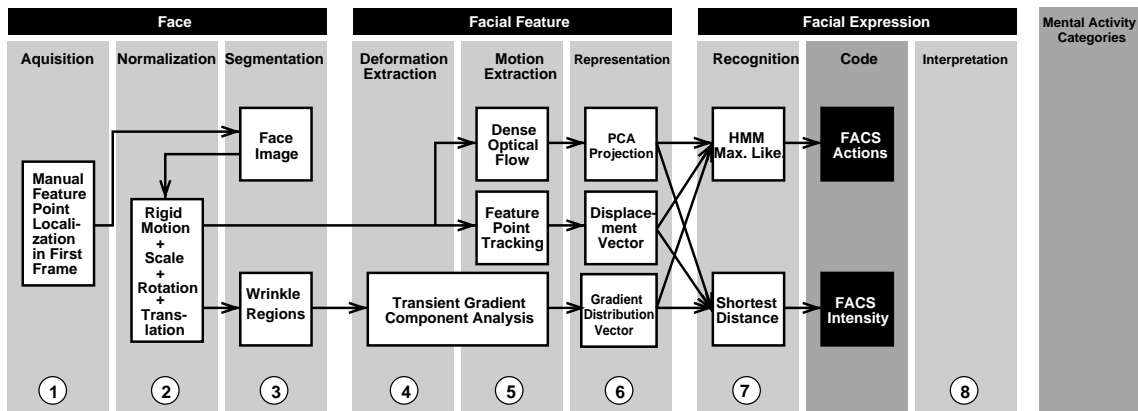


Fig. 9. Hybrid facial expression analysis system proposed by Lien et al. [76]. The whole facial expression analysis framework is situated in the spatio-temporal domain, including the classification stage that is driven by HMM.

Table 2

Selected facial expression recognition systems: Classification of facial actions using FACS

Authors	# Subjects		# Faces		# FACS AUs	Extraction methods	Classification methods	Rec. rate (%)
	Test	Train	Test	Train				
Bartlett [33]	4–20	4–20	111*5	111*5	6 <sup>a</sup> +6 <sup>b</sup>	Diff – Img. + G. Jets	Near. neigh.	96
	4–20	4–20	111*5	111*5	6 <sup>a</sup> +6 <sup>b</sup>	Diff. Img.	ICA + N. neigh.	96
	4–20	4–20	111*5	111*5	6 <sup>a</sup> +6 <sup>b</sup>	Optical flow	Motion Templ.	86
Fasel [55]	1	1	45	182	9 <sup>c</sup>	Diff. – Img.	ICA + Eucl. Dist.	83/41
Cohn [72]	30	N/A	N/A	N/A	15	Feat. point track.	HMM	86
Lien [76]	N/A	N/A	75*20	60*20	3 <sup>a</sup>	Feat. point track.	HMM	85
	N/A	N/A	150*20	120*20	6 <sup>b</sup>	Feat. point track.	HMM	88
	N/A	N/A	75*20	44*20	3 <sup>a</sup>	Optical flow	PCA + HMM	93
	N/A	N/A	160*20	100*20	4 <sup>a</sup>	High grad. comp.	HMM	85
	N/A	N/A	80*20	50*20	2 <sup>b</sup>	High grad. comp.	HMM	81
Pantic [41]	8	N/A	496	N/A	31	Multi feat. detect.	Expert rules	89

<sup>a</sup>Upper-face FACS coding.

<sup>b</sup>Lower-face FACS coding.

<sup>c</sup>Nine asymmetric FACS classes, each with five intensity levels on two face sides.

pyramidal multi-scale search approach was employed that is sensitive to subtle feature motion and also allowed to track large displacements of feature motion with sub-pixel accuracy. Holistic facial motion on the other hand was estimated by employing Wu's multi-resolution wavelet-based optical flow [66]. Forehead, cheek and chin regions were analyzed for transient facial features by using high-gradient component analysis based on horizontal, vertical and diagonal line and edge detectors in the spatial and frame comparisons in the temporal domain. The latter allowed to separate transient from intransient facial features and hair occlusion. Face tracking and face alignment were manually initialized by se-

lecting three facial feature points in the first frame of each test image sequence. Lien et al.'s system was trained to analyze both the activity and intensity of 15 FACS AUs situated in the brow, eye and mouth regions. The holistic dense optical flow approach gave the best average AU recognition rates, followed by feature point tracking and the high gradient component analysis approach, see also Table 2. Lien et al.'s facial expression analysis system performed well even with difficult sequences such as those containing baby faces. The latter differ from adult faces both in morphology and tissue texture. Unfortunately, heavy computational requirements arised with the use of optical flow.

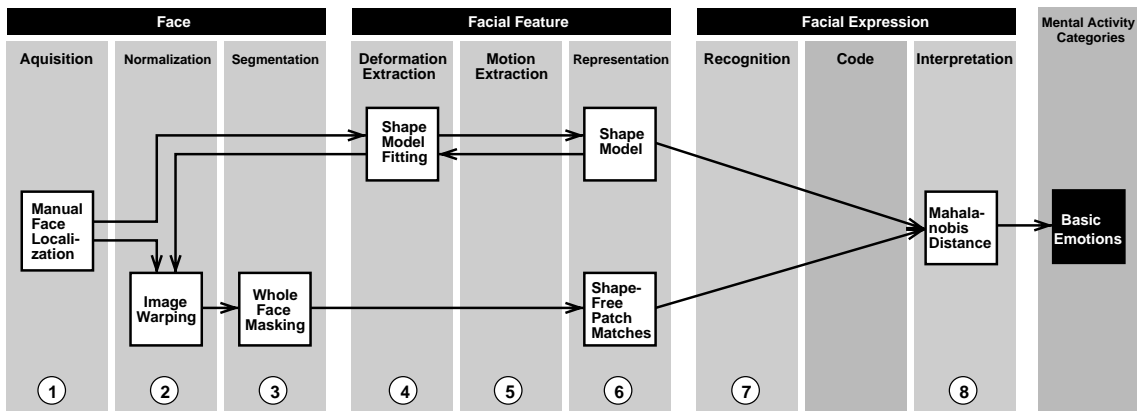


Fig. 10. Unified facial expression analysis framework proposed by Lanitis et al. [21]. It is based on active shape models (ASM). Shape information is also used to extract shape-free whole-face patches that represent texture information.

#### 4.4. Multimodal frameworks

Today, most facial expression analysis systems are of the *unimodal* type, as they focus only on facial expressions when determining mental activities. However, the evaluation of multiple communication channels may foster robustness as well as improve correct interpretation of facial expressions in ambiguous situations. At present, most attempts of channel fusion are of the *bimodal* type and integrate voice in addition to facial expressions. Vocal expressions are conveyed by prosodic features, which include the fundamental frequency, intensity and rhythm of the voice. Cohn and Katz [83] as well as Chen et al. [84] focused on the fundamental frequency, as it is an important voice feature for emotion recognition and can be easily extracted.

#### 4.5. Unified frameworks

Facial expression recognition may be improved by considering not only facial actions but also face characteristics such as identity, gender, age and ethnicity. Lanitis et al. [21] proposed a unified framework that performs multiple face processing tasks in parallel, albeit without influencing each other. They employed 2D active shape models that yield, once aligned to a test face, face appearance parameters from which it is possible to estimate 3D pose, identity, gender and facial expressions. During training, shape models are derived from a set of images by statistical analysis of the landmark point positions. These represent main facial features and are placed manually on training images prior to the shape model creation process, see Fig. 10. During testing, facial features are located in test image using ASM search [85], guided by the flexible shape models obtained during training. Gray-level profile information is collected at each model point and used to determine the best fitting shape model. Lanitis et al. then deformed test faces to the

mean face shape by using the previously registered shape information in order to extract holistic shape-less patches, which account for facial texture. Finally, Hong et al. [22] proposed an online facial expression recognition system, which is based on personalized galleries and uses identity information in conjunction with facial expression analysis. Faces of interest are detected and tracked in live video sequences with the aid of the PersonSpotter system [23] and recognition of facial expressions is achieved by performing elastic graph matching. The nodes of the employed graphs are comprised of Gabor jets and are fitted to given test faces. The obtained graph is first used to determine the identity of a given subject by choosing the closest match. In a second stage, the closest matching graph found in the personalized gallery of the identified person is used to determine the displayed facial expression, thus allowing for a better focus on intra-personal variations.

## 5. Discussion

In this survey on automatic facial expression analysis, we have discussed automatic face analysis with regard to different motion and deformation-based extraction methods, model and image-based representation techniques as well as recognition and interpretation-based classification approaches. It is not possible to directly compare facial expression recognition results of face analysis systems found in the literature due to varying facial action labeling and different test beds that were used for the assessment. Recently, a rather complete facial expression database has been presented by Kanade et al. [86]. However, there is a lack of publicly available, all encompassing facial expression databases that would allow for testing facial expression analysis methods in a more transparent way. Nonetheless, we tried to characterize a few selected systems with regard to the employed

Table 3

Selected facial expression interpretation systems: Classification of emotional displays. Note that the systems presented in the last two rows perform facial action recognition prior to using dictionaries in order to interpret facial actions

Authors	# Subjects		# Faces		# Em. class.	Extraction methods	Classification methods	Rec. rate (%)
	Test	Train	Test	Train				
Lyons [43]	9	N/A	193	N/A	7	G.Wav. + El.Gr.	LDA + PCA + Cl.	75–92
Kobayashi [39]	15	N/A	90	N/A	6	—	Feed Forw. NN	85
Rosenblum [53]	32	20/14	34 sq.	20/14 sq.	2	Optical flow	RBF NN	88
Padgett [35]	12	N/A	N/A	N/A	6	—	PCA + NN	86
Essa [24]	8	8	8 sq.	8 sq.	6	Opt.F. + 3D M.	Motion templ.	98
Lanitis [21]	30	30	300	390	7	Appear. mod.	Mahal. distance	74
Black [30]	40	N/A	70 sq.	N/A	6	Motion mod.	Expert rules	83–100
Pantic [41]	8	N/A	496	N/A	6	Mul. feat. det.	Expert rules	91

feature extraction and classification methods. Table 2 lists systems that perform facial expression recognition by classifying facial actions, while Table 3 presents systems that attempt both direct interpretation of emotional facial displays or indirect interpretation via facial expression dictionaries. The application of currently available automatic facial expression recognition systems is often very restricted due to the limited robustness and hard constraints imposed on the recording conditions. Many systems assume faces to be centered in the input image and seen from a near frontal view throughout the whole test sequence. Also, it is often taken for granted that there are only small rigid head motions between any two consecutive frames. In addition, most facial expression analysis systems require important *manual intervention* for the detection and accurate normalization of test faces, during the initialization of facial feature tracking approaches or for warping video sequences. Most facial expression analysis systems are limited to the analysis of either static images or image sequences. However, an ideal system should be capable of analyzing both static images as well as image sequences, as there are sometimes no image sequence available, respectively if there are, motion should be extracted in order to obtain directional information of skin and facial feature deformation. Furthermore, the measurement of *facial expression intensities* has only been addressed by a few systems [32,34,55]. It is of importance for the interpretation of facial expressions, especially when attempting to analyze the temporal evolution and timing of facial actions. Out-of-plane rotated faces are difficult to tackle and only a few approaches found in the literature were able to deal with this problem: Active appearance models [21,38], local parametric models [30,47], 3D motion models [48,49] and to some degree also feature point tracking approaches [34,50]. Hybrid facial expression analysis systems are also of interest, as they combine different face analysis methods and may thus give better recognition results than the individual methods applied on their own [33]. This is true, if the employed extraction algorithms focus on different facial features or the combined extraction, representation and recognition stages produce different error patterns.

## 6. Conclusion

Today, most facial expression analysis systems attempt to map facial expressions directly into basic emotional categories and are thus unable to handle facial actions caused by non-emotional mental and physiological activities. FACS may provide a solution to this dilemma, as it allows to classify facial actions prior to any interpretation attempts. So far, only marker-based systems are able to reliably code all FACS action unit activities and intensities [58]. More work has to be done in the field of automatic facial expression interpretation with regard to the integration of other communication channels such as voice and gestures. Although facial expressions often occur during conversations [87], none of the cited approaches did consider this possibility. If automatic facial expression analysis systems are to be operated autonomously, current feature extraction methods have to be improved and extended with regard to robustness in natural environments as well as independence of manual intervention during initialization and deployment.

## 7. Summary

In recent years, facial expression analysis has become an active research area. Various approaches have been made towards robust facial expression recognition, applying different image acquisition, analysis and classification methods. Facial expression analysis is an inherently multi-disciplinary field and it is important to look at it from all domains involved in order to gain insight on how to build reliable automated facial expression analysis systems. This fact has often been neglected in various implementations presented in the literature. Facial expressions reflect not only emotions, but also cognitive processes, social interaction and physiological signals. However, most facial expression analysis systems have attempted to map facial expressions directly towards basic emotions, which represents an ill-posed problem. Decoupling facial expression recognition and facial expression interpretation may provide a solution to this dilemma. This can be achieved by first coding facial expressions with an

appearance-based representation scheme such as facial action coding system (FACS) and then using facial expression dictionaries in order to translate recognized facial actions into mental activity categories.

In this survey, we have reviewed the most prominent automatic facial expression analysis methods and systems presented in the literature. Facial motion and deformation extraction approaches as well as facial feature representation and classification methods were discussed with respect to issues such as face normalization, facial expression dynamics and intensity, as well as robustness towards environmental changes.

The application of currently available automatic facial expression recognition systems to the analysis of natural scenes is often very restricted due to the limited robustness of these systems and the hard constraints posed on the subject and on the recording conditions. Especially out-of-plane rotated and partly occluded faces due to facial hair or sun glasses are difficult to handle. Furthermore, many analysis methods make the hypothesis that faces are centered in the image and seen from a near frontal view throughout test sequences. Often, small rigid head motion between any two consecutive frames is assumed. Most facial expression analysis systems need important manual intervention for the accurate normalization of test faces and initialization of extraction methods such as localization of facial feature points, facial feature template selection and manual warping of video sequences.

In this article, we also had a closer look at a few representative facial expression analysis systems and discussed deformation and motion-based feature extraction systems, hybrid systems based on multiple complementary face processing tasks and multimodal systems that integrate e.g. visual and acoustic signals. As we have seen, unified frameworks are of interest as well, as they allow to focus on multiple facial characteristics such as face identity and facial expression displays and thus allow for synergy effects between different modalities. We concluded this survey by summarizing recognition results and shortcomings of currently employed analysis methods and proposed possible future research directions.

Various applications using automatic facial expression analysis can be envisaged in the near future, fostering further interest in doing research in the fields of facial expression recognition, facial expression interpretation and the facial expression animation. Non-verbal information transmitted by facial expressions is of great importance in different areas, including image understanding, psychological studies, facial nerve grading in medicine, face image compression and synthetic face animation, more engaging human-machine interfaces, view-indexing, robotics as well as virtual reality.

## References

- [1] C. Darwin, *The Expression of the Emotions in Man and Animals*, J. Murray, London, 1872.

- [2] P. Ekman, W.V. Friesen, Constants across cultures in the face and emotion, *J. Personality Social Psychol.* 17 (2) (1971) 124–129.
- [3] M. Suwa, N. Sugie, K. Fujimora, A preliminary note on pattern recognition of human emotional expression, *Proceedings of the Fourth International Joint Conference on Pattern Recognition*, Kyoto, Japan, 1978, pp. 408–410.
- [4] K. Mase, A. Pentland, Recognition of facial expression from optical flow, *IEICE Trans. E* 74 (10) (1991) 3474–3483.
- [5] P. Dulguerov, F. Marchal, D. Wang, C. Gysin, P. Gidley, B. Gantz, J. Rubinstein, S. Seiff, L. Poon, K. Lun, Y. Ng, Review of objective topographic facial nerve evaluation methods, *Am. J. Otol.* 20 (5) (1999) 672–678.
- [6] R. Koenen, Mpeg-4 Project Overview, International Organisation for Standardisation, ISO/IEC JTC1/SC29/WG11, La Baule, 2000.
- [7] D. Messinger, A. Fogel, K.L. Dickson, What's in a smile? *Develop. Psychol.* 35 (3) (1999) 701–708.
- [8] G. Schwartz, P. Fair, P. Salt, M. Mandel, G. Klerman, Facial expression and imagery in depression: an electromyographic study, *Psychosomatic Med.* 38 (337–347) (1976).
- [9] P. Ekman, *Emotions in the Human Face*, Cambridge University Press, Cambridge, 1982.
- [10] P. Ekman, W.V. Friesen, *Facial Action Coding System: A Technique for the Measurement of Facial Movement*, Consulting Psychologists Press, Palo Alto, 1978.
- [11] W. Friesen, P. Ekman, *Emotional facial action coding system*, unpublished manual, 1984.
- [12] C. Izard, The maximally discriminative facial movement coding system (MAX), Available from Instructional Resource Center, University of Delaware, Newark, Delaware, 1979.
- [13] C. Izard, L. Dougherty, E. Hembree, A system for indentifying affect expressions by holistic judgments, unpublished manuscript, 1983.
- [14] N. Tsapatsoulis, K. Karpouzis, G. Stamou, A fuzzy system for emotion classification based on the MPEG-4 facial definition parameter, *European Association on Signal Processing EUSIPCO*, 2000.
- [15] M. Hoch, G. Fleischmann, B. Girod, Modeling and animation of facial expressions based on B-splines, *Visual Comput.* (1994) 87–95.
- [16] W. Friesen, P. Ekman, *Dictionary—interpretation of FACS scoring*, unpublished manuscript, 1987.
- [17] P. Ekman, E. Rosenberg, J. Hager, Facial action coding system affect interpretation database (FACSAID), <http://nirc.com/Expression/FACSAID/facsaid.html>, July 1998.
- [18] P. Ekman, Methods for measuring facial actions, in: K. Scherer, P. Ekman (Eds.), *Handbook of Methods in Nonverbal Behaviour Research*, Cambridge University Press, Cambridge, 1982, pp. 45–90.
- [19] D. Matsumoto, Cultural similarities and differences in display rules, *Motivation Emotion* 14 (3) (1990) 195–214.
- [20] D. Matsumoto, Ethnic differences in affect intensity, emotion judgments, display rules, and self-reported emotional expression, *Motivation Emotion* 17 (1993) 107–123.
- [21] A. Lanitis, C. Taylor, T. Cootes, Automatic interpretation and coding of face images using flexible models, *IEEE Trans. Pattern Anal. Mach. Intell.* 19 (7) (1997) 743–756.
- [22] H. Hong, H. Neven, C. Von der Malsburg, Online facial expression recognition based on personalized galleries, *Proceedings of the Second International Conference on*

- Automatic Face and Gesture Recognition, (FG'98), IEEE, Nara, Japan, 1998, pp. 354–359.
- [23] J. Steffens, E. Elagin, H. Neven, PersonSpotter—fast and robust system for human detection, tracking and recognition, Proceedings of the Second International Conference on Face and Gesture Recognition, (FG'98), Nara, Japan, 1998, pp. 516–521.
- [24] I. Essa, A. Pentland, Coding, analysis, interpretation and recognition of facial expressions, *IEEE Trans. Pattern Anal. Mach. Intell.* 19 (7) (1997) 757–763.
- [25] A. Pentland, B. Moghaddam, T. Starner, View-based and modular eigenspaces for face recognition, *IEEE Conference on Computer Vision and Pattern Recognition*, Seattle, NA, USA, 1994, pp. 84–91.
- [26] H. Rowley, S. Baluja, T. Kanade, Neural network-based face detection, *IEEE Trans. Pattern Anal. Mach. Intell.* 20 (1) (1998) 23–38.
- [27] W. Fellenz, J. Taylor, N. Tsapatsoulis, S. Kollias, Comparing template-based, feature-based and supervised classification of facial expressions from static images, Proceedings of Circuits, Systems, Communications and Computers (CSCC'99), Nugata, Japan, 1999, pp. 5331–5336.
- [28] P. Belhumeur, J. Hespanha, D. Kriegman, Eigenfaces vs. fisherfaces: recognition using class specific linear projection, *IEEE Trans. Pattern Anal. Mach. Intell.* 19 (7) (1997) 711–720.
- [29] M. Black, D. Fleet, Y. Yacoob, A framework for modeling appearance change in image sequences, Sixth International Conference on Computer Vision (ICCV'98), IEEE Computer Society Press, Silverspring, MD, 1998.
- [30] M. Black, Y. Yacoob, Recognizing facial expressions in image sequences using local parameterized models of image motion, *Internat. J. Comput. Vision* 25 (1) (1997) 23–48.
- [31] M. Dailey, G. Cottrell, PCA Gabor for expression recognition, Institution UCSD, Number CS-629, 1999.
- [32] C. Lisetti, D. Rumelhart, Facial expression recognition using a neural network, Proceedings of the 11th International Flairs Conference, AAAI Press, New York, 1998.
- [33] M. Bartlett, Face image analysis by unsupervised learning and redundancy reduction, Ph.D. Thesis, University of California, San Diego, 1998.
- [34] J. Lien, Automatic recognition of facial expression using hidden Markov models and estimation of expression intensity, Ph.D. Thesis, The Robotics Institute, CMU, April 1998.
- [35] C. Padgett, G. Cottrell, Representing face image for emotion classification, in: M. Mozer, M. Jordan, T. Petsche (Eds.), *Advances in Neural Information Processing Systems*, Vol. 9, MIT Press, Cambridge, MA, pp. 894–900.
- [36] G.W. Cottrell, J. Metcalfe, EMPATH: face, gender and emotion recognition using holons, in: R. Lippman, J. Moody, D. Touretzky (Eds.), *Advances in Neural Information Processing Systems*, Morgan Kaufman, San Mateo, CA, Vol. 3, 1991, pp. 564–571.
- [37] T. Cootes, G. Edwards, C. Taylor, Active appearance models, *IEEE PAMI* 23 (6) (2001) 681–685.
- [38] G. Edwards, T. Cootes, C. Taylor, Face recognition using active appearance models, Proceedings of the Fifth European Conference on Computer Vision (ECCV), Vol. 2, University of Freiburg, Germany, 1998, pp. 581–695.
- [39] H. Kobayashi, F. Hara, Facial interaction between animated 3D face robot and human beings, Proceedings of the International Conference on Systems, Man and Cybernetics, Orlando, FL, USA, 1997, pp. 3732–3737.
- [40] C. Huang, Y. Huang, Facial expression recognition using model-based feature extraction and action parameters classification, *J. Visual Commun. Image Representation* 8 (3) (1997) 278–290.
- [41] M. Pantic, L. Rothkrantz, Expert system for automatic analysis of facial expression, *Image Vision Comput. J.* 18 (11) (2000) 881–905.
- [42] Z. Zhang, M. Lyons, M. Schuster, S. Akamatsu, Comparison between geometry-based and Gabor-wavelets-based facial expression recognition using multi-layer perceptron, *IEEE Proceedings of the Second International Conference on Automatic Face and Gesture Recognition (FG'98)*, Nara, Japan, 1998, pp. 454–459.
- [43] M. Lyons, J. Budynek, S. Akamatsu, Automatic classification of single facial images, *IEEE Trans. Pattern Anal. Mach. Intell.* 21 (12) (1999).
- [44] T. Otsuka, J. Ohya, Spotting segments displaying facial expression from image sequences using HMM, *IEEE Proceedings of the Second International Conference on Automatic Face and Gesture Recognition (FG'98)*, Nara, Japan, 1998, pp. 442–447.
- [45] M. Yoneyama, Y. Iwano, A. Ohtake, K. Shirai, Facial expression recognition using discrete Hopfield neural networks, Proceedings of the International Conference on Image Processing (ICIP), Santa Barbara, CA, USA, Vol. 3, 1997, pp. 117–120.
- [46] P. Eisert, B. Girod, Facial expression analysis for model-based coding of video sequences, *Picture Coding Symposium*, Berlin, Germany, 1997, pp. 33–38.
- [47] Y. Yacoob, L.S. Davis, Recognizing human facial expression from long image sequences using optical flow, *IEEE Trans. Pattern Anal. Mach. Intell.* 18 (6) (1996) 636–642.
- [48] D. DeCarlo, D. Metaxas, The integration of optical flow and deformable models with applications to human face shape and motion estimation, Proceedings of the International Conference on Computer Vision and Pattern Recognition (CVPR'96), 1996, pp. 231–238.
- [49] S. Basu, N. Oliver, A. Pentland, 3D modeling and tracking of human lip motions, Proceedings of ICCV 98, Bombay, India, 1998.
- [50] Y. Tian, T. Kanade, J. Cohn, Recognizing action units for facial expression analysis, *IEEE Trans. Pattern Anal. Mach. Intell.* 23 (2) (2001) 97–115.
- [51] M. Wang, Y. Iwai, M. Yachida, Expression recognition from time-sequential facial images by use of expression change model, *IEEE Proceedings of the Second International Conference on Automatic Face and Gesture Recognition (FG'98)*, Nara, Japan, 1998, pp. 324–329.
- [52] T. Otsuka, J. Ohya, Extracting facial motion parameters by tracking feature points, Proceedings of First International Conference on Advanced Multimedia Content Processing, Osaka, Japan, 1998, pp. 442–453.
- [53] M. Rosenblum, Y. Yacoob, L. Davis, Human expression recognition from motion using a radial basis function network architecture, *IEEE Trans. Neural Networks* 7 (5) (1996) 1121–1138.
- [54] G. Donato, S. Bartlett C. Hager, P. Ekman, J. Sejnowski, Classifying facial actions, *IEEE Trans. Pattern Anal. Mach. Intell.* 21 (10) (1999) 974–989.
- [55] B. Fasel, J. Luetttin, Recognition of asymmetric facial action unit activities and intensities, Proceedings of the International

- Conference on Pattern Recognition (ICPR 2000), Barcelona, Spain, 2000.
- [56] T. Choudhury, A. Pentland, Motion field histograms for robust modeling of facial expressions, Proceedings of the International Conference on Pattern Recognition (ICPR 2000), Barcelona, Spain, 2000.
  - [57] B. Bascle, A. Blake, Separability of pose and expression in facial tracking and animation, Proceedings of the International Conference on Computer Vision, Bombay, India, 1998.
  - [58] S. Kaiser, T. Wehrle, Automated coding of facial behavior in human–computer interactions with FACS, *J. Nonverbal Behavior* 16 (2) (1992) 67–83.
  - [59] I. Essa, A. Pentland, Facial expression recognition using a dynamic model and motion energy, *IEEE Proceedings of the Fifth International Conference on Computer Vision (ICCV 1995)*, Cambridge, MA, 1995, pp. 360–367.
  - [60] K. Karpouzis, G. Votsis, G. Moschovitis, S. Kollias, Emotion recognition using feature extraction and 3-D models, Proceedings of IMACS International Multiconference on Circuits and Systems Communications and Computers (CSCC'99), Athens, Greece, 1999, pp. 5371–5376.
  - [61] W. Hardcastle, *Physiology of Speech Production*, Academic Press, New York, 1976.
  - [62] D. Pollen, S. Ronner, Phase relationship between adjacent simple cells in the visual cortex, *Science* 212 (1981) 1409–1411.
  - [63] J. Daugman, Complete discrete 2D Gabor transform by neural networks for image analysis and compression, *IEEE Trans. Acoustics, Speech Signal Process.* 36 (1988) 1169–1179.
  - [64] I. Matthews, Active shape model toolbox, University of East Anglia, Norwich, UK, Matlab Toolbox version 2.0.0, July 1997.
  - [65] H. Nagel, On the estimation of optical flow: relations between different approaches and some new results *Artif. Intell.* 33 (1987) 299–324.
  - [66] Y. Wu, T. Kanade, J. Cohn, C. Li, Optical flow estimation using wavelet motion model, *IEEE International Conference on Computer Vision*, Bombay, India, 1998, pp. 992–998.
  - [67] W. Cai, J. Wang, Adaptive multiresolution collocation methods for initial boundary value problems of nonlinear PDEs, *Soc. Ind. Appl. Math.* 33 (3) (1996) 937–970.
  - [68] E. Simoncelli, Distributed representation and analysis of visual motion, Ph.D Thesis, Massachusetts Institute of Technology, 1993.
  - [69] M. Abdel-Mottaleb, R. Chellappa, A. Rosenfeld, Binocular motion stereo using MAP estimation, *IEEE CVPR* (1993) 321–327.
  - [70] D. Terzopoulos, K. Waters, Analysis of facial images using physical and anatomical models, *Proceeding of the Third International Conference on Computer Vision*, Osaka, Japan, 1990, pp. 727–732.
  - [71] B. Schiele, J. Crowley, Probabilistic object recognition using multidimensional receptive field histograms, Proceedings of the International Conference on Pattern Recognition (ICPR 1996), Vienna, Austria, 1996.
  - [72] J. Cohn, A. Zlochow, J. Lien, Y. Wu, T. Kanade, Automated face coding: a computer-vision based method of facial expression analysis, *Seventh European Conference on Facial Expression Measurement and Meaning*, Salzburg, Austria, 1997, pp. 329–333.
  - [73] N. Oliver, A. Pentland, F. Berard, LAFTER: a real-time lips and face tracker with facial expression recognition, Proceedings of the IEEE Conference on Computer Vision (CVPR97), S. Juan, Puerto Rico, 1997.
  - [74] H. Kobayashi, F. Hara, Dynamic recognition of basic facial expressions by discrete-time recurrent neural network, Proceedings of the International Joint Conference on Neural Networks, 1993, pp. 155–158.
  - [75] J. Zhao, G. Kearney, Classifying facial emotions by backpropagation neural networks with fuzzy inputs, Proceedings of the International Conference on Neural Information Processing, Vol. 1, 1996, pp. 454–457.
  - [76] J. Lien, T. Kanade, J. Cohn, C. Li, Automated facial expression recognition based on FACS action units, *IEEE Proceedings of the Second International Conference on Automatic Face and Gesture Recognition (FG'98)*, Nara, Japan, 1998.
  - [77] C. Padgett, G. Cottrell, R. Adolphs, Categorical perception in facial emotion classification, Proceedings of the 18th Annual Conference of the Cognitive Science Society, San Diego, CA, USA, 1996.
  - [78] S. Kimura, M. Yachida, Facial expression recognition and its degree estimation, *IEEE Conference on Computer Vision and Pattern Recognition*, San Juan, Puerto Rico, 1997, pp. 295–300.
  - [79] C. Lisetti, D. Schiano, Automatic facial expression interpretation: where human–computer interaction, artificial intelligence and cognitive science intersect pragmatics and cognition (Special issue on facial information processing: a multidisciplinary perspective) *Pragmat. Cognition* 8 (1) (2000) 185–235.
  - [80] M. Lades, J. Vorbruggen, J. Buhmann, J. Lange, C. Von der Malsburg, Distortion invariant object recognition in the dynamic link architecture, *IEEE Trans. Comput.* 42 (1993) 300–311.
  - [81] J. Bergen, P. Anandan, K. Hanna, R. Hingorani, Hierarchical model-based motion estimation, in: G. Sandini (Ed.), Proceedings of the Second European Conference on Computer Vision, ECCV-92, Lecture Notes in Computer Science, Vol. 588, Springer, Berlin, 1992, pp. 237–252.
  - [82] M. Bartlett, P. Viola, T. Sejnowski, B. Golomb, J. Larsen, J. Hager, P. Ekman, Classifying facial action, *Advances in Neural Information Processing Systems*, Vol. 8, MIT Press, Cambridge, 1996.
  - [83] J. Cohn, G. Katz, Bimodal expression of emotion by face and voice, *Workshop on Face/Gesture Recognition and Their Applications*, Sixth ACM International Multimedia Conference, Bristol, UK, 1998.
  - [84] L. Chen, T. Huang, T. Miyasato, R. Nakatsu, Multimodal human emotion/expression recognition, Proceedings of the Second International Conference on Automatic Face and Gesture Recognition (FG'98), IEEE, Nara, Japan, 1998.
  - [85] T. Cootes, C. Taylor, A. Lanitis, Multi-resolution search using active shape models, *12th International Conference on Pattern Recognition*, Vol. 1, IEEE CS Press, Los Alamitos, CA, 1994, pp. 610–612.
  - [86] T. Kanade, J. Cohn, Y. Tian, Comprehensive database for facial expression analysis, *IEEE Proceedings of the Fourth International Conference on Automatic Face and Gesture Recognition (FG'00)*, Grenoble, France, 2000.
  - [87] P. Ekman, About brows: emotional and conversational signals in: J. Aschoff, M. Con Carnach, K. Foppa, W. Lepenies, D. Plog (Eds.), *Human Ethology*, Cambridge University Press, Cambridge, 1979, pp. 169–202.



**About the Author**—BEAT FASEL graduated from the Swiss Federal Institute of Technology Lausanne (EPFL) with a diploma in Communication Systems. He currently works towards a Ph.D. degree at IDIAP in Martigny, Switzerland. His research interests include computer vision, pattern recognition and artificial intelligence.

**About the Author**—JUERGEN LUETTIN received a Ph.D. degree in Electronic and Electrical Engineering from the University of Sheffield, UK, in the area of visual speech and speaker recognition. He joined IDIAP in Martigny, Switzerland, in 1996 as a research assistant where he worked on multimodal biometrics. From 1997 to 2000, he was head of the computer vision group at IDIAP, where he initiated and lead several European Community and Swiss SNF projects in the area of biometrics, speech recognition, face analysis and document recognition. In 2000, he joined Ascom AG in Maegenwil, Switzerland as head of the technology area Pattern Recognition. Dr. Luettin has been a visiting researcher at the Center for Language and Speech Processing at the Johns Hopkins University, Baltimore, in 1997 (large vocabulary conversational speech recognition) and 2000 (audio–visual speech recognition). His research interests include speech recognition, computer vision, biometrics, and multimodal recognition.