

# Chapter 19

## Facial Expression Recognition

Yingli Tian, Takeo Kanade, and Jeffrey F. Cohn

### 19.1 Introduction

Facial expressions are the facial changes in response to a person's internal emotional states, intentions, or social communications. Facial expression analysis has been an active research topic for behavioral scientists since the work of Darwin in 1872 [21, 26, 29, 83]. Suwa et al. [90] presented an early attempt to automatically analyze facial expressions by tracking the motion of 20 identified spots on an image sequence in 1978. After that, much progress has been made to build computer systems to help us understand and use this natural form of human communication [5, 7, 8, 17, 23, 32, 43, 45, 57, 64, 77, 92, 95, 106–108, 110].

In this chapter, facial expression analysis refers to computer systems that attempt to automatically analyze and recognize facial motions and facial feature changes from visual information. Sometimes the facial expression analysis has been confused with emotion analysis in the computer vision domain. For emotion analysis, higher level knowledge is required. For example, although facial expressions can convey emotion, they can also express intention, cognitive processes, physical effort, or other intra- or interpersonal meanings. Interpretation is aided by context, body gesture, voice, individual differences, and cultural factors as well as by facial configuration and timing [11, 79, 80]. Computer facial expression analysis systems need to analyze the facial actions regardless of context, culture, gender, and so on.

---

Y. Tian (✉)

Department of Electrical Engineering, The City College of New York, New York, NY 10031, USA

e-mail: [ytian@ccny.cuny.edu](mailto:ytian@ccny.cuny.edu)

T. Kanade

Robotics Institute, Carnegie Mellon University, Pittsburgh, PA 15213, USA

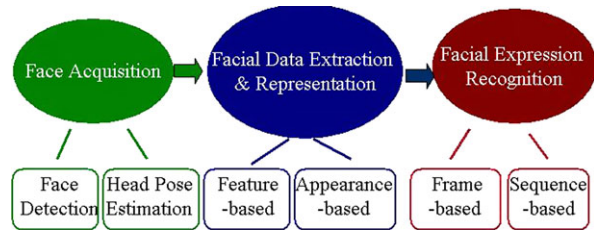
e-mail: [tk@cs.cmu.edu](mailto:tk@cs.cmu.edu)

J.F. Cohn

Department of Psychology, University of Pittsburgh, Pittsburgh, PA 15260, USA

e-mail: [jeffcohn@pitt.edu](mailto:jeffcohn@pitt.edu)

**Fig. 19.1** Basic structure of facial expression analysis systems



The accomplishments in the related areas such as psychological studies, human movement analysis, face detection, face tracking, and recognition make the automatic facial expression analysis possible. Automatic facial expression analysis can be applied in many areas such as emotion and paralinguistic communication, clinical psychology, psychiatry, neurology, pain assessment, lie detection, intelligent environments, and multimodal human computer interface (HCI).

## 19.2 Principles of Facial Expression Analysis

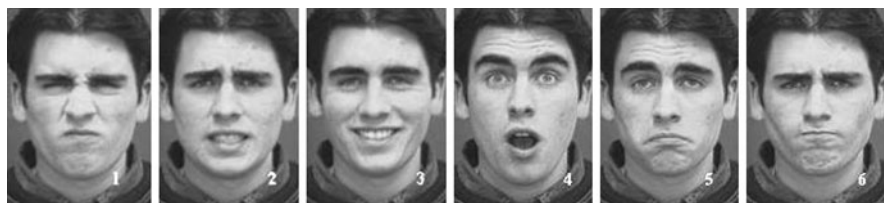
### 19.2.1 Basic Structure of Facial Expression Analysis Systems

Facial expression analysis includes both measurement of facial motion and recognition of expression. The general approach to automatic facial expression analysis (AFEA) consists of three steps (Fig. 19.1): face acquisition, facial data extraction and representation, and facial expression recognition.

Face acquisition is a processing stage to automatically find the face region for the input images or sequences. It can be a detector to detect face for each frame or just detect face in the first frame and then track the face in the remainder of the video sequence. To handle large head motion, the head finder, head tracking, and pose estimation can be applied to a facial expression analysis system.

After the face is located, the next step is to extract and represent the facial changes caused by facial expressions. In facial feature extraction for expression analysis, there are mainly two types of approaches: geometric feature-based methods and appearance-based methods. The geometric facial features present the shape and locations of facial components (including mouth, eyes, brows, nose, etc.). The facial components or facial feature points are extracted to form a feature vector that represents the face geometry. With appearance-based methods, image filters, such as Gabor wavelets, are applied to either the whole-face or specific regions in a face image to extract a feature vector. Depending on the different facial feature extraction methods, the effects of in-plane head rotation and different scales of the faces can be eliminated by face normalization before the feature extraction or by feature representation before the step of expression recognition.

Facial expression recognition is the last stage of AFEA systems. The facial changes can be identified as facial action units or prototypic emotional expressions (see Sect. 19.3.1 for definitions). Depending on whether the temporal information is



**Fig. 19.2** Emotion-specified facial expression (posed images from database [49]). 1, disgust; 2, fear; 3, joy; 4, surprise; 5, sadness; 6, anger

used, in this chapter we classify a recognition approach as frame-based or sequence-based.

## 19.2.2 Organization of the Chapter


















This chapter introduces recent advances in facial expression analysis. The first part discusses general structure of AFEA systems. The second part describes the problem space for facial expression analysis. This space includes multiple dimensions: level of description, individual differences in subjects, transitions among expressions, intensity of facial expression, deliberate versus spontaneous expression, head orientation and scene complexity, image acquisition and resolution, reliability of ground truth, databases, and the relation to other facial behaviors or nonfacial behaviors. We note that most work to date has been confined to a relatively restricted region of this space. The last part of this chapter is devoted to a description of more specific approaches and the techniques used in recent advances. They include the techniques for face acquisition, facial data extraction and representation, facial expression recognition, and multimodal expression analysis. The chapter concludes with a discussion assessing the current status, future possibilities, and open questions about automatic facial expression analysis.

## 19.3 Problem Space for Facial Expression Analysis

### 19.3.1 Level of Description

With few exceptions [17, 23, 34, 95], most AFEA systems attempt to recognize a small set of prototypic emotional expressions as shown in Fig. 19.2 (i.e., disgust, fear, joy, surprise, sadness, anger). This practice may follow from the work of Darwin [21] and more recently Ekman and Friesen [27, 28] and Izard et al. [48] who proposed that emotion-specified expressions have corresponding prototypic facial expressions. In everyday life, however, such prototypic expressions occur relatively infrequently. Instead, emotion more often is communicated by subtle changes in one

**Table 19.1** FACS action units (AU). AUs with “\*” indicate that the criteria have changed for this AU, that is, AU 25, 26, and 27 are now coded according to criteria of intensity (25A-E), and AU 41, 42, and 43 are now coded according to criteria of intensity

Upper Face Action Units					
AU 1	AU 2	AU 4	AU 5	AU 6	AU 7
					
Inner Brow Raiser	Outer Brow Raiser	Brow Lowerer	Upper Lid Raiser	Cheek Raiser	Lid Tightener
*AU 41	*AU 42	*AU 43	AU 44	AU 45	AU 46
					
Lid Droop	Slit	Eyes Closed	Squint	Blink	Wink
Lower Face Action Units					
AU 9	AU 10	AU 11	AU 12	AU 13	AU 14
					
Nose Wrinkler	Upper Lip Raiser	Nasolabial Deepener	Lip Corner Puller	Cheek Puffer	Dimpler
AU 15	AU 16	AU 17	AU 18	AU 20	AU 22
					
Lip Corner Depressor	Lower Lip Depressor	Chin Raiser	Lip Puckerer	Lip Stretcher	Lip Funneler
AU 23	AU 24	*AU 25	*AU 26	*AU 27	AU 28
					
Lip Tightener	Lip Pressor	Lips Part	Jaw Drop	Mouth Stretch	Lip Suck





















or a few discrete facial features, such as tightening of the lips in anger or obliquely lowering the lip corners in sadness [12]. Change in isolated features, especially in the area of the eyebrows or eyelids, is typical of paralinguistic displays; for instance, raising the brows signals greeting [25]. To capture such subtlety of human emotion and paralinguistic communication, automated recognition of fine-grained changes in facial expression is needed. The facial action coding system (FACS: [29]) is a human-observer-based system designed to detect subtle changes in facial features. Viewing videotaped facial behavior in slow motion, trained observers can manually FACS code all possible facial displays, which are referred to as action units and may occur individually or in combinations.

FACS consists of 44 action units. Thirty are anatomically related to contraction of a specific set of facial muscles (Table 19.1) [22]. The anatomic basis of the remaining 14 is unspecified (Table 19.2). These 14 are referred to in FACS as miscellaneous actions. Many action units may be coded as symmetrical or asymmetrical. For action units that vary in intensity, a 5-point ordinal scale is used to measure the degree of muscle contraction. Table 19.3 shows some examples of combinations of FACS action units.

**Table 19.2** Miscellaneous actions

AU	Description
8	Lips toward
19	Tongue show
21	Neck tighten
29	Jaw thrust
30	Jaw sideways
31	Jaw clench
32	Bite lip
33	Blow
34	Puff
35	Cheek suck
36	Tongue bulge
37	Lip wipe
38	Nostril dilate
39	Nostril compress

**Table 19.3** Some examples of combination of FACS action units

AU 1+2	AU 1+4	AU 4+5	AU 1+2+4	AU 1+2+5
				
AU 1+6	AU 6+7	AU 1+2+5+6+7	AU 23+24	AU 9+17
				
AU 9+25	AU 9+17+23+24	AU 10+17	AU 10+25	AU 10+15+17
				
AU 12+25	AU 12+26	AU 15+17	AU 17+23+24	AU 20+25
				

Although Ekman and Friesen proposed that specific combinations of FACS action units represent prototypic expressions of emotion, emotion-specified expressions are not part of FACS; they are coded in separate systems, such as the emotional facial action system (EMFACS) [41]. FACS itself is purely descriptive and includes no inferential labels. By converting FACS codes to EMFACS or similar systems, face images may be coded for emotion-specified expressions (e.g., joy or anger) as well as for more molar categories of positive or negative emotion [65].

### 19.3.2 Individual Differences in Subjects

Face shape, texture, color, and facial and scalp hair vary with sex, ethnic background, and age [33, 119]. Infants, for instance, have smoother, less textured skin and often lack facial hair in the brows or scalp. The eye opening and contrast between iris and sclera differ markedly between Asians and Northern Europeans, which may affect the robustness of eye tracking and facial feature analysis more generally. Beards, eyeglasses, or jewelry may obscure facial features. Such individual differences in appearance may have important consequences for face analysis. Few attempts to study their influence exist. An exception was a study by Zlochowier et al. [119], who found that algorithms for optical flow and high-gradient component detection that had been optimized for young adults performed less well when used in infants. The reduced texture of infants' skin, their increased fatty tissue, juvenile facial conformation, and lack of transient furrows may all have contributed to the differences observed in face analysis between infants and adults.

In addition to individual differences in appearance, there are individual differences in expressiveness, which refers to the degree of facial plasticity, morphology, frequency of intense expression, and overall rate of expression. Individual differences in these characteristics are well established and are an important aspect of individual identity [61] (these individual differences in expressiveness and in biases for particular facial actions are sufficiently strong that they may be used as a biometric to augment the accuracy of face recognition algorithms [19]). An extreme example of variability in expressiveness occurs in individuals who have incurred damage either to the facial nerve or central nervous system [75, 99]. To develop algorithms that are robust to individual differences in facial features and behavior, it is essential to include a large sample of varying ethnic background, age, and sex, which includes people who have facial hair and wear jewelry or eyeglasses and both normal and clinically impaired individuals.

### 19.3.3 Transitions Among Expressions

A simplifying assumption in facial expression analysis is that expressions are singular and begin and end with a neutral position. In reality, facial expression is more complex, especially at the level of action units. Action units may occur in combinations or show serial dependence. Transitions from action units or combination of actions to another may involve no intervening neutral state. Parsing the stream of behavior is an essential requirement of a robust facial analysis system, and training data are needed that include dynamic combinations of action units, which may be either additive or nonadditive.

As shown in Table 19.3, an example of an additive combination is smiling (AU 12) with mouth open, which would be coded as AU 12 + 25, AU 12 + 26, or AU 12 + 27 depending on the degree of lip parting and whether and how far the mandible was lowered. In the case of AU 12 + 27, for instance, the facial analysis

system would need to detect transitions among all three levels of mouth opening while continuing to recognize AU 12, which may be simultaneously changing in intensity.

Nonadditive combinations represent further complexity. Following usage in speech science, we refer to these interactions as co-articulation effects. An example is the combination AU 12 + 15, which often occurs during embarrassment. Although AU 12 raises the cheeks and lip corners, its action on the lip corners is modified by the downward action of AU 15. The resulting appearance change is highly dependent on timing. The downward action of the lip corners may occur simultaneously or sequentially. The latter appears to be more common [85]. To be comprehensive, a database should include individual action units and both additive and nonadditive combinations, especially those that involve co-articulation effects. A classifier trained only on single action units may perform poorly for combinations in which co-articulation effects occur.

### ***19.3.4 Intensity of Facial Expression***

Facial actions can vary in intensity. Manual FACS coding, for instance, uses a 3- or more recently a 5-point intensity scale to describe intensity variation of action units (for psychometric data, see Sayette et al. [82]). Some related action units, moreover, function as sets to represent intensity variation. In the eye region, action units 41, 42, and 43 or 45 can represent intensity variation from slightly drooped to closed eyes. Several computer vision researchers proposed methods to represent intensity variation automatically. Essa and Pentland [32] represented intensity variation in smiling using optical flow. Kimura and Yachida [50] and Lien et al. [56] quantified intensity variation in emotion-specified expression and in action units, respectively. These authors did not, however, attempt the more challenging step of discriminating intensity variation within types of facial actions. Instead, they used intensity measures for the more limited purpose of discriminating between different types of facial actions. Tian et al. [94] compared manual and automatic coding of intensity variation. Using Gabor features and an artificial neural network, they discriminated intensity variation in eye closure as reliably as human coders did. Recently, Bartlett and colleagues [5] investigated action unit intensity by analyzing facial expression dynamics. They performed a correlation analysis to explicitly measure the relationship between the output margin of the learned classifiers and expression intensity. Yang et al. [111] converted the problem of intensity estimation to a ranking problem, which is modeled by the RankBoost. They employed the output ranking score for intensity estimation. These findings suggest that it is feasible to automatically recognize intensity variation within types of facial actions. Regardless of whether investigators attempt to discriminate intensity variation within facial actions, it is important that the range of variation be described adequately. Methods that work for intense expressions may generalize poorly to ones of low intensity.

### ***19.3.5 Deliberate Versus Spontaneous Expression***

Most face expression data have been collected by asking subjects to perform a series of expressions. These directed facial action tasks may differ in appearance and timing from spontaneously occurring behavior [30]. Deliberate and spontaneous facial behavior are mediated by separate motor pathways, the pyramidal and extrapyramidal motor tracks, respectively [75]. As a consequence, fine-motor control of deliberate facial actions is often inferior and less symmetrical than what occurs spontaneously. Many people, for instance, are able to raise their outer brows spontaneously while leaving their inner brows at rest; few can perform this action voluntarily. Spontaneous depression of the lip corners (AU 15) and raising and narrowing the inner corners of the brow (AU 1 + 4) are common signs of sadness. Without training, few people can perform these actions deliberately, which incidentally is an aid to lie detection [30]. Differences in the temporal organization of spontaneous and deliberate facial actions are particularly important in that many pattern recognition approaches, such as hidden Markov modeling, are highly dependent on the timing of the appearance change. Unless a database includes both deliberate and spontaneous facial actions, it will likely prove inadequate for developing face expression methods that are robust to these differences.

### ***19.3.6 Head Orientation and Scene Complexity***

Face orientation relative to the camera, the presence and actions of other people, and background conditions may influence face analysis. In the face recognition literature, face orientation has received deliberate attention. The FERET database [76], for instance, includes both frontal and oblique views, and several specialized databases have been collected to try to develop methods of face recognition that are invariant to moderate change in face orientation [100]. In the face expression literature, use of multiple perspectives is rare; and relatively less attention has been focused on the problem of pose invariance. Most researchers assume that face orientation is limited to in-plane variation [3] or that out-of-plane rotation is small [57, 68, 77, 95]. In reality, large out-of-plane rotation in head position is common and often accompanies change in expression. Kraut and Johnson [54] found that smiling typically occurs while turning toward another person. Camras et al. [10] showed that infant surprise expressions often occur as the infant pitches her head back. To develop pose invariant methods of face expression analysis, image data are needed in which facial expression changes in combination with significant non-planar change in pose. Some efforts have been made to handle large out-of-plane rotation in head position [5, 20, 97, 104].

Scene complexity, such as background and the presence of other people, potentially influences accuracy of face detection, feature tracking, and expression recognition. Most databases use image data in which the background is neutral or has a



consistent pattern and only a single person is present in the scene. In natural environments, multiple people interacting with each other are likely, and their effects need to be understood. Unless this variation is represented in training data, it will be difficult to develop and test algorithms that are robust to such variation.





### ***19.3.7 Image Acquisition and Resolution***

The image acquisition procedure includes several issues, such as the properties and number of video cameras and digitizer, the size of the face image relative to total image dimensions, and the ambient lighting. All of these factors may influence facial expression analysis. Images acquired in low light or at coarse resolution can provide less information about facial features. Similarly, when the face image size is small relative to the total image size, less information is available. NTSC cameras record images at 30 frames per second, The implications of down-sampling from this rate are unknown. Many algorithms for optical flow assume that pixel displacement between adjacent frames is small. Unless they are tested at a range of sampling rates, the robustness to sampling rate and resolution cannot be assessed.

Within an image sequence, changes in head position relative to the light source and variation in ambient lighting have potentially significant effects on face expression analysis. A light source above the subject's head causes shadows to fall below the brows, which can obscure the eyes, especially for subjects with pronounced bone structure or hair. Methods that work well in studio lighting may perform poorly in more natural lighting (e.g., through an exterior window) when the angle of lighting changes across an image sequence. Most investigators use single-camera setups, which is problematic when a frontal orientation is not required. With image data from a single camera, out-of-plane rotation may be difficult to standardize. For large out-of-plane rotation, multiple cameras may be required. Multiple camera setups can support three dimensional (3D) modeling and in some cases ground truth with which to assess the accuracy of image alignment. Pantic and Rothkrantz [70] were the first to use two cameras mounted on a headphone-like device; one camera is placed in front of the face and the other on the right side of the face. The cameras are moving together with the head to eliminate the scale and orientation variance of the acquired face images.

Image resolution is another concern. Professional grade PAL cameras, for instance, provide very high resolution images. By contrast, security cameras provide images that are seriously degraded. Although postprocessing may improve image resolution, the degree of potential improvement is likely limited. Also the effects of post processing for expression recognition are not known. Table 19.4 shows a face at different resolutions. Most automated face processing tasks should be possible for a  $69 \times 93$  pixel image. At  $48 \times 64$  pixels the facial features such as the corners of the eyes and the mouth become hard to detect. The facial expressions may be recognized at  $48 \times 64$  and are not recognized at  $24 \times 32$  pixels. Algorithms that work well at optimal resolutions of full face frontal images and studio lighting

**Table 19.4** A face at different resolutions. All images are enlarged to the same size. At  $48 \times 64$  pixels the facial features such as the corners of the eyes and the mouth become hard to detect. Facial expressions are not recognized at  $24 \times 32$  pixels [97]

				
Face Process	96 x 128	69 x 93	48 x 64	24 x 32
Detect?	Yes	Yes	Yes	Yes
Pose?	Yes	Yes	Yes	Yes
Recognize?	Yes	Yes	Yes	Maybe
Features?	Yes	Yes	Maybe	No
Expressions?	Yes	Yes	Maybe	No

can be expected to perform poorly when recording conditions are degraded or images are compressed. Without knowing the boundary conditions of face expression algorithms, comparative performance is difficult to assess. Algorithms that appear superior within one set of boundary conditions may perform more poorly across the range of potential applications. Appropriate data with which these factors can be tested are needed.

### 19.3.8 Reliability of Ground Truth

When training a system to recognize facial expression, the investigator assumes that training and test data are accurately labeled. This assumption may or may not be accurate. Asking subjects to perform a given action is no guarantee that they will. To ensure internal validity, expression data must be manually coded, and the reliability of the coding verified. Interobserver reliability can be improved by providing rigorous training to observers and monitoring their performance. FACS coders must pass a standardized test, which ensures (initially) uniform coding among international laboratories. Monitoring is best achieved by having observers independently code a portion of the same data. As a general rule, 15% to 20% of data should be comparison-coded. To guard against drift in coding criteria [62], restandardization is important. When assessing reliability, coefficient kappa [36] is preferable to raw percentage of agreement, which may be inflated by the marginal frequencies of codes. Kappa quantifies interobserver agreement after correcting for the level of agreement expected by chance.

### 19.3.9 Databases

Because most investigators have used relatively limited data sets, the generalizability of different approaches to facial expression analysis remains unknown. In most data sets, only relatively global facial expressions (e.g., joy or anger) have been considered, subjects have been few in number and homogeneous with respect to age

and ethnic background, and recording conditions have been optimized. Approaches to facial expression analysis that have been developed in this way may transfer poorly to applications in which expressions, subjects, contexts, or image properties are more variable. In the absence of comparative tests on common data, the relative strengths and weaknesses of different approaches are difficult to determine. In the areas of face and speech recognition, comparative tests have proven valuable [76], and similar benefits would likely accrue in the study of facial expression analysis. A large, representative test-bed is needed with which to evaluate different approaches. We list several databases for facial expression analysis in Sect. 19.4.5.

### ***19.3.10 Relation to Other Facial Behavior or Nonfacial Behavior***

Facial expression is one of several channels of nonverbal communication. Contraction of the muscle *zygomaticus major* (AU 12), for instance, is often associated with positive or happy vocalizations, and smiling tends to increase vocal fundamental frequency [16]. Also facial expressions often occur during conversations. Both expressions and conversations can cause facial changes. Few research groups, however, have attempted to integrate gesture recognition broadly defined across multiple channels of communication [44, 45]. An important question is whether there are advantages to early rather than late integration [38]. Databases containing multimodal expressive behavior afford the opportunity for integrated approaches to analysis of facial expression, prosody, gesture, and kinetic expression.

### ***19.3.11 Summary and Ideal Facial Expression Analysis Systems***

The problem space for facial expression includes multiple dimensions. An ideal facial expression analysis system has to address all these dimensions, and it outputs accurate recognition results. In addition, the ideal facial expression analysis system must perform automatically and in real-time for all stages (Fig. 19.1). So far, several systems can recognize expressions in real time [53, 68, 97]. We summarize the properties of an ideal facial expression analysis system in Table 19.5.

## **19.4 Recent Advances**

For automatic facial expression analysis, Suwa et al. [90] presented an early attempt in 1978 to analyze facial expressions by tracking the motion of 20 identified spots on an image sequence. Considerable progress had been made since 1990 in related technologies such as image analysis and pattern recognition that make AFEA possible. Samal and Iyengar [81] surveyed the early work (before 1990) about automatic recognition and analysis of human face and facial expression. Two survey papers

**Table 19.5** Properties of an ideal facial expression analysis system

Robustness	
Rb1	Deal with subjects of different age, gender, ethnicity
Rb2	Handle lighting changes
Rb3	Handle large head motion
Rb4	Handle occlusion
Rb5	Handle different image resolution
Rb6	Recognize all possible expressions
Rb7	Recognize expressions with different intensity
Rb8	Recognize asymmetrical expressions
Rb9	Recognize spontaneous expressions
Automatic process	
Am1	Automatic face acquisition
Am2	Automatic facial feature extraction
Am3	Automatic expression recognition
Real-time process	
Rt1	Real-time face acquisition
Rt2	Real-time facial feature extraction
Rt3	Real-time expression recognition
Autonomic Process	
An1	Output recognition with confidence
An2	Adaptive to different level outputs based on input images

summarized the work (before year 1999) of facial expression analysis [35, 69]. Recently, Zeng et al. [114] surveyed the work (before year 2007) for affect recognition methods including audio, visual and spontaneous expressions. In this chapter, instead of giving a comprehensive survey of facial expression analysis literature, we explore the recent advances in facial expression analysis based on four problems: (1) face acquisition, (2) facial feature extraction and representation, (3) facial expression recognition, and (4) multimodal expression analysis. In addition, we list the public available databases for expression analysis.

Many efforts have been made for facial expression analysis [4, 5, 8, 13, 15, 18, 20, 23, 32, 34, 35, 37, 45, 58–60, 67, 69, 70, 87, 95, 96, 102, 104, 107, 110–115, 117]. Because most of the work are summarized in the survey papers [35, 69, 114], here we focus on the recent research in automatic facial expression analysis which tends to follow these directions:

- Build more robust systems for face acquisition, facial data extraction and representation, and facial expression recognition to handle head motion (in-plane and out-of-plane), occlusion, lighting changes, and lower intensity of expressions
- Employ more facial features to recognize more expressions and to achieve a higher recognition rate
- Recognize facial action units and their combinations rather than emotion-specified expressions
- Recognize action units as they occur spontaneously
- Develop fully automatic and real-time AFEA systems
- Analyze emotion portrayals by combining multimodal features such as facial expression, vocal expression, gestures, and body movements

### ***19.4.1 Face Acquisition***

With few exceptions, most AFEA research attempts to recognize facial expressions only from frontal-view or near frontal-view faces [51, 70]. Kleck and Mendolia [51] first studied the decoding of profile versus full-face expressions of affect by using three perspectives (a frontal face, a 90° right profile, and a 90° left profile). Forty-eight decoders viewed the expressions from 64 subjects in one of the three facial perspectives. They found that the frontal faces elicited higher intensity ratings than profile views for negative expressions. The opposite was found for positive expressions. Pantic and Rothkrantz [70] used dual-view facial images (a full-face and a 90° right profile) which are acquired by two cameras mounted on the user's head. They did not compare the recognition results by using only the frontal view and the profile. So far, it is unclear how many expressions can be recognized by side-view or profile faces. Because the frontal-view face is not always available in real environments, the face acquisition methods should detect both frontal and nonfrontal view faces in an arbitrary scene.

To handle out-of-plane head motion, face can be obtained by face detection, 2D or 3D face tracking, or head pose detection. Nonfrontal view faces are warped or normalized to frontal view for expression analysis.

#### **19.4.1.1 Face Detection**

Many face detection methods have been developed to detect faces in an arbitrary scene [47, 55, 72, 78, 86, 89, 101]. Most of them can detect only frontal and near-frontal views of faces. Heisele et al. [47] developed a component-based, trainable system for detecting frontal and near-frontal views of faces in still gray images. Rowley et al. [78] developed a neural network based system to detect frontal-view face. Viola and Jones [101] developed a robust real-time face detector based on a set of rectangle features.

To handle out-of-plane head motion, some researchers developed face detectors to detect face from different views [55, 72, 86]. Pentland et al. [72] detected faces by

using the view-based and modular eigenspace method. Their method runs real-time and can handle varying head positions. Schneiderman and Kanade [86] proposed a statistical method for 3D object detection that can reliably detect human faces with out-of-plane rotation. They represent the statistics of both *object* appearance and *nonobject* appearance using a product of histograms. Each histogram represents the joint statistics of a subset of wavelet coefficients and their position on the object. Li et al. [55] developed an AdaBoost-like approach to detect faces with multiple views. A detail survey about face detection can be found in paper [109].

Some facial expression analysis systems use the face detector which developed by Viola et al. [101] to detect face for each frame [37]. Some systems [18, 67, 94–96, 104] assume that the first frame of the sequence is frontal and expressionless. They detect faces only in the first frame and then perform feature tracking or head tracking for the remaining frames of the sequence.

#### 19.4.1.2 Head Pose Estimation

In a real environment, out-of-plane head motion is common for facial expression analysis. To handle the out-of-plane head motion, head pose estimation can be employed. The methods for estimating head pose can be classified as 3D model-based methods [1, 91, 98, 104] and 2D image-based methods [9, 97, 103, 118].

**3D Model-Based Method** Many systems employ a 3D model based method to estimate head pose [4, 5, 15, 18, 20, 67, 102, 104]. Bartlett et al. [4, 5] used a canonical wire-mesh face model to estimate face geometry and 3D pose from hand-labeled feature points. In papers [15, 102], the authors used an explicit 3D wireframe face model to track geometric facial features defined on the model [91]. The 3D model is fitted to the first frame of the sequence by manually selecting landmark facial features such as corners of the eyes and mouth. The generic face model, which consists of 16 surface patches, is warped to fit the selected facial features. To estimate the head motion and deformations of facial features, a two-step process is used. The 2D image motion is tracked using template matching between frames at different resolutions. From the 2D motions of many points on the face model, the 3D head motion then is estimated by solving an overdetermined system of equations of the projective motions in the least-squares sense [15].

In paper [104], a cylindrical head model is used to automatically estimate the 6 degrees of freedom (dof) of head motion in realtime. An active appearance model (AAM) method is used to automatically map the cylindrical head model to the face region, which is detected by face detection [78], as the initial appearance template. For any given frame, the template is the head image in the previous frame that is projected onto the cylindrical model. Then the template is registered with the head appearance in the given frame to recover the full motion of the head. They first use the iteratively reweighted least squares technique [6] to deal with nonrigid motion and occlusion. Second, they update the template dynamically in order to deal with gradual changes in lighting and self-occlusion. This enables the system to work well











**Fig. 19.3** Example of 3D head tracking, including re-registration after losing the head

even when most of the face is occluded. Because head poses are recovered using templates that are constantly updated and the pose estimated for the current frame is used in estimating the pose in the next frame, errors would accumulate unless otherwise prevented. To solve this problem, the system automatically selects and stores one or more frames and associated head poses from the tracked images in the sequence (usually including the initial frame and pose) as references. Whenever the difference between the estimated head pose and that of a reference frame is less than a preset threshold, the system rectifies the current pose estimate by re-registering this frame with the reference. The reregistration prevents errors from accumulating and enables the system to recover head pose when the head reappears after occlusion, such as when the head moves momentarily out of the camera's view. On-line tests suggest that the system could work robustly for an indefinite period of time. It was also quantitatively evaluated in image sequences that include maximum pitch and yaw as large as 40 and 75 degrees, respectively. The precision of recovered motion was evaluated with respect to the ground truth obtained by a precise position and orientation measurement device with markers attached to the head and found to be highly consistent (e.g., for maximum yaw of 75 degrees, absolute error averaged 3.86 degrees). An example of the 3D head tracking is shown in Fig. 19.3 including re-registration after losing the head. More details can be found in paper [104].

**2D Image-Based Method** To handle the full range of head motion for expression analysis, Tian et al. [97] detected the head instead of the face. The head detection uses the smoothed silhouette of the foreground object as segmented using background subtraction and computing the *negative curvature minima* (NCM) points of the silhouette. Other head detection techniques that use silhouettes can be found elsewhere [42, 46].

**Table 19.6** Definitions and examples of the three head pose classes: frontal or near frontal view, side view or profile, and others, such as back of the head or occluded faces. The expression analysis process is applied to only the frontal and near-frontal view faces [9, 97]

Poses	Frontal or near frontal			Side view or profile		Others		
Definitions	Both eyes and lip corners are visible			One eye or one lip corner is occluded		No enough facial features		
Examples								

After the head is located, the head image is converted to gray-scale, histogram-equalized, and resized to the estimated resolution. Then a three-layer neural network (NN) is employed to estimate the head pose. The inputs to the network are the processed head image. The outputs are the three head poses: (1) frontal or near frontal view, (2) side view or profile, (3) others, such as back of the head or occluded face (Table 19.6). In the frontal or near frontal view, both eyes and lip corners are visible. In the side view or profile, at least one eye or one corner of the mouth becomes self-occluded because of the head. The expression analysis process is applied only to the frontal and near-frontal view faces. Their system performs well even with very low resolution of face images.

19.4.2 Facial Feature Extraction and Representation

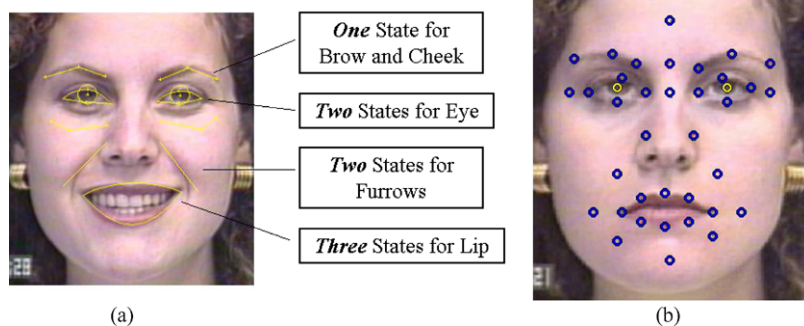
After the face is obtained, the next step is to extract facial features. Two types of features can be extracted: geometric features and appearance features. Geometric features present the shape and locations of facial components (including mouth, eyes, brows, and nose). The facial components or facial feature points are extracted to form a feature vector that represents the face geometry. The appearance features present the appearance (skin texture) changes of the face, such as wrinkles and furrows. The appearance features can be extracted on either the whole-face or specific regions in a face image.

To recognize facial expressions, an AEFA system can use geometric features only [15, 20, 70], appearance features only [5, 37, 59], or hybrid features (both geometric and appearance features) [23, 95, 96, 102]. The research shows that using hybrid features can achieve better results for some expressions.

To remove the effects of variation in face scale, motion, lighting, and other factors, one can first align and normalize the face to a standard face (2D or 3D) manually or automatically [23, 37, 57, 102], and then obtain normalized feature measurements by using a reference image (neutral face) [95].



## Multi-State Models for Geometric Feature Extraction



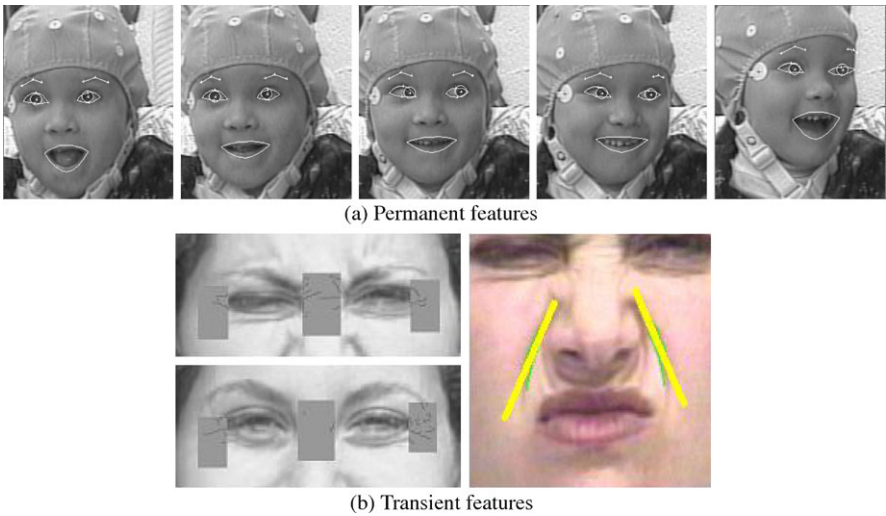
**Fig. 19.4** Facial feature extraction for expression analysis [95]. **a** Multistate models for geometric feature extraction. **b** Locations for calculating appearance features

### 19.4.2.1 Geometric Feature Extraction

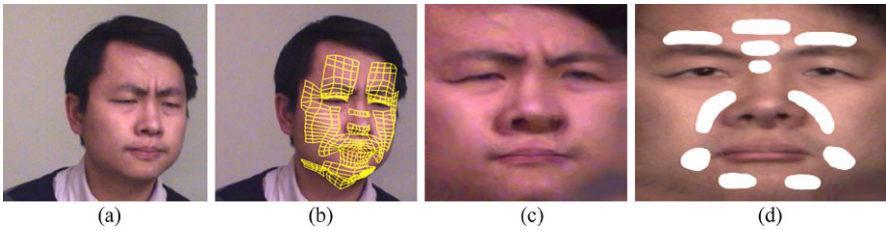
As shown in Fig. 19.4, in order to detect and track changes of facial components in near frontal face images, Tian et al. develop multi-state models to extract the geometric facial features. A three-state lip model describes the lip state: open, closed, tightly closed. A two-state model (open or closed) is used for each of the eyes. Each brow and cheek has a one-state model. Some appearance features, such as *nasolabial furrows* and *crows-feet wrinkles* (Fig. 19.5b), are represented explicitly by using two states: present and absent. Given an image sequence, the region of the face and approximate location of individual face features are detected automatically in the initial frame [78]. The contours of the face features and components then are adjusted manually in the initial frame. After the initialization, all face feature changes are automatically detected and tracked in the image sequence. The system groups 15 parameters for the upper face and 9 parameters for the lower face, which describe shape, motion, and state of face components and furrows. To remove the effects of variation in planar head motion and scale between image sequences in face size, all parameters are computed as ratios of their current values to that in the reference frame. Details of geometric feature extraction and representation can be found in paper [95].

Automatic active appearance model (AAM) mapping can be employed to reduce the manual preprocessing of the geometric feature initialization [66, 105]. Xiao et al. [104] performed the 3D head tracking to handle large out-of plane head motion (Sect. 19.4.1) and track nonrigid features. Once the head pose is recovered, the face region is stabilized by transforming the image to a common orientation for expression recognition [18, 67].

The systems in [15, 102] use an explicit 3D wireframe face model to track geometric facial features defined on the model [91]. The 3D model is fitted to the first frame of the sequence by manually selecting landmark facial features such as corners of the eyes and mouth. The generic face model, which consists of 16 surface



**Fig. 19.5** Example results of feature extraction [95]. **a** Permanent feature extraction (eyes, brows, and mouth). **b** Transient feature extraction (crows-feet wrinkles, wrinkles at nasal root, and nasolabial furrows)



**Fig. 19.6** Example of feature extraction [102]. **a** Input video frame. **b** Snapshot of the geometric tracking system. **c** Extracted texture map. **d** Selected facial regions for appearance feature extraction [102]

patches, is warped to fit the selected facial features. Figure 19.6b shows an example of the geometric feature extraction of paper [102].

**19.4.2.2 Appearance Feature Extraction**

Gabor wavelets [22] are widely used to extract the facial appearance changes as a set of multiscale and multiorientation coefficients. The Gabor filter may be applied to specific locations on a face [59, 94, 96, 116] or to the whole face image [4, 23, 37]. Zhang et al. [116] was the first to compare two type of features to recognize expressions, the geometric positions of 34 fiducial points on a face and 612 Gabor wavelet coefficients extracted from the face image at these 34 fiducial points. The recognition rates for six emotion-specified expressions (e.g., joy and anger) were

significantly higher for Gabor wavelet coefficients. Donato et al. [23] compared several techniques for recognizing six single upper face AUs and six lower face AUs. These techniques include optical flow, principal component analysis, independent component analysis, local feature analysis, and Gabor wavelet representation. The best performances were obtained using a Gabor wavelet representation and independent component analysis. All of these systems [23, 116] used a manual step to align each input image with a standard face image using the center of the eyes and mouth.

Tian et al. [96] studied geometric features and Gabor coefficients to recognize single AU and AU combinations. In their system, they used 480 Gabor coefficients in the upper face for 20 locations and 432 Gabor coefficients in the lower face for 18 locations (Fig. 19.4). They found that Gabor wavelets work well for single AU recognition for homogeneous subjects without head motion. However, for recognition of AU combinations when image sequences include nonhomogeneous subjects with small head motions, the recognition results are relatively poor if we use only Gabor appearance features. Several factors may account for the difference. First, the previous studies used homogeneous subjects. For instance, Zhang et al. [116] included only Japanese and Donato et al. [23] included only Euro-Americans. Tian et al. use Cohn–Kanade database which contains diverse subjects of European, African, and Asian ancestry. Second, the previous studies recognized emotion-specified expressions or only single AUs. Tian et al. tested the Gabor-wavelet-based method on both single AUs and AU combinations, including nonadditive combinations in which the occurrence of one AU modifies another. Third, the previous studies manually aligned and cropped face images. System [96] omitted this pre-processing step. In summary, using Gabor wavelets alone, recognition is adequate only for AU6, AU43, and AU0. Using geometric features alone, recognition is consistently good and shows high AU recognition rates with the exception of AU7. Combining both Gabor wavelet coefficients and geometric features, the recognition performance increased for all AUs.

In system [4], 3D pose and face geometry is estimated from hand-labeled feature points by using a canonical wire-mesh face model [73]. Once the 3D pose is estimated, faces are rotated to the frontal view and warped to a canonical face geometry. Then, the face images are automatically scaled and cropped to a standard face with a fixed distance between the two eyes. Difference images are obtained by subtracting a neutral expression face. They employed a family of Gabor wavelets at five spatial frequencies and eight orientations to a different image. Instead of specific locations on a face, they apply the Gabor filter to the whole face image. To provide robustness to lighting conditions and to image shifts they employed a representation in which the outputs of two Gabor filters in quadrature are squared and then summed. This representation is known as Gabor energy filters and it models complex cells of the primary visual cortex. Recently, Bartlett and her colleagues extend the system by using fully automatic face and eye detection. For facial expression analysis, they continue employ Gabor wavelets as appearance features [5].

Wen and Huang [102] use the ratio-image based method to extract appearance features, which is independent of a person's face albedo. To limit the effects of the

noise in tracking and individual variation, they extracted the appearance features in facial regions instead of points, and then used the weighted average as the final feature for each region. Eleven regions were defined on the geometric-motion-free texture map of the face (Fig. 19.6d). Gabor wavelets with two spatial frequency and six orientations are used to calculate Gabor coefficients. A 12-dimension appearance feature vector is computed in each of the 11 selected regions by weighted averaging of the Gabor coefficients. To track the face appearance variations, an appearance model (texture image) is trained using a Gaussian mixture model based on exemplars. Then an online adaption algorithm is employed to progressively adapt the appearance model to new conditions such as lighting changes or differences in new individuals. See [102] for details.

### 19.4.3 Facial Expression Recognition

The last step of AFEA systems is to recognize facial expression based on the extracted features. Many classifiers have been applied to expression recognition such as neural network (NN), support vector machines (SVM), linear discriminant analysis (LDA), K-nearest neighbor, multinomial logistic ridge regression (MLR), hidden Markov models (HMM), tree augmented naive Bayes, RankBoost, and others. Some systems use only a rule-based classification based on the definition of the facial actions. Here, we summarize the expression recognition methods to frame-based and sequence-based expression recognition methods. The frame-based recognition method uses only the current frame with or without a reference image (it is mainly a neutral face image) to recognize the expressions of the frame. The sequence-based recognition method uses the temporal information of the sequences to recognize the expressions for one or more frames. Table 19.7 summarizes the recognition methods, recognition rates, recognition outputs, and the databases used in the most recent systems. For the systems that used more classifiers, the best performance for person-independent test has been selected.

**Frame-Based Expression Recognition** Frame-based expression recognition does not use temporal information for the input images. It uses the information of current input image with/without a reference frame. The input image can be a static image or a frame of a sequence that is treated independently. Several methods can be found in the literature for facial expression recognition such as *neural networks* [95, 96, 116], *support vector machines* [4, 37], *linear discriminant analysis* [17], *Bayesian network* [15], and *rule-based classifiers* [70].

Tian et al. [96] employed a neural network-based recognizer to recognize FACS AUs. They used three-layer neural networks with one hidden layer to recognize AUs by a standard back-propagation method [78]. Separate networks are used for the upper and lower face. The inputs can be the normalized geometric features, the appearance feature, or both. The outputs are the recognized AUs. The network is trained to respond to the designated AUs whether they occur alone or in combination. When

**Table 19.7** FACS AU or expression recognition of recent advances. SVM, support vector machines; MLR, multinomial logistic ridge regression; HMM, hidden Markov models; BN, Bayesian network; GMM, Gaussian mixture model; RegRankBoost, RankBoost with l1 regularization

Systems	Recognition methods	Recognition rate	Recognized outputs	Databases
[94–96]	Neural network (frame)	95.5%	16 single AUs and their combinations	Ekman–Hager [31], Cohn–Kanade [49]
[18, 67]	Rule-based (sequence)	100% 57%	Blink, nonblink Brow up, down, and non-motion	Frank–Ekman [40]
[37]	SVM + MLR (frame)	91.5%	6 Basic expressions	Cohn–Kanade [49]
[5]	Adaboost + SVM (sequence)	80.1%	20 facial actions	Frank–Ekman [40]
[15]	BN + HMM (frame & sequence)	73.22% 66.53%	6 Basic expressions 6 Basic expressions	Cohn–Kanade [49] UIUC–Chen [14]
[102]	NN + GMM (frame)	71%	6 Basic expressions	Cohn–Kanade [49]
[111]	RegRankBoost (frame)	88%	6 Basic expressions	Cohn–Kanade [49]

AUs occur in combination, multiple output nodes are excited. To our knowledge, system of [96] was the first system to handle AU combinations. Although several other systems tried to recognize AU combinations [17, 23, 57], they treated each combination as if it were a separate AU. More than 7000 different AU combinations have been observed [83], and a system that can handle AU combinations is more efficient. A overall recognition rate of 95.5% had been achieved for neutral expression and 16 AUs whether they occurred individually or in combinations.

In [37], a two-stage classifier was employed to recognize neutral expression and six emotion-specified expressions. First, SVMs were used for the pairwise classifiers, that is, each SVM is trained to distinguish two emotions. Then they tested several approaches, such as nearest neighbor, a simple voting scheme, and multinomial logistic ridge regression (MLR) to convert the representation produced by the first stage into a probability distribution over six emotion-specified expressions and neutral. The best performance at 91.5% was achieved by MLR.

Wen and Huang [102] also employed a two-stage classifier to recognize neutral expression and six emotion-specified expressions. First, a neural network is used to classify *neutral* and *nonneutral*-like [93]. Then Gaussian mixture models (GMMs) were used for the remaining expressions. The overall average recognition rate was 71% for a people-independent test.

Yang et al. [111] employ RankBoost with l1 regularization for expression recognition. They also estimate the intensity of expressions by using the output ranking scores. For six emotion-specified expressions in Cohn–Kanade database, they achieved 88% recognition rate.

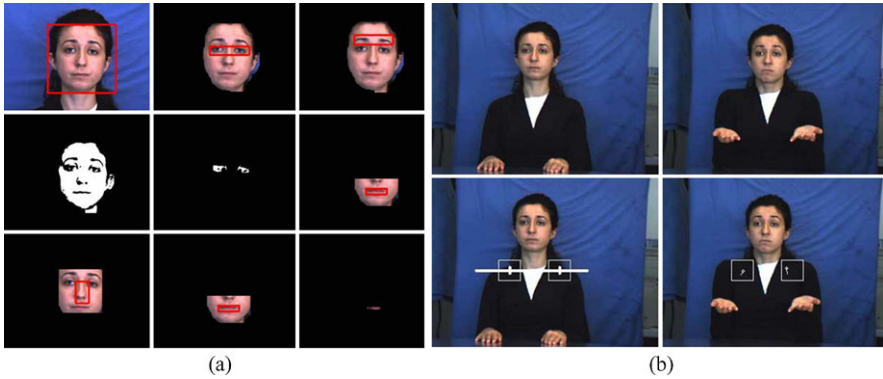
**Sequence-Based Expression Recognition** The sequence-based recognition method uses the temporal information of the sequences to recognize the expressions of one or more frames. To use the temporal information, the techniques such as HMM [4, 15, 17, 57], recurrent neural networks [52, 77], and rule-based classifier [18] were employed in facial expression analysis. The systems of [4, 15, 18] employed a sequence-based classifier. Note that the systems of [4] and [18] are comparative studies for FACS AU recognition in spontaneously occurring behavior by using the same database [40]. In that database, subjects were ethnically diverse, AUs occurred during speech, and out-of-plane motion and occlusion from head motion and glasses were common. So far, only several systems tried to recognize AUs or expression in spontaneously occurring behavior [4, 5, 18, 97].

The system [18] employed a rule-based classifier to recognize AUs of eye and brow in spontaneously occurring behavior by using a number of frames in the sequence. The algorithm achieved an overall accuracy of 98% for three eye behaviors: blink (AU 45), flutter, and no blink (AU 0). *Flutter* is defined as two or more rapidly repeating blinks (AU 45) with only partial eye opening (AU 41 or AU 42) between them. 100% accuracy is achieved between blinks and non-blinks. Accuracy across the three categories in the brow region (brow-up, brow-down, nonbrow motion) was 57%. The number of brow-down actions was too small for reliable point estimates. Omitting brow-down from the analysis, recognition accuracy would increase to 80%. Human FACS coders had similar difficulty with brow-down, agreeing only about 50% in this database. The small number of occurrences was no doubt a factor for FACS coders as well. The combination of occlusion from eyeglasses and correlation of forward head pitch with brow-down complicated FACS coding.

System [4] first employed SVMs to recognize AUs by using Gabor representations. Then they used hidden Markov models (HMMs) to deal with AU dynamics. HMMs were applied in two ways: (1) taking Gabor representations as input, and (2) taking the outputs of SVM as input. When they use Gabor representations as input to train HMMs, the Gabor coefficients were reduced to 100 dimensions per image using PCA. Two HMMs, one for blinks and one for nonblinks were trained and tested using leave-one-out cross-validation. A best performance of 95.7% recognition rate was obtained using five states and three Gaussians. They achieved a 98.1% recognition rate for blink and non-blink using SVM outputs as input to train HMMs for five states and three Gaussians. Accuracy across the three categories in the brow region (brow-up, brow-down, nonbrow motion) was 70.1% (HMMs trained on PCA-reduced Gabors) and 66.9% (HMMs trained on SVM outputs) respectively. Omitting brow-down, the accuracy increases to 90.9% and 89.5%, respectively.

Cohen et al. [15] first evaluated Bayesian network (frame-based) classifiers such as Gaussian naive Bayes (NB-Gaussian), Cauchy naive Bayes (NB-Cauchy), and tree-augmented-naive Bayes (TAN), focusing on changes in distribution assumptions and feature dependency structures. They also proposed a new architecture of HMMs to segment and recognize neutral and six emotion-specified expressions from video sequences. For the person-independent test in the Cohn-Kanade database [49], the best performance at recognition rate of 73.2% was achieved by the TAN classifier. See details in Cohen et al. [15].





**Fig. 19.7** Example of the face and body feature extraction employed in the FABO system [45]. **a** Face features. **b** Body features—shoulder extraction procedure. Shoulder regions found and marked on the neutral frame (*first row*), estimating the movement within the shoulder regions using optical flow (*second row*)

#### 19.4.4 Multimodal Expression Analysis

Facial expression is one of several modes of nonverbal communication. The message value of various modes may differ depending on context and may be congruent or discrepant with each other. Recently, several researchers integrated facial expression analysis with other modes such as gesture, prosody, and speech [20, 44, 45, 84]. Cohn et al. [20] investigated the relation between facial actions and vocal prosody for depression detection. They achieved the same accuracy rate at 79% by using facial actions and vocal prosody respectively. No results are reported for combination. Gunes and Piccardi [45] combined facial actions and body gestures for 9 expression recognition. They found that recognition from fused face and body modalities performs better than that from the face or the body modality alone.

For facial feature extraction in [45], following frame-by-frame face detection, a combination of appearance (e.g., wrinkles) and geometric features (e.g., feature points) is extracted from the face videos. A reference frame with neutral expression is employed for feature comparison. For body feature extraction and tracking, they detected and tracked head, shoulders and hands by using meanshift method from the body videos. Figure 19.7 shows examples of the face and body feature extraction in [45]. A total of 152 features for face modality and 170 features for body modality were used for the detection of face and body temporal segments with various classifiers including both frame-based and sequence-based methods. They tested the system on FABO database [44] and achieved recognition rate at 35.22% by only using face features and 76.87% by only using body features. The recognition rate increased to 85% with combination of both face and body features. More details can be found at [45].

**Table 19.8** Summary of databases for facial expression analysis

Databases	Images/ Videos	Subjects	Expressions	Neutral	Spontaneous	Multimodal	3D data
Cohn–Kanade [49]	videos	210	basic expressions single AUs AU combina- tions	yes	no	frontal face 30° face 30° face	no
FABO [44]	videos	23	9 expressions hand gestures	yes	no	frontal face upper body	no
JAFFE [59]	images	10	6 basic expressions	yes	no	frontal face	no
MMI [71]	images videos	19	single AUs AU combina- tions	yes	no	frontal face profile face	no
RU-FACS [5]	videos	100	AU combina- tions AU	yes	yes	4 face poses speech	no
BU-3DFE [112]	static	100	6 basic expressions	yes	no	face	yes
BU-4DFE [113]	dynamic	101	6 basic expressions	yes	no	face	yes

**19.4.5 Databases for Facial Expression Analysis**

Standard databases play important roles to train, evaluate, and compare different methods and systems for facial expression analysis. There are some public available databases (images or videos) of expression analysis for conducting comparative tests [5, 24, 40, 44, 49, 59, 63, 71, 74, 88, 112, 113]. In this chapter, we summarize several common used standard databases for facial expression analysis in Table 19.8.

Cohn–Kanade AU-Coded Face Expression Database (Cohn–Kanade) [49] is the most commonly used comprehensive database in research on automated facial expression analysis. In Cohn–Kanade database, facial behavior was recorded for two views of faces (frontal view and 30-degree view) in 210 adults between the ages of 18 and 50 years. They were 69% female, 31% male, 81% Euro-American, 13% Afro-American, and 6% other groups. In the database, 1917 image sequences from frontal view videos for 182 subjects have been FACS coded for either target action units or the entire sequence. Japanese Female Facial Expression (JAFPE) Database [59] contains 213 images of 6 basic facial expressions and neutral posed by 10 Japanese female subjects. It is the first downloadable database for facial expression analysis. MMI Facial Expression Database (MMI) [71] contains more than 1500 samples of both static images and image sequences of faces from 19 subjects in frontal and profile views displaying various facial expressions of emotion, single AUs, and AU combinations. It also includes the identification of the temporal



segments (onset, apex, offset) of shown AU and emotion facial displays. The Bi-modal Face and Body Gesture Database (FABO) [44] contains image sequences captured by two synchronized cameras (one for frontal view facial actions, and another for frontal view upper body gestures as shown in Fig. 19.7) from 23 subjects. The database is coded to neutral and nine general expressions (uncertainty, anger, surprise, fear, anxiety, happiness, disgust, boredom, and sadness) based on facial actions and body gestures. The RU-FACS Spontaneous Expression Database (RU-FACS) [5] is a dataset of spontaneous facial behavior with rigorous FACS coding. The dataset consists of 100 subjects participating in a ‘false opinion’ paradigm with speech-related mouth movements and out-of-plane head rotations from four views of face (frontal, left 45°, right 45°, and up about 22°). To date, image sequences from frontal view of 33 subjects have been FACS-coded. The database is being prepared for release. The Binghamton University 3D Facial Expression Database (BU-3DFE) [112] contains 2500 3D facial expression models including neutral and 6 basic expressions from 100 subjects. Associated with each 3D expression model, there are two corresponding facial texture images captured at two views (about +45° and −45°). The BU-4DFE database [113] is extended from a static 3D space (BU-3DFE database) to a dynamic 3D space at a video rate of 25 frames per second. BU-4DFE database contains 606 3D facial expression sequences captured from 101 subjects. Associated with each 3D expression sequence, there is a facial texture video with high resolution of  $1040 \times 1329$  pixels per frame.

## 19.5 Open Questions

Although many recent advances and successes in automatic facial expression analysis have been achieved, as described in the previous sections, many questions remain open, for which answers must be found. Some major points are considered here.

### 1. *How do humans correctly recognize facial expressions?*

Research on human perception and cognition has been conducted for many years, but it is still unclear how humans recognize facial expressions. Which types of parameters are used by humans and how are they processed? By comparing human and automatic facial expression recognition we may be able to advance our understanding of each and discover new ways of improving automatic facial expression recognition.

### 2. *Is it always better to analyze finer levels of expression?*

Although it is often assumed that more fine-grained recognition is preferable, the answer depends on both the quality of the face images and the type of application. Ideally, an AFEA system should recognize all action units and their combinations. In high quality images, this goal seems achievable; emotion-specified expressions then can be identified based on emotion prototypes identified in the psychology literature. For each emotion, prototypic action units have been identified. In lower quality image data, only a subset of action units and emotion-specified expression may be recognized. Recognition of emotion-specified expressions directly may be needed. We seek systems that become “self aware”

about the degree of recognition that is possible based on the information of given images and adjust processing and outputs accordingly. Recognition from coarse-to-fine, for example from emotion-specified expressions to subtle action units, depends on image quality and the type of application. Indeed, for some purposes, it may be sufficient that a system is able to distinguish between positive, neutral, and negative expression, or recognize only a limited number of target action units, such as brow lowering to signal confusion, cognitive effort, or negative affect.

3. *Is there any better way to code facial expressions for computer systems?*

Almost all the existing works have focused on recognition of facial expression, either emotion-specified expressions or FACS coded action units. The emotion-specified expressions describe expressions at a coarse level and are not sufficient for some applications. Although the FACS was designed to detect subtle changes in facial features, it is a human-observer-based system with only limited ability to distinguish intensity variation. Intensity variation is scored at an ordinal level; the interval level measurement is not defined and anchor points may be subjective. Challenges remain in designing a computer-based facial expression coding system with more quantitative definitions.

4. *How do we obtain reliable ground truth?*

Whereas some approaches have used FACS, which is a criterion measure widely used in the psychology community for facial expression analysis, most vision-based work uses emotion-specified expressions. A problem is that emotion-specified expressions are not well defined. The same label may apply to very different facial expressions, and different labels may refer to the same expressions, which confounds system comparisons. Another problem is that the reliability of labels typically is unknown. With few exceptions, investigators have failed to report interobserver reliability and the validity of the facial expressions they have analyzed. Often there is no way to know whether subjects actually showed the target expression or whether two or more judges would agree that the subject showed the target expression. At a minimum, investigators should make explicit labeling criteria and report interobserver agreement for the labels. When the dynamics of facial expression are of interest, temporal resolution should be reported as well. Because intensity and duration measurements are critical, it is important to include descriptive data on these features as well. Unless adequate data about stimuli are reported, discrepancies across studies are difficult to interpret. Such discrepancies could be due to algorithms or to errors in ground truth determination.

5. *How do we recognize facial expressions in real life?*

Real-life facial expression analysis is much more difficult than the posed actions studied predominantly to date. Head motion, low resolution input images, absence of a neutral face for comparison, and low intensity expressions are among the factors that complicate facial expression analysis. Recent works in 3D modeling of spontaneous head motion and action unit recognition in spontaneous facial behavior are exciting developments. How elaborate a head model is required to be in such work is as yet a research question. A cylindrical model

is relatively robust and has proven effective as a part of blink detection system [104], but highly parametric, generic, or even custom-fitted head models may prove necessary for more complete action unit recognition.

Most works to date have used a single, passive camera. Although there are clear advantages to approaches that require only a single passive camera or video source, multiple cameras are feasible in a number of settings and can be expected to provide improved accuracy. Active cameras can be used to acquire high resolution face images [46]. Also, the techniques of super-resolution can be used to obtain higher resolution images from multiple low resolution images [2]. At present, it is an open question how to recognize expressions in situations in which a neutral face is unavailable, expressions are of low intensity, or other facial or nonverbal behaviors, such as occlusion by the hands, are present.

6. *How do we best use the temporal information?*

Almost all works have emphasized recognition of discrete facial expressions, regardless of being defined as emotion-specified expressions or action units. The timing of facial actions may be as important as their configuration. Recent work by our group has shown that intensity and duration of expression vary with context and that the timing of these parameters is highly consistent with automatic movement [85]. Related work suggests that spontaneous and deliberate facial expressions may be discriminated in terms of timing parameters [19], which is consistent with neuropsychological models [75] and may be important to lie detection efforts. Attention to timing is also important in guiding the behavior of computer avatars. Without veridical timing, believable avatars and ones that convey intended emotions and communicative intents may be difficult to achieve.

7. *How may we integrate facial expression analysis with other modalities?*

Facial expression is one of several modes of nonverbal communication. The message value of various modes may differ depending on context and may be congruent or discrepant with each other. An interesting research topic is the integration of facial expression analysis with that of gesture, prosody, and speech. Combining facial features with acoustic features would help to separate the effects of facial actions due to facial expression and those due to speech-related movements. The combination of facial expression and speech can be used to improve speech recognition and multimodal person identification [39].

## 19.6 Conclusions

Five recent trends in automatic facial expression analysis are (1) diversity of facial features in an effort to increase the number of expressions that may be recognized; (2) recognition of facial action units and their combinations rather than more global and easily identified emotion-specified expressions; (3) more robust systems for face acquisition, facial data extraction and representation, and facial expression recognition to handle head motion (both in-plane and out-of-plane), occlusion, lighting change, and low intensity expressions, all of which are common in spontaneous facial behavior in naturalistic environments; (4) fully automatic and real-time AFEA

systems; and (5) combination of facial actions with other modes such as gesture, prosody, and speech. All of these developments move AFEA toward real-life applications. Several databases that addresses most problems for deliberate facial expression analysis have been released to researchers to conduct comparative tests of their methods. Databases with ground-truth labels, preferably both action units and emotion-specified expressions, are needed for the next generation of systems, which are intended for naturally occurring behavior (spontaneous and multimodal) in real-life settings. Work in spontaneous facial expression analysis is just now emerging and potentially will have significant impact across a range of theoretical and applied topics.

**Acknowledgements** We sincerely thank Zhen Wen and Hatice Gunes for providing pictures and their permission to use them in this chapter.

## References

1. Ahlberg, J., Forchheimer, R.: Face tracking for model-based coding and face animation. *Int. J. Imaging Syst. Technol.* **13**(1), 8–22 (2003)
2. Baker, S., Kanade, T.: Limits on super-resolution and how to break them. *IEEE Trans. Pattern Anal. Mach. Intell.* **24**(9), 1167–1183 (2002)
3. Bartlett, M., Hager, J., Ekman, P., Sejnowski, T.: Measuring facial expressions by computer image analysis. *Psychophysiology* **36**, 253–264 (1999)
4. Bartlett, M., Braathen, B., Littlewort-Ford, G., Hershey, J., Fasel, I., Marks, T., Smith, E., Sejnowski, T., Movellan, J.R.: Automatic analysis of spontaneous facial behavior: A final project report. Technical Report INC-MPLab-TR-2001.08, Machine Perception Lab, Institute for Neural Computation, University of California, San Diego (2001)
5. Bartlett, M., Littlewort, G., Frank, M., Lainscsek, C., Fasel, I., Movellan, J.: Automatic recognition of facial actions in spontaneous expressions. *J. Multimed.* **1**(6), 22–35 (2006)
6. Black, M.: Robust incremental optical flow. PhD thesis, Yale University (1992)
7. Black, M., Yacoob, Y.: Tracking and recognizing rigid and non-rigid facial motions using local parametric models of image motion. In: *Proc. of International Conference on Computer Vision*, pp. 374–381 (1995)
8. Black, M., Yacoob, Y.: Recognizing facial expressions in image sequences using local parameterized models of image motion. *Int. J. Comput. Vis.* **25**(1), 23–48 (1997)
9. Brown, L., Tian, Y.-L.: Comparative study of coarse head pose estimation. In: *IEEE Workshop on Motion and Video Computing*, Orlando (2002)
10. Camras, L., Lambrecht, L., Michel, G.: Infant surprise expressions as coordinative motor structures. *J. Nonverbal Behav.* **20**, 183–195 (1966)
11. Carroll, J., Russell, J.: Do facial expression signal specific emotions? *J. Pers. Soc. Psychol.* **70**, 205–218 (1996)
12. Carroll, J., Russell, J.: Facial expression in Hollywood's portrayal of emotion. *J. Pers. Soc. Psychol.* **72**, 164–176 (1997)
13. Chang, Y., Hu, C., Feris, R., Turk, M.: Manifold based analysis of facial expression. *Image Vis. Comput.* **24**(6), 605–614 (2006)
14. Chen, L.: Joint processing of audio-visual information for the recognition of emotional expressions in human-computer interaction. PhD thesis, University of Illinois at Urbana-Champaign, Department of Electrical Engineering (2000)
15. Cohen, I., Sebe, N., Cozman, F., Cirelo, M., Huang, T.: Coding, analysis, interpretation, and recognition of facial expressions. *J. Comput. Vis. Image Underst.* (2003). Special Issue on Face Recognition

16. Cohn, J., Katz, G.: Bimodal expression of emotion by face and voice. In: ACM and ATR Workshop on Face/Gesture Recognition and Their Applications, pp. 41–44 (1998)
17. Cohn, J., Zlochow, A., Lien, J., Kanade, T.: Automated face analysis by feature point tracking has high concurrent validity with manual facs coding. *Psychophysiology* **36**, 35–43 (1999)
18. Cohn, J., Kanade, T., Moriyama, T., Ambadar, Z., Xiao, J., Gao, J., Imamura, H.: A comparative study of alternative facs coding algorithms. Technical Report CMU-RI-TR-02-06, Robotics Institute, Carnegie Mellon University, Pittsburgh, November 2001
19. Cohn, J., Schmidt, K., Gross, R., Ekman, P.: Individual differences in facial expression: stability over time, relation to self-reported emotion, and ability to inform person identification. In: *Proceedings of the International Conference on Multimodal User Interfaces (ICMI 2002)*, pp. 491–496 (2002)
20. Cohn, J., Kreuz, T., Yang, Y., Nguyen, M., Padilla, M., Zhou, F., Fernando, D.: Detecting depression from facial actions and vocal prosody. In: *International Conference on Affective Computing and Intelligent Interaction (ACII2009)* (2009)
21. Darwin, C.: *The Expression of Emotions in Man and Animals*. Murray, London (1872), reprinted by University of Chicago Press, 1965
22. Daugmen, J.: Complete discrete 2d Gabor transforms by neural networks for image analysis and compression. *IEEE Trans. Acoust. Speech Signal Process.* **36**(7), 1169–1179 (1988)
23. Donato, G., Bartlett, M., Hager, J., Ekman, P., Sejnowski, T.: Classifying facial actions. *IEEE Trans. Pattern Anal. Mach. Intell.* **21**(10), 974–989 (1999)
24. Douglas-Cowie, E., Cowie, R., Schroeder, M.: The description of naturally occurring emotional speech. In: *International Conference of Phonetic Sciences* (2003)
25. Eibl-Eibesfeldt, I.: *Human Ethology*. Aldine de Gruyter, New York (1989)
26. Ekman, P.: *The Argument and Evidence about Universals in Facial Expressions of Emotion*, vol. 58, pp. 143–164. Wiley, New York (1989)
27. Ekman, P.: Facial expression and emotion. *Am. Psychol.* **48**, 384–392 (1993)
28. Ekman, P., Friesen, W.: *Pictures of Facial Affect*. Consulting Psychologist, Palo Alto (1976)
29. Ekman, P., Friesen, W.: *The Facial Action Coding System: A Technique for the Measurement of Facial Movement*. Consulting Psychologists Press, San Francisco (1978)
30. Ekman, P., Rosenberg, E.E.: *What the Face Reveals*. Oxford University, New York (1997)
31. Ekman, P., Hager, J., Methvin, C., Irwin, W.: Ekman–Hager facial action exemplars. Human Interaction Laboratory, University of California, San Francisco
32. Essa, I., Pentland, A.: Coding, analysis, interpretation, and recognition of facial expressions. *IEEE Trans. Pattern Anal. Mach. Intell.* **19**(7), 757–763 (1997)
33. Farkas, L., Munro, I.: *Anthropometric Facial Proportions in Medicine*. Charles C Thomas, Springfield (1987)
34. Fasel, B., Luttin, J.: Recognition of asymmetric facial action unit activities and intensities. In: *Proceedings of International Conference of Pattern Recognition* (2000)
35. Fasel, B., Luttin, J.: Automatic facial expression analysis: Survey. *Pattern Recognit.* **36**(1), 259–275 (2003)
36. Fleiss, J.: *Statistical Methods for Rates and Proportions*. Wiley, New York (1981)
37. Ford, G.: Fully automatic coding of basic expressions from video. Technical Report INC-MPLab-TR-2002.03, Machine Perception Lab, Institute for Neural Computation, University of California, San Diego (2002)
38. Fox, N., Reilly, R.: Audio-visual speaker identification. In: *Proc. of the 4th International Conference on Audio- and Video-Based Biometric Person Authentication* (2003)
39. Fox, N., Gross, R., de Chazal, P., Cohn, J., Reilly, R.: Person identification using multi-modal features: speech, lip, and face. In: *Proc. of ACM Multimedia Workshop in Biometrics Methods and Applications (WBMA 2003)*, CA (2003)
40. Frank, M., Ekman, P.: The ability to detect deceit generalizes across different types of high-stake lies. *Pers. Soc. Psychol.* **72**, 1429–1439 (1997)
41. Friesen, W., Ekman, P.: *Emfacs-7: emotional facial action coding system*. Unpublished manuscript, University of California at San Francisco (1983)

42. Fujiyoshi, H., Lipton, A.: Real-time human motion analysis by image skeletonization. In: *Proc. of the Workshop on Application of Computer Vision* (1998)
43. Fukui, K., Yamaguchi, O.: Facial feature point extraction method based on combination of shape extraction and pattern matching. *Syst. Comput. Jpn.* **29**(6), 49–58 (1998)
44. Gunes, H., Piccardi, M.: A bimodal face and body gesture database for automatic analysis of human nonverbal affective behavior. In: *International Conference on Pattern Recognition (ICPR)*, pp. 1148–1153 (2006)
45. Gunes, H., Piccardi, M.: Automatic temporal segment detection and affect recognition from face and body display. *IEEE Trans. Syst. Man Cybern., Part B, Cybern.* **39**(1), 64–84 (2009)
46. Hampapur, A., Pankanti, S., Senior, A., Tian, Y., Brown, L., Bolle, R.: Face cataloger: multi-scale imaging for relating identity to location. In: *Proceedings of IEEE Conference on Advanced Video and Signal Based Surveillance* (2003)
47. Heisele, B., Serre, T., Pontil, M., Poggio, T.: Component-based face detection. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recogn. (CVPR)* (2001)
48. Izard, C., Dougherty, L., Hembree, E.A.: A system for identifying affect expressions by holistic judgments. Unpublished Manuscript, University of Delaware (1983)
49. Kanade, T., Cohn, J., Tian, Y.-L.: Comprehensive database for facial expression analysis. In: *Proceedings of International Conference on Face and Gesture Recognition*, pp. 46–53 (2000)
50. Kimura, S., Yachida, M.: Facial expression recognition and its degree estimation. In: *Proc. of the International Conference on Computer Vision and Pattern Recognition*, pp. 295–300 (1997)
51. Kleck, R., Mendolia, M.: Decoding of profile versus full-face expressions of affect. *J. Non-verbal Behav.* **14**(1), 35–49 (1990)
52. Kobayashi, H., Tange, K., Hara, F.: Dynamic recognition of six basic facial expressions by discrete-time recurrent neural network. In: *Proc. of the International Joint Conference on Neural Networks*, pp. 155–158 (1993)
53. Kobayashi, H., Tange, K., Hara, F.: Real-time recognition of six basic facial expressions. In: *Proc. IEEE Workshop on Robot and Human Communication*, pp. 179–186 (1995)
54. Kraut, R., Johnson, R.: Social and emotional messages of smiling: an ethological approach. *J. Pers. Soc. Psychol.* **37**, 1539–1523 (1979)
55. Li, S., Gu, L.: Real-time multi-view face detection, tracking, pose estimation, alignment, and recognition. In: *IEEE Conf. on Computer Vision and Pattern Recognition Demo Summary* (2001)
56. Lien, J.-J., Kanade, T., Cohn, J., Li, C.: Subtly different facial expression recognition and expression intensity estimation. In: *Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition*, pp. 853–859 (1998)
57. Lien, J.-J., Kanade, T., Cohn, J., Li, C.: Detection, tracking, and classification of action units in facial expression. *J. Robot. Auton. Syst.* **31**, 131–146 (2000)
58. Lucey, S., Wang, Y., Cox, M., Sridharan, S., Cohn, J.: Efficient constrained local model fitting for non-rigid face alignment. *Image Vis. Comput.* **27**(12), 1804–1813 (2009)
59. Lyons, M., Akamasku, S., Kamachi, M., Gyoba, J.: Coding facial expressions with Gabor wavelets. In: *Proceedings of International Conference on Face and Gesture Recognition* (1998)
60. Mahoor, M., Cadavid, S., Messinger, D., Cohn, J.: A framework for automated measurement of the intensity of non-posed facial action units. In: *IEEE Workshop on CVPR for Human Communicative Behavior Analysis*, pp. 74–80 (2009)
61. Manstead, A.: Expressiveness as an Individual Difference, pp. 285–328. Cambridge University Press, Cambridge (1991)
62. Martin, P., Bateson, P.: *Measuring Behavior: An Introductory Guide*. Cambridge University Press, Cambridge (1986)
63. Martinez, A., Benavente, R.: The ar face database. CVC Technical Report, No. 24 (1998)
64. Mase, K.: Recognition of facial expression from optical flow. *IEICE Trans. Electron.* **74**(10), 3474–3483 (1991)
65. Matias, R., Cohn, J., Ross, S.: A comparison of two systems to code infants' affective expression. *Dev. Psychol.* **25**, 483–489 (1989)

66. [Matthews, I., Baker, S.: Active appearance models revisited. \*Int. J. Comput. Vis.\* \*\*60\*\*\(2\), 135–164 \(2004\)](#)
67. [Moriyama, T., Kanade, T., Cohn, J., Xiao, J., Ambadar, Z., Gao, J., Imanura, M.: Automatic recognition of eye blinking in spontaneously occurring behavior. In: \*Proceedings of the 16th International Conference on Pattern Recognition \(ICPR '2002\)\*, vol. 4, pp. 78–81 \(2002\)](#)
68. [Moses, Y., Reynard, D., Blake, A.: Determining facial expressions in real time. In: \*Proc. of Int. Conf. On Automatic Face and Gesture Recognition\*, pp. 332–337 \(1995\)](#)
69. [Pantic, M., Rothkrantz, L.: Automatic analysis of facial expressions: the state of the art. \*IEEE Trans. Pattern Anal. Mach. Intell.\* \*\*22\*\*\(12\), 1424–1445 \(2000\)](#)
70. [Pantic, M., Rothkrantz, L.: Expert system for automatic analysis of facial expression. \*Image Vis. Comput.\* \*\*18\*\*\(11\), 881–905 \(2000\)](#)
71. [Pantic, M., Valstar, M., Rademaker, R., Maat, L.: Web-based database for facial expression analysis. In: \*International conference on Multimedia and Expo \(ICME05\)\* \(2005\)](#)
72. [Pentland, A., Moghaddam, B., Starner, T.: View-based and modular eigenspaces for face recognition. In: \*Proc. IEEE Conf. Computer Vision and Pattern Recognition\*, pp. 84–91 \(1994\)](#)
73. [Pighin, F., Szeliski, H., Salesin, D.: Synthesizing realistic facial expressions from photographs. In: \*Proc of SIGGRAPH\* \(1998\)](#)
74. [Pilz, K., Thornton, I., Bülthoff, H.: A search advantage for faces learned in motion. In: \*Experimental Brain Research\* \(2006\)](#)
75. [Rinn, W.: The neuropsychology of facial expression: a review of the neurological and psychological mechanisms for producing facial expressions. \*Psychol. Bull.\* \*\*95\*\*, 52–77 \(1984\)](#)
76. [Rizvi, S., Phillips, P., Moon, H.: The Feret verification testing protocol for face recognition algorithms. In: \*Proceedings of the Third International Conference on Automatic Face and Gesture Recognition\*, pp. 48–55 \(1998\)](#)
77. [Rosenblum, M., Yacoob, Y., Davis, L.: Human expression recognition from motion using a radial basis function network architecture. \*IEEE Trans. Neural Netw.\* \*\*7\*\*\(5\), 1121–1138 \(1996\)](#)
78. [Rowley, H., Baluja, S., Kanade, T.: Neural network-based face detection. \*IEEE Trans. Pattern Anal. Mach. Intell.\* \*\*20\*\*\(1\), 23–38 \(1998\)](#)
79. [Russell, J.: Culture and the categorization. \*Psychol. Bull.\* \*\*110\*\*, 426–450 \(1991\)](#)
80. [Russell, J.: Is there universal recognition of emotion from facial expression? A review of the cross-cultural studies. \*Psychol. Bull.\* \*\*115\*\*, 102–141 \(1991\)](#)
81. [Samal, A., Iyengar, P.: Automatic recognition and analysis of human faces and facial expressions: A survey. \*Pattern Recognit.\* \*\*25\*\*\(1\), 65–77 \(1992\)](#)
82. [Sayette, M., Cohn, J., Wertz, J., Perrott, M., Parrott, D.: A psychometric evaluation of the facial action coding system for assessing spontaneous expression. \*J. Nonverbal Behav.\* \*\*25\*\*, 167–186 \(2001\)](#)
83. [Scherer, K., Ekman, P.: \*Handbook of Methods in Nonverbal Behavior Research\*. Cambridge University Press, Cambridge \(1982\)](#)
84. [Scherer, K., Ellgring, H.: Multimodal expression of emotion: Affect programs or componential appraisal patterns? \*Emotion\* \*\*7\*\*\(1\), 158–171 \(2007\)](#)
85. [Schmidt, K., Cohn, J.F., Tian, Y.-L.: Signal characteristics of spontaneous facial expressions: Automatic movement in solitary and social smiles. \*Biol. Psychol.\* \(2003\)](#)
86. [Schneiderman, H., Kanade, T.: A statistical model for 3d object detection applied to faces and cars. In: \*IEEE Conference on Computer Vision and Pattern Recognition\*. IEEE, New York \(2000\)](#)
87. [Shan, C., Gong, S., McOwan, P.: Facial expression recognition based on local binary patterns: A comprehensive study. \*Image Vis. Comput.\* \*\*27\*\*\(6\), 803–816 \(2009\)](#)
88. [Sim, T., Baker, S., Bsat, M.: The cmu pose, illumination, and expression database. \*IEEE Trans. Pattern Anal. Mach. Intell.\* \*\*25\*\*\(12\), 1615–1618 \(2003\)](#)
89. [Sung, K., Poggio, T.: Example-based learning for view-based human face detection. \*IEEE Trans. Pattern Anal. Mach. Intell.\* \*\*20\*\*\(1\), 39–51 \(1998\)](#)



90. Suwa, M., Sugie, N., Fujimora, K.: A preliminary note on pattern recognition of human emotional expression. In: *International Joint Conference on Pattern Recognition*, pp. 408–410 (1978)
91. Tao, H., Huang, T.: Explanation-based facial motion tracking using a piecewise Bezier volume deformation model. In: *Proc. IEEE Conf. Computer Vision and Pattern Recognition (1999)*
92. Terzopoulos, D., Waters, K.: Analysis of facial images using physical and anatomical models. In: *IEEE International Conference on Computer Vision*, pp. 727–732 (1990)
93. Tian, Y.-L., Bolle, R.: Automatic detecting neutral face for face authentication. In: *Proceedings of AAAI-03 Spring Symposium on Intelligent Multimedia Knowledge Management, CA (2003)*
94. Tian, Y.-L., Kanade, T., Cohn, J.: Eye-state action unit detection by Gabor wavelets. In: *Proceedings of International Conference on Multi-modal Interfaces (ICMI 2000)*, pp. 143–150, September 2000
95. Tian, Y.-L., Kanade, T., Cohn, J.: Recognizing action units for facial expression analysis. *IEEE Trans. Pattern Anal. Mach. Intell.* **23**(2), 1–19 (2001)
96. Tian, Y.-L., Kanade, T., Cohn, J.: Evaluation of Gabor-wavelet-based facial action unit recognition in image sequences of increasing complexity. In: *Proceedings of the 5th IEEE International Conference on Automatic Face and Gesture Recognition (FG'02)*, Washington, DC (2002)
97. Tian, Y.-L., Brown, L., Hampapur, A., Pankanti, S., Senior, A., Bolle, R.: Real world real-time automatic recognition of facial expressions. In: *Proceedings of IEEE Workshop on Performance Evaluation of Tracking and Surveillance*, Graz, Austria (2003)
98. Toyama, K.: “Look, ma—no hands!” hands-free cursor control with real-time 3d face tracking. In: *Proc. Workshop on Perceptual User Interfaces (PUI'98)* (1998)
99. VanSwearingen, J., Cohn, J., Bajaj-Luthra, A.: Specific impairment of smiling increases severity of depressive symptoms in patients with facial neuromuscular disorders. *J. Aesthet. Plast. Surg.* **23**, 416–423 (1999)
100. Vetter, T.: Learning novel views to a single face image. In: *Proceedings of the IEEE International Conference on Automatic Face and Gesture Recognition*, pp. 22–29 (1995)
101. Viola, P., Jones, M.: Robust real-time object detection. In: *International Workshop on Statistical and Computational Theories of Vision—Modeling, Learning, Computing, and Sampling (2001)*
102. Wen, Z., Huang, T.: Capturing subtle facial motions in 3d face tracking. In: *Proc. of Int. Conf. on Computer Vision (2003)*
103. Wu, Y., Toyama, K.: Wide-range person and illumination-insensitive head orientation estimation. In: *Proceedings of International Conference on Automatic Face and Gesture Recognition*, pp. 183–188 (2000)
104. Xiao, J., Moriyama, T., Kanade, T., Cohn, J.: Robust full-motion recovery of head by dynamic templates and re-registration techniques. *Int. J. Imaging Syst. Technol.* (2003)
105. Xiao, J., Baker, S., Matthews, I., Kanade, T.: Real-time combined 2d + 3d active appearance models. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 535–542 (2004)
106. Yacoob, Y., Black, M.: Parameterized modeling and recognition of activities. In: *Proc. 6th IEEE Int. Conf. on Computer Vision*, pp. 120–127, Bombay (1998)
107. Yacoob, Y., Davis, L.: Recognizing human facial expression from long image sequences using optical flow. *IEEE Trans. Pattern Anal. Mach. Intell.* **18**(6), 636–642 (1996)
108. Yacoob, Y., Lam, H.-M., Davis, L.: Recognizing faces showing expressions. In: *Proc. Int. Workshop on Automatic Face- and Gesture-Recognition*, pp. 278–283, Zurich, Switzerland (1995)
109. Yang, M., Kriegman, D., Ahuja, N.: Detecting faces in images: a survey. *IEEE Trans. Pattern Anal. Mach. Intell.* **24**(1) (2002)
110. Yang, P., Liu, Q., Metaxas, D.: Facial expression recognition using encoded dynamic features. In: *International conference on Computer Vision and Pattern Recognition (CVPR) (2008)*



111. [Yang, P., Liu, Q., Cui, X., Metaxas, D.: Rankboost with l1 regularization for facial expression recognition and intensity estimation. In: International conference on Computer Vision \(ICCV\) \(2009\)](#)
112. [Yin, L., Wei, X., Sun, Y., Wang, J., Rosato, M.: A 3d facial expression database for facial behavior research. In: International Conference on Automatic Face and Gesture Recognition, pp. 211–216 \(2006\)](#)
113. [Yin, L., Chen, X., Sun, Y., Worm, T., Reale, M.: A high-resolution 3d dynamic facial expression database. In: International Conference on Automatic Face and Gesture Recognition \(2008\)](#)
114. [Zeng, Z., Pantic, G.R.M., Huang, T.: A survey of affect recognition methods: Audio, visual, and spontaneous expressions. IEEE Trans. Pattern Anal. Mach. Intell. \*\*31\*\*\(1\), 39–58 \(2009\)](#)
115. [Zhang, Y., Ji, Q.: Facial expression recognition with dynamic Bayesian networks. IEEE Trans. Pattern Anal. Mach. Intell. \*\*27\*\*\(5\) \(2005\)](#)
116. [Zhang, Z., Lyons, M., Schuster, M., Akamatsu, S.: Comparison between geometry-based and Gabor-wavelets-based facial expression recognition using multi-layer perceptron. In: International Workshop on Automatic Face and Gesture Recognition, pp. 454–459 \(1998\)](#)
117. [Zhang, Y., Ji, Q., Zhu, Z., Yi, B.: Dynamic facial expression analysis and synthesis with mpeg-4 facial animation parameters. IEEE Trans. Circuits Syst. Video Technol. \*\*18\*\*\(10\), 1383–1396 \(2008\)](#)
118. [Zhao, L., Pingali, G., Carlbom, I.: Real-time head orientation estimation using neural networks. In: Proc of the 6th International Conference on Image Processing \(2002\)](#)
119. [Zlochower, A., Cohn, J., Lien, J., Kanade, T.: A computer vision based method of facial expression analysis in parent-infant interaction. In: International Conference on Infant Studies, Atlanta \(1998\)](#)