

A Survey of Affect Recognition Methods: Audio, Visual, and Spontaneous Expressions

Zhihong Zeng, *Member, IEEE*, Maja Pantic, *Senior Member, IEEE*, Glenn I. Roisman and Thomas S. Huang, *Fellow, IEEE*

Abstract— Automated analysis of human affective behavior has attracted increasing attention from researchers in psychology, computer science, linguistics, neuroscience, and related disciplines. However, the existing methods typically handle only deliberately displayed and exaggerated expressions of prototypical emotions despite the fact that deliberate behaviour differs in visual appearance, audio profile, and timing from spontaneously occurring behaviour. To address this problem, efforts to develop algorithms that can process naturally occurring human affective behaviour have recently emerged. Moreover, an increasing number of efforts are reported toward multimodal fusion for human affect analysis including audiovisual fusion, linguistic and paralinguistic fusion, and multi-cue visual fusion based on facial expressions, head movements, and body gestures. This paper introduces and surveys these recent advances. We first discuss human emotion perception from a psychological perspective. Next we examine available approaches to solving the problem of machine understanding of human affective behavior, and discuss important issues like the collection and availability of training and test data. We finally outline some of the scientific and engineering challenges to advancing human affect sensing technology.

Index Terms— Evaluation/methodology, human-centered computing, introductory and survey.

1 INTRODUCTION

A widely accepted prediction is that computing will move to the background, weaving itself into the fabric of our everyday living spaces and projecting the human user into the foreground. Consequently, the future “ubiquitous computing” environments will need to have human-centered designs instead of computer-centered designs [26], [31], [100], [107], [109]. Current human-computer interaction (HCI) designs, however, usually involve traditional interface devices such as the keyboard and mouse, and are constructed to emphasize the transmission of explicit messages while ignoring implicit information about the user, such as changes in affective state. Yet, a change in the user’s affective state is a fundamental component of human-human communication. Some affective states motivate human actions and others enrich the meaning of human communication. Consequently, the traditional HCI that ignores the user’s affective states filters out a large portion of the information available in the interaction process. As a result, such interactions are frequently perceived as cold, incompetent and socially inept. Human Computing paradigm suggests that user interfaces of the future need to be anticipatory and human-centered, built for humans, and based on naturally occurring multimodal human communication

[100], [109]. Specifically, human-centered interfaces must have the ability to detect subtleties of and changes in the user’s behavior, especially his or her affective behavior, and to initiate interactions based on this information, rather than simply responding to the user’s commands.

Examples of affect-sensitive, multimodal HCI systems include the system of Lisetti and Nasoz [85], which combines facial expression and physiological signals to recognize the user’s emotion like fear and anger and then to adapt an animated interface agent to mirror the user’s emotion, the multimodal system of Duric et al. [39], which applies a model of embodied cognition that can be seen as a detailed mapping between the user’s affective states and the types of interface adaptations, the proactive HCI tool of Maat and Pantic [89] capable of learning and analyzing the user’s context-dependent behavioral patterns from multi-sensory data and of adapting the interaction accordingly, the automated Learning Companion of Kapoor et al. [72] that combines information from cameras, a sensing chair and mouse, wireless skin sensor, and task state to detect frustration in order to predict when the user need help, and the multimodal computer-aided learning system¹ at Beckman Institute UIUC where the computer avatar offers an appropriate tutoring strategy based on the information of user’s facial expression, keywords, eye movement and task state. These systems represent initial efforts towards the future human-centered, multimodal HCI.

Except in standard HCI scenarios, potential commercial applications of automatic human affect recognition include affect-sensitive systems for customer services, call centers, intelligent automobile system, and game and entertainment industry. These systems will change the ways

- Zhihong Zeng and Thomas S. Huang are with Beckman Institute, University of Illinois at Urbana-Champaign. 405 N Mathews Av., Urbana, 61801. E-mail: {zhzeng.huang}@ifp.uiuc.edu..
- Maja Pantic is with Imperial College London, Department of Computing, 180 Queen’s Gate, London SW7 2AZ, UK, and with University of Twente, Faculty of EEMCS, the Netherlands. E-mail: m.pantic@imperial.ac.uk.
- Glenn I. Roisman is with Psychology Department, University of Illinois at Urbana-Champaign. 603 East Daniel St., Champaign, IL 61820. E-mail: roisman@uiuc.edu.

in which we interact with computer systems. For example, an automatic service call center with an affect detector would be able to make appropriate response or pass control over to human operators [83], and an intelligent automobile system with a fatigue detector could monitor the vigilance of the driver and apply appropriate action to avoid accidents [69].

Another important application of automated systems for human affect recognition is in affect-related research (e.g. in psychology, psychiatry, behavioral and neuroscience), where such systems can improve the quality of the research by improving the reliability of measurements and speeding up the currently tedious, manual task of processing data on human affective behavior [47]. The research areas that would reap substantial benefits from such automatic tools include social and emotional development research [111], mother-infant interaction [29], tutoring [54], psychiatric disorders [45], and studies on affective expressions (e.g., deception) [65], [47]. Automated detectors of affective states and moods including fatigue, depression, and anxiety, could also form an important step toward personal wellness and assistive technologies [100].

Because of this practical importance and the theoretical interest of cognitive scientists, automatic human affect analysis has attracted the interest of many researchers in the past three decades. Suwa et al. [127] presented an early attempt in 1978 to automatically analyze facial expressions. The vocal emotion analysis has an even longer history, starting with the study of Williams and Stevens from 1972 [145]. Since late 90s, an increasing number of efforts toward automatic affect recognition were reported in the literature. Early efforts toward machine affect recognition from face images include those of Mase [90], and Kobayashi and Hara [76] from 1991. Early efforts toward machine analysis of basic emotions from vocal cues include studies like that of Dellaert et al. in 1996 [33]. The study of Chen et al. in 1998 [22] represents an early attempt toward audiovisual affect recognition. For exhaustive surveys of the past work in machine analysis of affective expressions, readers are referred to [115], [31], [102], [49], [96], [105], [130], [121], [98] that were published in 1992 to 2007 respectively.

Overall, most of the existing approaches to automatic human affect analysis are:

- trained and tested on deliberately displayed series of exaggerated expressions of affective behavior,
- aimed at recognition of a small number of prototypical (basic) expressions of emotion (i.e., happiness, sadness, anger, fear, surprise, and disgust),
- single-modal: information processed by the computer system is limited to either face images or the speech signals.

Accordingly, reviewing the efforts toward single-modal analysis of artificial affective expressions have been the focus in the previously published survey papers among which the papers of Cowie et al. in 2001 [31] and of Pantic and Rothkrantz in 2003 [102] have been the most comprehensive and widely cited in this field to date. At that time when these surveys were written, most of the

available datasets of affective displays were small, and contained only deliberate affective displays (mainly of the six prototypical emotions) recorded under highly constrained conditions. Multimedia data were rare, and there was no 3D data on facial affective behavior, no data of combined face and body displays of affective behavior, and it was rare to find data that included spontaneous displays of affective behavior.

Hence, while automatic detection of the six basic emotions in posed, controlled audio or visual displays can be done with reasonably high accuracy, detecting these expressions or any expression of human affective behavior in less constrained settings is still a very challenging problem due to the fact that deliberate behaviour differs in visual appearance, audio profile, and timing from spontaneously occurring behaviour. Due to this criticism received from both cognitive and computer scientists, the focus of the research in the field started to shift to automatic analysis of spontaneously displayed affective behavior. Several studies have recently emerged on machine analysis of spontaneous facial expressions (e.g., [10], [28], [135], [4]) and vocal expressions (e.g., [12], [83]).

Also, it has been shown by several experimental studies that integrating the information from audio and video leads to an improved performance of affective behavior recognition. The improved reliability of audiovisual approaches in comparison to single-modal approaches can be explained as follows. Current techniques for detection and tracking of facial expressions are sensitive to head pose, clutter, and variations in lighting conditions, while current techniques for speech processing are sensitive to auditory noise. Audiovisual fusion can make use of the complementary information from these two channels. In addition, many psychological studies have theoretically and empirically demonstrated the importance of integration of information from multiple modalities (vocal and visual expression in this paper) to yield a coherent representation and inference of emotions [1], [113], [117]. As a result, an increased number of studies on audiovisual human affect recognition have emerged in recent years (e.g., [17], [53], [151]).

This paper introduces and surveys these recent advances in the research on human affect recognition. In contrast to previously published survey papers in the field, it focuses on the approaches which can handle audio and/or visual recordings of *spontaneous* (as opposed to posed) displays of affective states. It also examines the state-of-the-art methods that have not been reviewed in previous survey papers, but are important specifically for advancing human affect sensing technology. Finally, we discuss the collection and availability of training and test data in detail. The paper is organized as follows. Section 2 describes human perception of affect from a psychological perspective. Section 3 provides a detailed review of the related studies, including multimedia emotion databases and existing human affect recognition methods. Section 4 discusses some of the challenges that researchers face in this field. A summary and closing remarks conclude the paper.

2 HUMAN AFFECT (EMOTION) PERCEPTION

Automatic affect recognition is inherently a multi-disciplinary enterprise involving different research fields, including psychology, linguistics, computer vision, speech analysis, and machine learning. There is no doubt that the progress in automatic affect recognition is contingent on the progress of the research in each of those fields [44].

2.1 The Description of Affect

We begin by briefly introducing three primary ways that affect has been conceptualized in psychological research. Research on the basic structure and description of affect is important in that these conceptualizations provide information about the affective displays that automatic emotion recognition systems are designed to detect.

Perhaps the most longstanding way that affect has been described by psychologists is in terms of discrete categories, an approach that is rooted in the language of daily life [40], [41], [131]. The most popular example of this description is the prototypical (basic) emotion categories, which include happiness, sadness, fear, anger, disgust, and surprise. This description of basic emotions was supported especially by the cross-cultural studies conducted by Ekman [40], [42] indicating that humans perceive certain basic emotions with respect to facial expression in the same way regardless of culture. This influence of basic emotion theory has resulted in the fact that most of existing studies of automatic affect recognition focus on recognizing these basic emotions. The main advantage of a category representation is that people use this categorical scheme to describe observed emotional displays in daily life. The labeling scheme based on category is very intuitive and thus matches people's experience. However, discrete lists of emotions fail to describe the range of emotions that occur in natural communication settings. For example, although prototypical emotions are key points of emotion reference, they cover a rather small part of our daily emotional displays. Selection of affect categories that can describe the wide variety of affective displays that people show in daily interpersonal interactions needs to be done in a pragmatic and context-dependent manner [102], [105].

An alternative to categorical description of human affect is the dimensional description [58], [114], [140], where an affective state is characterized in terms of a small number of latent dimensions, rather than in terms of a small number of discrete emotion categories. These dimensions include evaluation, activation, control, power, etc. In particular, the evaluation and activation dimensions are expected to reflect the main aspects of emotion. The evaluation dimension measures how human feels, from positive to negative. The activation dimension measures whether humans are more or less likely to take an action under the emotional state, from active to passive. In contrast to categorical representation, dimensional representation enables raters to label a range of emotions. However, the projection of the high-dimensional emotional states onto a rudimentary 2D

space results to some degree in the loss of information. Some emotions become indistinguishable (e.g., fear and anger) and some emotions lie outside the space (e.g., surprise). This representation is not intuitive and raters need special training to use the dimensional labeling system (e.g., Feeltrace system [30]). In automatic emotion recognition systems that are based on the 2D dimensional emotion representation (e.g., [17], [53]), the problem is often further simplified to 2-class (positive vs. negative and active vs. passive) or 4-class (quadrants of 2D space) classification.

One of the most influential emotion theories in modern psychology is the appraisal-based approach [117] that can be regarded as the extension of the dimensional approach described above. In this representation, emotion is described through a set of stimulus evaluation checks, including the novelty, intrinsic pleasantness, goal-based significance, coping potential, and compatibility with standards. However, translating this scheme into one engineering framework for the purposes of automatic emotion recognition remains challenging [116].

2.2 Association between Affect, Audio and Visual Signals

Affective arousal modulates all human communicative signals. Psychologists and linguists have various opinions about the importance of different cues (audio and visual cues in this paper) in human affect judgment. Ekman [41] found that the relative contributions of facial expression, speech, and body gestures to affect judgment depend both on the affective state and the environment where the affective behavior occurs while some studies (e.g., [1], [92]) indicated that a facial expression in the visual channel is the most important affective cue and correlates well with body as well as voice. Many studies have theoretically and empirically demonstrated the advantage of integration of multiple modalities (vocal and visual expression) in human affect perception over single modalities [1], [113], [117].

Different from the traditional message judgment in which the aim is to infer what underlies a displayed behavior, such as affect or personality, another major approach to human behavior measurement is the sign judgment [26]. The aim of sign judgment is to describe the appearance rather than meaning of the shown behavior, such as facial signal, body gesture or speech rate. While message judgment is focused on interpretation, sign judgment attempts to be objective description, leaving the inference about the conveyed message to high-level decision making. As indicated by Cohn [26], most commonly used sign judgment method used for manual labeling of facial behavior is the Facial Action Coding System (FACS) proposed by Ekman et al. [43]. FACS is a comprehensive and anatomically based system that is used to measure all visually discernible facial movements in terms of atomic facial actions called Action Units (AUs). As AUs are independent of interpretation, they can be used for any high-level decision making process including recognition of basic emotions according to Emotional FACS (EMFACS) rules², recognition of various affective states according to

FACS Affect Interpretation Database (FACSAID)² introduced by Ekman et al. [43], as well as for recognition of other complex psychological states such as depression [47] or pain [144]. AUs of the FACS are very suitable to be used in studies on human naturalistic facial behavior as the thousands of anatomically possible facial expressions (independently of their high-level interpretation) can be described as combinations of 27 basic AUs and a number of AU descriptors. It is not surprising, therefore, that an increasing number of studies on human spontaneous facial behavior are based on automatic AU recognition (e.g., [10], [27], [135], [87], [134]).

Speech is another important communicative modality in human-human interaction. Speech conveys affective information through explicit (linguistic) messages, and implicit (paralinguistic) messages that reflect the way the words are spoken. As the linguistic content is concerned, some information about the speaker's affective state can be inferred directly from the surface features of words, which were summarized in some affective word dictionaries and lexical affinity [110], [142], and the rest of affective information lies below the text surface and can only be detected when the semantic context (e.g., discourse information) is taken into account. However, findings in basic research [1], [55] indicate that linguistic messages are rather unreliable means to analyze human (affective) behavior, and it is very difficult to anticipate a person's word choice and the associated intent in affective expressions. In addition, the association between linguistic content and emotion is language-dependent and generalizing from one language to another is very difficult to achieve.

When it comes to implicit, paralinguistic messages that convey affective information, basic researchers have not identified an optimal set of voice cues that reliably discriminate among emotions. Nonetheless, listeners seem to be accurate in decoding some basic emotions from prosody [70] and some non-basic affective states such as distress, anxiety, boredom, and sexual interest from nonlinguistic vocalizations like laughs, cries, sighs, and yawns [113]. Cowie et al. [31] provided a comprehensive summary of qualitative acoustic correlations for prototypical emotions.

In a summary, a large number of studies in psychology and linguistics confirm the correlation between some affective displays (especially prototypical emotions) and specific audio and visual signals (e.g., [1], [47], [113]). The human judgment agreement is typically higher for facial expression modality than it is for vocal expression modality. However, the amount of the agreement drops considerably when the stimuli are spontaneously displayed expressions of affective behavior rather than posed exaggerated displays. In addition, facial expression and vocal expression of emotion are often studied separately. This precludes finding evidence of the temporal correlation between them. On the other hand, a growing body of research in cognitive sciences argues that the dynamics of human behavior are crucial for its interpretation (e.g., [47], [113], [116], [117]). For example, it has been shown

that temporal dynamics of facial behavior represents a critical factor for distinction between spontaneous and posed facial behavior (e.g., [28], [47], [135], [134]) as well as for categorization of complex behaviors like pain, shame, and amusement (e.g., [47], [144], [4], [87]). Based on these findings, we may expect that temporal dynamics of each modality separately (facial and vocal) and temporal correlations between the two modalities play an important role in interpretation of human naturalistic, audiovisual affective behavior. However, these are virtually unexplored areas of research.

Another largely unexplored area of research is that of context dependency. The interpretation of human behavioral signals is context dependent. For example a smile can be a display of politeness, irony, joy, or greeting. To interpret a behavioral signal, it is important to know the context in which this signal has been displayed – where the expresser is (e.g., inside, on the street, in the car), what his or her current task is, who the receiver is, and who the expresser is [113].

3 THE STATE OF THE ART

Rather than providing exhaustive coverage of all past efforts in the field of automatic recognition of human affect, we focus here on the efforts recently proposed in the literature that have not been reviewed elsewhere, that represent multimodal approaches to the problem of human affect recognition, that address the problem of automatic analysis of spontaneous affective behavior, or that represent exemplary approaches to treating a specific problem relevant for achieving a better human affect sensing technology. Due to limitation of our knowledge and page, we sincerely apologize to those authors whose work is not included in this paper.

For exhaustive surveys of the past efforts in the field, readers are referred to the following articles:

- overviews of early work on facial expression analysis: Samal and Iyengar 1992 [115], Pantic and Rothkrantz 2000 [101], and Fasel and Luttin 2003, [49],
- surveys of techniques for automatic facial muscle action recognition and facial expression analysis: Tian et al. 2005 [130], and Pantic and Bartlett 2007 [98], and
- overviews of multimodal affect recognition methods: Cowie et al. 2001 [31], Pantic and Rothkrantz 2003 [102], Pantic et al. 2005 [105], Sebe et al. 2005 [121], Jaimes and Sebe 2005 [68], and Zeng et al. 2007 [152] (this is a short, preliminary version of the survey presented in the current paper).

In this section we first offer an overview of the existing databases of audio and/or visual recordings of human affective displays, which provide the basis of automatic affect analysis. Next we examine available computing methods for automatic human affect recognition.

3.1 Databases

Having enough labeled data of human affective expressions is a prerequisite in designing automatic affect recognizer. Authentic affective expressions are difficult to collect because they are relatively rare and short lived,

² <http://face-and-emotion.com/dataface/general/homepage.jsp>

and filled with subtle context-based changes that make it difficult to elicit affective displays without influencing the results. In addition, manual labeling of spontaneous emotional expressions for ground truth is very time consuming, error prone, and expensive. This state of affairs makes automatic analysis of spontaneous emotional expression a very difficult task. Due to these difficulties, most of the existing studies on automatic analysis of human affective displays have been based on the “artificial” material of deliberately expressed emotions, elicited by asking the subjects to perform a series of emotional expressions in front of a camera and/or microphone.

However, increasing evidence suggests that deliberate behaviour differs in visual appearance, audio profile, and timing from spontaneously occurring behaviour. For example, Whissell shows that the posed nature of emotions in spoken language may differ in the choice of words and timing from corresponding performances in natural settings [142]. When it comes to facial behavior, there is a large body of research in psychology and neuroscience demonstrating that spontaneous and deliberately displayed facial behavior has differences both in utilized facial muscles and their dynamics (e.g., [47]). For instance, many types of spontaneous smiles (e.g., polite) are smaller in amplitude, longer in total duration, and slower in onset and offset time than posed smiles (e.g., [28], [47], [134]). Similarly, it has been shown that spontaneous brow actions (AU1, AU2 and AU4 in the FACS system) have different morphological and temporal characteristics (intensity, duration, and occurrence order) than posed brow actions [135]. It is not surprising, therefore, that methods of automated human affect analysis that have been trained on deliberate and often exaggerated behaviours usually fail to generalize to the subtlety and complexity of spontaneous affective behaviour.

In addition, most of the current human affect recognizers are evaluated using clear and constrained input (e.g., high quality visual and audio recording, non-occluded, and front-view or profile-view face), which is different from the input coming from a natural setting. In addition, most of the emotion expressions that occur in a realistic interpersonal or human-computer interaction are non-basic emotions [32]. Yet, the majority of the existing systems for human affect recognition aim at classifying the input expression as the basic emotion category (e.g., [31], [102], [105]).

These findings and the general lack of a comprehensive, reference set of audio and/or visual recordings of human affective displays motivated several efforts aimed at development of datasets that could be used for training and test of automatic systems for human affect analysis. Table 1 lists some noteworthy audio, visual, and audiovisual data resources that were reported in the literature. For each database, we provide the following information: affect elicitation method (i.e., whether the elicited affective displays are posed or spontaneous), size (the number of subjects and available data samples), modality (audio and/or visual), affect description (category or dimension), labeling scheme, and public accessibility. For other surveys of existing databases of human affective behav-

ior, the readers are referred to [32], [59], [106].

As far as the databases of deliberate affective behavior are concerned, the following databases need to be mentioned. The Cohn-Kanade facial expression database [71] is the most widely used database for facial expression recognition. The BU-3DFE database of Yin and colleagues [148] contains 3D range data of six prototypical facial expressions displayed at four different levels of intensity. The FABO database of Gunes and Piccardi [63] contains videos of facial expressions and body gestures portraying posed displays of basic and non-basic affective states (six prototypical emotions, uncertainty, anxiety, boredom, and neutral). The MMI facial expression database [106], [98] is to our knowledge the most comprehensive dataset of facial behavior recordings to date. It contains both posed expressions and spontaneous expressions of facial behavior. The available recordings of deliberate facial behavior are both static images and videos, where a large part of video recordings were recorded in both the frontal and the profile view of the face. The database represents a facial behavior data repository that is available, searchable, and downloadable via the Internet³. Although there are many databases of acted emotional speech⁴, a large majority of these datasets contain unlabeled data, which makes them unsuitable for research on automatic vocal affect recognition. The Banse-Scherer vocal affect database [8] and the Danish Emotional Speech database⁵ are the two most widely used databases in the research on vocal affect recognition from acted emotional speech. Finally, the Chen-Huang audiovisual database [21] is to our knowledge the largest multimedia databases containing facial and vocal deliberate displays of basic emotions and 4 cognitive states (interest, puzzlement, frustration and boredom).

The existing datasets of spontaneous affective behavior were collected in one of the following scenarios: human-human conversation, human-computer interaction, and use of a video kiosk. Human-human conversation scenarios include face-to-face interviews (e.g., [10], [38], [111], [65]), phone conversations (e.g., [34]), and meetings (e.g., [15], AMI⁶). Human computer interaction scenarios include Wizard of OZ scenarios (e.g., [13], SAL⁷), and computer-based dialogue systems (e.g., [83], [86]). In the video kiosk settings (e.g., [95], [98], [123]), the subjects’ affective reactions are recorded while the subjects are watching emotion-inducing videos.

In most of the existing databases discrete emotion categories are used as the emotion descriptors. The labels of prototypical emotions are often used, especially in the databases of deliberate affective behavior. In databases of spontaneous affective behavior, coarse affective states like positive vs. negative (e.g., [15], [83]), dimensional descriptions in the evaluation-activation space (e.g., SAL⁷), and some application-dependent affective states are usually

³ <http://www.mmifacedb.com/>

⁴ <http://emotion-research.net/wiki/Databases>

⁵ <http://cpk.auc.dk/~tb/speech/Emotions/>

⁶ <http://corpus.amiproject.org/>

⁷ <http://emotion-research.net/toolbox/toolboxdatabase.2006-09-26.5667892524>

used as the data labels. Interest, boredom, confusion, frustration, fatigue, empathy, stress, irony, annoyance, amusement, helplessness, panic, shame, reprehension, and rebelliousness are some typical examples of the used application-dependent affect-interpretative labels (e.g., [95], [63], [13], [111]).

TABLE 1. AUDIO AND/OR VISUAL DATABASES OF HUMAN AFFECTIVE BEHAVIOR
Legend: A – audio, V – video, AV- audiovisual, N/A – not available, Y – yes, N – not yet

References	Elicitation method	Size	A/V	Emotion description	Labeling	Accessibility
Cohn-Kanade (CK) '00 [71]	Posed	210 adults, 3 races; Available: 480 videos	V	Category: 6 basic emotions, and AUs	FACS	Y
Sebe et al. (SD) '04 [119]	Natural: Subjects watched emotion-inducing videos	28 adults	V	Category: Neutral, happy, surprise, disgust	Self-report	N
MMI '05 ³ [106], [98]	Posed: static images, videos recorded simultaneously in frontal and profile view; Natural: Children interacted with a comedian. Adults watched emotion-inducing videos	Posed: 61 adults Natural: 11 children and 18 adults. Overall: 3 races Available: 1250 videos, 600 static images	V	Category: 6 basic emotions, single AU and multiple AUs activation	FACS, Observers' judgment	Y
UT Dallas '06 [95]	Natural: Subjects watched emotion-inducing videos	229 adults	V	Category: 6 basic emotions, puzzle, laughter, boredom, disbelief	Observers' judgment	Y
BU-3DFE (BU)'06 [148]	Posed: 3D range data by using 3DMD digitizer.	100 adults Mixed races	V	Category: 6 basic emotions. Four levels of intensity	N/A	Y
FABO face and body gesture [63]	Posed: two cameras to record facial expressions and body gestures respectively	23 adults Mixed races Available: 210 videos	V	Category: 6 basic emotions, neutral, uncertainty, anxiety, boredom	N/A	Y
Banse-Scherer '96 [8]	Posed	6 actors & 6 actresses Available: 1344 audio samples	A	Category: hot/cold anger, panic fear, anxiety, despair, sadness, elation, happiness, interest, boredom, shame, pride, disgust, contempt.	Listeners' judgment	Y
Danish Emotional Speech Database '96 ⁵	Posed	2 actors & 2 actresses; 2 words, 9 sentences, 2 passages; 10 min of audio data.	A	Category: neutral, surprise, happiness, sadness, anger	Listeners' judgment	Y
ISL meeting corpus '02 [15]	Natural: meeting corpus	18 meetings; Available: data of 5 participants per meeting averagely	A	Category: Positive, neutral, negative [3], [90]	Listeners' judgment	Y
CSC corpus [65]	Natural: subject was motivated to tell the truth and deceive the interviewers in different tasks	32 adults, 15.2 h, 3882 speaking turns, 9687 SUs	A	Deceptive, non-deceptive speech	Self-report	N
Automatic call center (ACC)'05 [83]	Natural: Human-computer dialogue at a commercial call system	1187 calls 7200 utterances	A	Category: Negative, non-negative	listeners' judgment	N
Bank and Stock Service 04 [34]	Natural: human-human dialogue at call center	350 dialogues, 10000 speaking turns	A	Category: fear, anger, stress	Listeners' judgment	N
AIBO database '04 [13]	Natural: children and robot interaction	110 dialogues, 29200 words	A	Category: joyful, emphatic, surprised, ironic, helpless, touchy, angry, bored, motherese, reprimanding, rest	Listeners' judgment	N
Chen-Huang (CH) '00 [21]	Posed	100 adults, 9900 visual and AV expressions	AV	Category: 6 basic emotions, and 4 cognitive states (interest, puzzle, bore, frustration)	N/A	N
Adult Attachment Interview (AAI)'04[111]	Natural: subjects were interviewed to describe the childhood experience	60 adults Each interview last 30-60min	AV	Category: 6 basic emotions, embarrassment, contempt, shame, general positive and negative.	FACS	N
RU-FACS (RU) '05 [10]	Natural: subjects were tried to convince the interviewers they were telling the truth	100 adults	AV	Category: 33 AUs	FACS	N
SAL '05 ⁷	Induced: subjects interacted with artificial listener with different personalities	24 adults 10h	AV	Dimensional labeling/categorical labeling	FEEL-TRACE	Y
Belfast database (BE) '03 [38]	Natural: clips taken from television and realistic interviews with research team	125 subjects. 209 sequences from TV, 30 from interview	AV	Dimensional labeling/categorical labeling	FEEL-TRACE	Y

TABLE 2. VISION-BASED AFFECT RECOGNITION

Legend: exp – Spontaneous/Posed expression, per – person Dependent/Independent, Im/Vi – Image/Video-based, cues – Other Cues besides the face (Head/Body/Eye/Skin/Posture/TaskState/pressureMouse/UserDdefinedClasses/otherContext), rea – realtime (Y–yes, N–no), class – number of classes, sub – number of subjects, samp – sample size, acc – Accuracy, AUs – Action Units corresponding to AU detection, min – minutes, EER – equal error rate, FAR – false acceptance rate, GP – Gaussian process. AAI, BU, CH, CK, FABO, MMI, RU and SD are the database names listed in Table 1. EH – Ekman-Hager database, OD – Other database, ? – missing entry.

References	Facial Feature	Classifier	Performance							
			exp	per	cues	rea	class	sub	samp	acc (%)
Ashraf et al. 07 [4]	AAM	SVM	S	I		N	2	21	?	Im: EER:19%
Bartlett et al. 04 [9]	Gabor wavelets	SVM+HMM	S	I		N	3 AUs	17	Vi: 230+ (OD)	Im/Vi: 75-98
Bartlett et al. 05 [10], [11],	Gabor wavelets	Adaboost SVM	S, P	I		Y	17 AUs	CK+EH: 119, RU:12	Im: 2568(CK+EH) 1689 (RU)	Im: 93.4(CK+EH), 90.5 (RU)
Cohen et al. 03 [25]	12 motion units	Tree-augmented DBN, HMM	P	D, I		Y	6	CH:5 CK:53	Vi: 30 (CH), 53 (CK)	Im: 66.53(CH), 73.22(CK) Vi: 58.63(CH)
Cohn et al. 04 [27]	shape models, Gabor wavelets	LDC	S	I	H	N	3 AUs	21	Im: 99 (OD)	Im: 76 (3-class)
El Kaliouby & Robinson 04 [48]	24 facial points	DBN	P	D	H	Y	6	30	Vi: 164 (OD)	Vi: 77.4
Fasel et al. 04 [50]	Gray-level intensity	NN	P	?	C	?	7	?	Im: 503 (CK)	Im: 38-68
Gunes & Piccardi 05 [61]	Shape features, optical flow	C4.5, Bayes-Net	P	?	B	N	8	FABO:4	Im: 206 (FABO)	Im: 80-100 (various fusion)
Ioannou et al. 05 [67]	FAPs	neurofuzzy network	S	I		N	3	?	Im: 984 (OD)	Im: 78
Ji et al. 06 [69]	Shape features	DBN	S	?	H,E,C	Y	2	8	Vi: 320min (OD)	Correlation coefficient: 95.3
Kapoor & Picard 05 [73]	Facial and head gesture	GP, SVM HMM, NN	S	?	E, P, T	?	2	8	Vi: 136 (OD)	Vi: 86
Kapoor et al. 07 [72]	Pixel difference of mouth region	Same as in [73]	S	I	E P S T M	?	2	24	Vi: 24 (OD)	Vi: 79.17 Baseline: 58
Lee & Elgammal [81]	Pixel intensity of face region	decomposable model	P	I		N	6	CK: 8 OD: 16	Vi: 48 (CK), 80 (OD)	Vi: 39.58 Im: 61.85
Littlewort et al. 07 [87]	Gabor wavelets	Adaboost SVM	S	I		Y	2	26	Vi: 312	Vi: 72
Lucey et al. 07 [88]	AAM	SVM	S,P	I		N	AUs: CK: 15 OD: 4	CK: 100 OD: ?	?	Im: 95 (CK) with 16.66% FAR, 70.47 (OD)
Pantic & Patras 06 [99]	Facial profile points	Rule-based	P	I		N	27 AUs	MMI: 19	Vi: 119 (MMI)	Vi: 86.3
Pantic & Rothkrantz 04 [103]	frontal and profile facial points	Rule-based	P	I		N	32 AUs	MMI: 25	Im: 454 (MMI)	Im: 86
Pantic & Rothkrantz [104]	same as in [103]	Rule-based, case-based	P	I	U	N	9	MMI: 8	Im: 196 (MMI)	Im: 83
Sebe et al. 04 [123]	12 motion units	kNN	S	I			4	CK: 53 SD: 28	Vi: ? (SD), 212+ (CK)	Im: 93 (CK) 95 (SD)
Tong et al. 07 [132]	Gabor wavelets	Adaboost, DBN	P	I		?	14 AUs	CK: 100 OD: 10	Im: 14000 (CK+OD)	Vi: 93.2 (OD), 93.3(CK)
Valstar et al. 04 [136]	Motion history images	SNoW kNN	P	I		N	15 AUs	MMI: 19 CK: 100	Vi: 344 (CK), 253 (MMI)	Vi: 61 (MMI) 68 (CK)
Valstar et al. 06 [135]	8 facial points	gentle boost, SVM	S, P	I		N	2	?	Vi: 60(MMI) 59(CK),70(OD)	Vi: 90.7
Valstar et al. 07 [134]	same as in [103]	GentleSVM-sigmoid	S, P	?	H, B	N	2	MMI: ?	Vi: 100 (P), 102 (S)	Vi: 94%
Wang & Ahuja 03 [137]	Shape and gray-level texture	NN with HOSVD	S	?		?	7	14	Im: 110 (OD)	Im: 84.58
Wang et al. 06 [139]	3D surface labels	LDA	P	I		N	6	BU: 60	Im: 720 (BU)	Im: 83.6
Wen & Huang 03 [141]	Geometric, ratio-image	Exemplars with GMM	P	I		N	4	CK: 47	Im: 2981 (CK)	Im:75.37
Whitehill & Omlin 06 [143]	Haar features	Adaboost	P	I		Y	11 AUs	?	Im: 580 (OD)	Im: 92.35
Yeasin et al. 06 [147]	Pixel intensity of face	kNN + HMM	P, S			N	6	CK: 97 OD:21	Vi: 488 (CK) 108 (OD)	Vi: 90.7 (CK) 72-82 (OD)
Zeng et al. 06 [149]	Texture with LPP	SVDD	S	D		N	2	AAI: 2	Female:7857 Male: 5230	Im: 79(male), 87(female)

As explained above, AUs are very suitable to describe the richness of spontaneous facial behavior, as the thousands of anatomically possible facial expressions can be represented as combination of few dozens of AUs. Hence, the labeling schemes used to code data include FACS AUs (e.g., [10], [71], [106], [98], [111]), Feeltrace system for evaluation-activation dimensional description (e.g., [38], SAL7), self-report (e.g., [123], [65]), and human-observer judgment (e.g., [13], [15], [83], [95], [98]).

The current situation of emotion database research is considerably different from what was described in the comprehensive surveys written by Pantic and Rothkrantz in 2003 [102] and Cowie et al. in 2001 [31]. The current state of the art is advanced and can be summarized as follows (Table 1):

- a database of 3D recordings of acted facial affect [148] and a database of face-and-body recordings of acted affective displays [63] have been made available,
- a collection of acted facial affect displays made from profile-view is shared on Internet [106], [98],
- several large audio, visual and audiovisual sets of human spontaneous affective behavior have been collected, some of which are released for public use.

The existence of these datasets of spontaneous affective behavior is very promising and we expect that this will produce a major shift in the course of the research in the field – from analysis of exaggerated expressions of basic emotions to analysis of naturalistic affective behavior. We also expect subsequent shifts in research in various related fields such as ambient intelligence, transportation, and personal wellness technologies.

3.2 Vision-based Affect Recognition

Because of the importance of face in emotion expression and perception, most of vision-based affect recognition studies focus on facial expression analysis. We can distinguish two main streams in the current research on machine analysis of facial expressions [26], [98]: recognition of affect and recognition of facial muscle action (facial action units). As explained above, facial action units are relatively objective description of facial signals, and can be mapped to the emotion categories based on a high-level mapping such as EMFACS and FACSaid, or to any other set of high-order interpretation categories including complex affective states like depression [47] or pain [144].

As far as automatic facial affect recognition is concerned, most of the existing efforts studied the expressions of the six basic emotions due to their universal properties, their marked reference representation in our affective lives, and the availability of the relevant training and test material (e.g., [71]). There are a few tentative efforts to detect non-basic affective states from deliberately displayed facial expressions including fatigue [60], [69], and mental states like agreeing, concentrated, disagreeing, interested, thinking, confused and frustration (e.g., [48], [72], [73], [129], [147]).

Most of the existing works on automatic facial ex-

pression recognition are based on deliberate and often exaggerated facial displays (e.g., [130]). However, several efforts have been recently reported on automatic analysis of spontaneous facial expression data (e.g., [9], [10], [11], [27], [28], [67], [88], [123], [135], [149], [87], [4], [134]). Some of them study automatic recognition of AUs rather than emotions from spontaneous facial displays (e.g., [9], [10], [11], [27], [28], [135], [134]). Studies reported in [28], [135], [134] and [87] investigated explicitly the difference between spontaneous and deliberate facial behavior. In particular, the studies of Valstar et al. [135], [134], and the study of Littlewort et al. [87] are the first reported efforts to date to automatically discern posed from spontaneous facial behavior. It is interesting to note that, confirming with research findings in psychology (e.g., [47]), the systems proposed by Valstar et al. were built to characterize temporal dynamics of facial actions and employ parameters like speed, intensity, duration, and the co-occurrence of facial muscles activations to classify facial behavior present in a video as either deliberate or spontaneous.

Some of the studies on machine analysis of spontaneous facial behavior were conducted using the datasets listed in Table 1 (e.g., [10], [149], [134]). For other studies new datasets were collected. Overall, the utilized data were collected in the following data-elicitation scenarios: human-human conversation (e.g., [10], [11], [28], [135], [149], [4]), Wizard of OZ scenario (e.g., [67]), or TV broadcast (e.g., [147]). Studies reported in [123], [147] explored automatic recognition of a subset of basic emotional expressions. The study of Zeng et al. [149] investigated separating emotional state from non-emotional states during the Adult Attachment Interview. Studies on separating posed from genuine smiles were reported in [28] and [134] and studies on recognition of pain from facial behavior were reported in [4] and [87].

Most of the existing facial expression recognizers employ various pattern recognition approaches, and are based on 2D spatio-temporal facial features. The usually extracted facial features are either geometric features such as the shapes of the facial components (eyes, mouth, etc.) and the location of facial salient points (corners of the eyes, mouth, etc.) or appearance features representing the facial texture including wrinkles, bulges, and furrows. Typical examples of geometric-feature-based methods are those of Chang et al. [19], who used a shape model defined by 58 facial landmarks, of Pantic and her colleagues [98], [99], [103], [135], [134], who used a set of facial characteristic points around the mouth, eyes, eyebrows, nose, and chin, and of Kotsia and Pitas [77], who used Candide grid. Typical examples of appearance-feature-based methods are those of Bartlett et al. [9], [10], [11], [87], and Guo and Dyer [64], who used Gabor wavelets, of Whitehill and Omlin [143] who used Haar features, of Anderson and McOwen [2], who used a holistic spatial ratio face template, of Valstar et al. [136], who used temporal templates, and of Chang et al. [18], who built

a probabilistic recognition algorithm based on the manifold subspace of aligned face appearances. As suggested in several studies (e.g., [99]), using both geometric and appearance features might be the best choice to design automatic facial expression recognizer. Typical examples of hybrid, geometric- and appearance-feature-based methods, are those proposed by Tian et al. (e.g., [130]), who used facial component shapes and the transient features like crow-feet wrinkles and nasal-labial furrows, and that of Zhang and Ji [158], who used 26 facial points around the eyes, eyebrows, and mouth, and the transient features proposed by Tian et al.. Another example of such a method is

that proposed by Lucey et al. [88], who uses Active Appearance Model (AAM) to capture the characteristics of the facial appearance and the shape of facial expressions.

Most of the existing 2D-feature-based methods are suitable for analysis of facial expressions under a small range of head motions. Thus, most of these methods focus on recognition of facial expressions in near-frontal-view recordings. An exemplar exception is the study of Pantic and Patras [99], who explored automatic analysis of facial expressions from the profile-view of the face.

TABLE 3. AUDIO-BASED AFFECT RECOGNITION

Legend: exp – Spontaneous/Posed expression, per – person Dependent/Independent, cont – contextual information (Subject/Gender/Task/SpeakerRole/SpeakerDependentFeature), class – number of classes, sub – number of subjects, samp – sample size (number of utterances), acc – accuracy, ? – missing entry, BL – Baseline, EER – equal error rate, NPN – negative/neutral/positive, NnN – Negative/non-negative, EnE – emotional/non-emotional, M – male, F – female, A – actor data, R – reading data, W – data of Wizard of OZ. ACC, AIBO, CSC and ISL are the database names listed in Table 1, OD – other database.

References	Feature	Classifier	Performance							
			exp	per	cont	class	sub	samp	acc (%)	other
Ang et al. 02 [3]	Prosody, LM features, position, repeats/ correction	Decision tree	S	I		2	837	21899	64-93	Various label and feature conditions
Austermann et al. 05 [6]	Prosody	Fuzzy rules	S	D, I		5	D: 4 I: 4	D: 280 I: 260	D: 84 I: 60	Robot head data
Batliner et al. 03 [12]	prosody, POS, DA, repetitions, corrections, etc.	MLP, LDA	S, P	I		2	A: 1 R: 19 W: 24	A: 10316 R: 13053 W: 28649	A: 95.7 R: 79.6 W: 74.2	AIBO data
Devillers & Vasilescu 06 [35]	Lexical cues, prosody, spectrum, disfluency, etc.	SVM	S	I		4	680	2258	Lexical: 78 paralinguistic: 60	Medical emergency center data
Forbes-Riley & Litman 04 [52]	prosodic, lexical, syntactic, dialogue features, etc.	boost decision tree	S	I	Su, G, T	3	17	453	84.75	computer tutor data
Graciarena et al. 06 [57]	Prosodic, acoustic, lexical	SVM, GMM	S	D		2	32	9328	64.4	CSC data
Hirschberg et al. 05 [65]	Prosodic, acoustic, lexical	Ripper rule-induction	S	D	Dep	2	32	9491	66.4	CSC data
Kwon et al. 03 [79]	Prosody, MFCC	QDA, SVM, HMM LDA	P	D, I		2, 4, 5	OD: 9; AIBO: 14	OD: 8820; AIBO: 3534	2 class: 96 4 class: 70.1 5 class: 42.3	AIBO and OD data
Lee & Narayanan 03 [82]	Prosody	Fuzzy inference	S	I		2	?	F: 776; M: 591	F: 73 M: 63	
Lee & Narayanan 05 [83]	prosody, lexical, and discourse	LDC, kNN	S	I	G	2	ACC: 1187 calls	7200	M: 89.55 F: 92.1	M: 76.5%BL F: 74.1%BL
Liscombe et al. 05 [84]	Acoustic-prosodic	C4.5 with Adaboost	S	I		3	17	6778 turns	76.42	60% BL
Litman & Forbes-Riley 04 [86]	Acoustic-prosodic, lexical	Boost decision tree	S	I	Su, G, T	2, 3	10	333	NPN: 47-67 NnN: 64-72 EnE: 52-75	Various label and feature conditions
Matos et al. 06 [91]	MFCC	HMM	S	I		2	19	Train: 2473 Test: 2155	82	
Neiberg et al. 06 [94]	MFCC, MFCC-low, pitch	GMM	S	I		3		OD: 7619 ISL: 12479	OD: 90 ISL: 80	Swedish, English
Schuller et al. 05 [120]	Acoustic-Prosodic, linguistic	StackingC MLR, NB, ND SVM, C4.5	S, P	D, I		7	13+	4336	I: 76.4 D: 94.8	
Steidl et al. 05 [125]	Prosodic, POS features	?	S	I		4	AIBO: 51	6071	60	Entropy measure
Truong & van Leeuwen 07 [133]	Spectral, prosodic	GMM+ SVM	S	I		2	OD1: 34 OD2: 8	OD1: 6838 OD2: 335	EER: 2.9-7.5	English, Dutch
Vasilescu & Devillers 05 [36]	prosodic, spectral, disfluency, etc.	SVM, logistic model tree	S	I	R	2	404	800	82	Same database as [35]
Zhang et al. 04 [157]	Lexical, prosodic, spectral, syntactic	CART tree	S	I		3	OD: 17	714	91.3	

TABLE 4. AUDIOVISUAL AFFECT RECOGNITION

Legend: Fusion – Feature/Decision/Model-level, exp – Spontaneous/Posed expression, per – person-Dependent/Independent, class – number of classes, sub – number of subjects, samp – sample size (number of utterances), cue – other Cues (Lexical/Body), acc – Accuracy, RR – mean with weighted recall values, FAP – facial animation parameter, ? – missing entry. AAL, CH, SAL and SD are the database names listed in Table 1

References	Feature	Fusion	Classifier	Performance							
				exp	per	cue	class	sub	samp	acc (%)	other
Busso et al. 04 [16]	102 markers, prosody	F, D	SVM	P	D		4	1	256 sentences	89	
Caridakis et al. 06 [17]	facial points, prosody	M	RNN	S	I		4	SAL 4	1000 tunes	79	
Fragopanagos and Taylor 05 [53]	17 FAPs, prosody	M	ANNA	S	I	L	4	SAL 4	500 epochs	44-71	various labels/features
Go et al. 03 [56]	Eigenfaces, MFCC	D	LDA	P	I		6	20	360 utterances	95-98	
Hoch et al. 05 [66]	Gabor feature, prosody	D	SVM	P	D		3	7	840 sequences	90.7	car setting
Karpouzis et al. 07 [74]	19 FPs, prosody	M	RNN	S	I	B	4	SAL 4	1000 tunes	82	
Pal et al. 06 [97]	Vertical gray level, F0-F3	D	Rules, k-means	S	D		5	1	?	75.2	
Petridis & Pantic 08 [108]	facial points, prosody	F, D	Adaboost + NN	S	I	B	2	8	96 laughter/speech episodes	86.9	AMI data ⁶
Schuller et al. 07 [118]	AAM, prosody, articulatory, voice quality, lexical	F	SVM	S	I	B	3	21	10.5 hours	recall: 41.7-63.9 (RR)	balance training
Sebe et al. 06 [122]	12 motion units, prosody	M	BN	P	D		11	SD 38	1254 sentences	90	
Song et al. 04 [124]	54 FAPs, prosody	M	THMM	P	?		7	?	?	84.7	
Wang & Guan 05 [138]	Gabor wavelets, prosody, MFCC, formants.	D	FLDA	P	I		6	8	500 sentences	82.14	6 languages
Zeng et al. 06 [150]	12 motion units, prosody	M	MFHMM	P	I		11	CH 20	660 sentences	83	
Zeng et al. 07 [151]	Texture with LLP, prosody	D	Adaboost + MHMM	S	D		2	AAI 2	137 utterances	89	
Zeng et al. 04 [153]	motion units, prosody	D	SNoW	P	D		11	CH 38	1254 sentences	89-90	
Zeng et al. 05 [154]	motion units, prosody	M	MFHMM	P	I		11	CH 20	660 sentences	80.61	
Zeng et al. 07 [155]	motion units, prosody, formants	D	HMM	P	D, I		11	CH 20	660 sentences	I: 72.42 D: 96.3	
Zeng et al. 05 [156]	motion units, prosody	F	Fisher-Boosting	P	D		4	CH 20	660 sentences	84-87	

A few approaches to automatic facial expression analysis are based on 3D face models. Huang and his colleagues (i.e., [25], [123], [141], [149]) used features extracted by a 3D face tracker called Piecewise Bezier Volume Deformation Tracker [128]. Cohn et al. [27] focused on analysis of brow action units and head movement based on a cylindrical head model [146]. Chang et al. [19] and Yin et al. [139], [148] used 3D expression data for facial expression recognition. The progress of the methodology based on 3D face models may yield view-independent facial expression recognition, which is important for spontaneous facial expression recognition because the subject can be recorded in less controlled, real-world settings.

Some efforts are reported to decompose multiple factors (e.g., the facial expression, face style, or pose) from face images. Typical examples are those of Wang and Ahuja [137], who used multi-linear subspace method, and of Lee and Elgammal [81], who proposed decomposable

nonlinear manifold, to estimate facial expression and face style simultaneously. The study of Zhu and Ji [160] used a normalized SVD decomposition to recover facial expression and pose.

Relatively few studies investigated the fusion of the information from facial expressions and head movement (e.g., [27], [69], [158], [160], [134]), the fusion of facial expression and body gesture (e.g., [7], [61], [62], [134]), and the fusion of facial expressions and postures from a sensor chair (e.g., [72], [73]), with the aim at improvement of affect recognition performance.

Finally, virtually all present approaches to automatic facial expression analysis are context insensitive. Exceptions from this overall state of the art in the field include just a few studies. For example, Pantic and Rothkrantz [104] and Fasel et al. [50] investigated interpretation of facial expressions in terms of user-defined interpretation labels. Ji et al. [69] investigated the influence of context (work condition, sleeping quality, circadian rhythm, and

environment, physical condition) on fatigue detection, and Kapoor and Picard [73] investigated the influence of the task states (difficulty level and game state) on interest detection.

Table 2 provides an overview of the currently existing, exemplar systems for vision-based affect recognition with respect to the utilized facial features, classifier, and performance. While summarizing the performance of the surveyed systems, we also mention a number of relevant aspects including the type of the utilized data (spontaneous or posed, number of different subjects; sample size), whether the system is person-dependent/ independent, whether it performs in real time condition, what is the number of target classification categories, whether and which other cues besides the face have been used in the classification (head/ body/ eye/ posture/ task state/ other context), whether the system processes still images or videos, and how accurately it performs the target classification. A missing entry means that the matter at issue was not reported or it remained unclear from the available literature. For instance, some studies did not explicitly indicate whether the recordings of the same subjects were used as both the testing data and the training data. Hence, it remains unclear whether these systems perform in a subject-independent manner. It is important to stress that we cannot rank the performances of the surveyed systems because each of the relevant studies has been conducted under different experimental conditions using different data, different testing methods (such as person-dependent/independent), and different performance measurements (accuracy, equal error rate, etc.).

The research in machine analysis of facial affect has seen a lot of progress when compared to that described in the survey paper of Pantic and Rothkrantz from 2003 [102]. The current state of the art in the field is as follows:

- Methods have been proposed to detect attitudinal and non-basic affective states such as confusion, boredom, agreement, fatigue, frustration, and pain from facial expressions (e.g., [69], [72], [129], [147], [87]).
- Initial efforts were conducted to analyze and automatically discern posed (deliberate) facial displays from genuine (spontaneous) displays (e.g., [135], [134]).
- First attempts are reported towards vision-based analysis of spontaneous human behavior based on 3D face models (e.g., [123], [149]), based on fusing the information from facial expressions and head gestures (e.g., [27], [134]), and based on fusing the information from facial expressions and body gestures (e.g., [61]).
- Few attempts have been also made towards context-dependent interpretation of the observed facial behavior (e.g., [50], [69], [72], [104]).
- Advanced techniques in feature extraction and classification have been applied and extended in this field. A few real-time robust systems have been built (e.g., [11]) thanks to the advance of relevant techniques such as real-time face detection and object tracking.

3.3 Audio-based Affect Recognition

Research in vocal affect recognition is also largely influenced by basic emotion theory. In turn, most of the exist-

ing efforts in this direction aim at recognition of a subset of basic emotions from speech signals. However, a few tentative studies were published recently on interpretation of speech signals in terms of certain application-dependent affective states. These studies are those of Hirschberg et al. [65] and Graciarena et al. [57], who attempted deception detection, of Liscombe et al. [84], who focused on detecting certainty, of Kwon et al. [79], who reported on stress detection, of Zhang et al. [157], who investigated speech-based analysis of confidence, confusion, and frustration, of Batliner et al. [12], who aimed at detecting trouble, of Ang et al. [3], who explored speech-based recognition of annoyance and frustration, and of Steidl et al. [125], who conducted studies on detection of empathy. In addition, few efforts towards automatic recognition of nonlinguistic vocalizations like laughters [133], coughs [91] and cries [97] have also been reported recently. This is of particular importance for the research in machine analysis of human affects since recent studies in cognitive sciences showed that listeners seem to be rather accurate in decoding some non-basic affective states such as distress, anxiety, boredom, and sexual interest from nonlinguistic vocalizations like laughs, cries, sighs, and yawns (e.g., [113]).

Most of the existing systems for automatic vocal affect recognition were trained and tested on speech data that was collected by asking actors to speak prescribed utterances with certain emotions (e.g., [6], [79]). As the utterances are isolated from the interaction context, this experimental strategy precludes finding and using correlations between the paralinguistic displays and the linguistic content, which seem to play an important role for affect recognition in daily interpersonal interactions.

Based on the above consideration, researchers started to focus on affect recognition in naturalistic audio recordings collected in call centers (e.g., [35], [82], [83], [94]), meetings (e.g., [94]), wizard of OZ scenarios (e.g., [12]), interview (e.g., [65]), and other dialogue systems (e.g., [14], [86]). In these natural interaction data, affect displays are often subtle, and basic emotion expressions seldom occur. It is therefore not surprising that recent studies in the field, which are based on such data, attempt to detect either coarse affective states, i.e., positive, negative, and neutral states (e.g., [82], [83], [86], [94]), or application-dependent states mentioned above, rather than basic emotions.

Most of the existing approaches to vocal affect recognition used acoustic features as classification input, based on the acoustic correlation for emotion expressions that was summarized in [31]. The popular features are prosodic features (e.g., pitch-related feature, energy-related features, speech rate), and spectral features (e.g., MFCC, cepstral features). Many studies show that pitch and energy among these features contribute most to affect recognition (e.g., [79]). An exemplar effort is that of Vasilescu and Devillers [36], who show the relevance of speech disfluencies (e.g., filler and silence pauses) to affect recognition.

With the research shift towards analysis of spontaneous human behavior, analysis of acoustic information

only will not suffice for identifying subtle changes in vocal affect expression. As indicated by Batliner et al. [12], “the closer we get to a realistic scenario, the less reliable is prosody as an indicator of the speaker’s emotional state”. In the preliminary experiments of Devillers and Vidrascu [35], using lexical cues resulted in a better performance than using paralinguistic cues to detect relief, anger, fear and sadness in human-human medical call conversations. In turn, several studies investigated the combination of acoustic features and linguistic features (language and discourse) to improve vocal affect recognition performance. Typical examples of linguistic-paralinguistic-fusion methods are those of Litman et al. [86] and Schuller et al. [120], who used spoken words and acoustic features, of Lee and Narayanan [83], who used prosodic features, spoken words and information of repetition, of Graciararena et al. [57], who combined prosodic, lexical and cepstral features, and of Bartliner et al. [12], who used prosodic features, Part-of-speech (POS), dialogue act (DA), repetitions, corrections, and syntactic-prosodic boundary to infer the emotion. Litman et al. [86] and Forbes-Riley and Litman [52] investigated also the role of the context information (e.g. subject, gender, turn-level features representing local and global aspects of the dialogue) on audio affective recognition.

Although the above studies indicated recognition improvement by using information of language, discourse and context, automatic extraction of these related features is a difficult problem. First, existing automatic speech recognition systems cannot reliably recognize the verbal content of emotional speech (e.g., [5]). Second, extracting semantic discourse information is even more challenging. Most of these features are typically extracted manually or directly from transcripts.

Table 3 provides an overview of the currently existing, exemplar systems for audio-based affect recognition with respect to the utilized auditory features, classifier, and performance. As in Table 2, we specify relevant aspects in Table 3 to summarize the reported performance of surveyed systems.

The current state of the art in the research field of automatic audio-based affect recognition can be summarized as follows:

- Methods have been proposed to detect non-basic affective states, including coarse affective states such as negative and non-negative states (e.g., [83]), application-dependent affective states (e.g., [3], [12], [65], [79], [157], [125]), and nonlinguistic vocalizations like laughter, cry (e.g., [133], [91], [97]).
- A few efforts have been made to integrate paralinguistic features and linguistic features such as lexical, dialogic, and discourse feature (e.g., [12], [35], [57], [83], [86], [120]).
- Few investigations have been conducted to make use of contextual information to improve the affect recognition performance (e.g., [52], [86]).
- Few studies have been reported to recognize the affective states across languages (e.g., [94], [133]).
- Some studies have investigated influence of ambiguity of human labeling on recognition performance (e.g., [3]

[86]), and proposed measures to compare human labelers and machine classifiers (e.g., [125]).

- Advanced techniques in feature extraction, classification and natural language processing have been applied and extended in this field. Some studies have been tested on commercial call data (e.g., [83], [35]).

3.4 Audiovisual Affect Recognition

In the survey of Pantic and Rothkrantz in 2003, [102], only four studies were found that were focused on audiovisual affect recognition. Since then, an increasing number of efforts are reported in this direction. Similar to the state of the art in single-modal affect recognition, most of the existing audio-visual affect recognition studies investigated recognition of the basic emotions from deliberate displays. Relatively few efforts have been reported toward detection of non-basic affective states from deliberate displays. Those include the work of Zeng et al. [150], [153], [154], [155], and that of Sebe et al. [122], who added 4 cognitive states (interest, puzzlement, frustration and boredom) considering the importance of these cognitive states in human computer interaction. Related studies conducted on naturalistic data include that of Pal et al. [97], who designed a system to detect hunger and pain as well as sadness, anger, and fear from infant facial expressions and cries, and that of Petridis and Pantic [108], who investigated separating speech from laughter episodes based on both facial and vocal expression.

Most of the existing methods for audiovisual affect analysis are based on deliberately posed affect displays (e.g., [16], [56], [66], [122], [124], [138], [150], [153], [154], [155]). Recently a few exceptional studies have been reported toward audiovisual affect analysis in spontaneous affect displays (e.g., [17], [53], [74], [97], [151], [108]). Zeng et al. [151], used the data collected in psychological research interview (Adult Attachment Interview), Pal et al. [97] used recordings of infants [97], Petridis and Pantic [108] used the recordings of people engaged in meetings (AMI corpus⁶), while Fragopanagos and Taylor [53], Caridakis et al. [17], and Karpouzis et al. [74], used the data collected in Wizard of OZ scenarios. Since the available data were usually insufficient to build a robust machine learning system for recognition of fine-grained affective states (e.g., basic emotions), recognition of coarse affective states was attempted in most of the aforementioned studies. Studies of Zeng et al. focus on audiovisual recognition of positive and negative affect [151], while other studies report on classification of audiovisual input data into the quadrants in evaluation-activation space [17], [53], [74]. The studies reported in [17], [53], [74] applied the FeelTrace system that enables raters to continuously label changes in affective expressions. However, note that the study discussed in [53] reported on a considerable labeling variation among four human raters due to the subjectivity of audio-visual affect judgment. More specifically, one of the raters mainly relied on audio information when making judgments while another rater mainly relied on visual information. This experiment actually also reflects the asynchronization of audio and visual expression. In order to reduce this variation of human

labels, the studies of Zeng et al. [151] made the assumption that facial expression and vocal expression has the same coarse emotional states (positive and negative), and then directly used FACS-based labels of facial expressions as audio-visual expression labels.

The data fusion strategies utilized in the current studies on audiovisual affect recognition are either feature-level or decision-level or model-level fusion. Typical examples of feature-level fusion are those reported in [16], [118], [156] which concatenated the prosodic features and facial features to construct joint feature vectors which are then used to build an affect recognizer. However, the different time scales and metric levels of features coming from different modalities, as well as increasing feature-vector dimensions influence the performance of a affect recognizer based on a feature-level fusion. The vast majority of studies on bimodal affect recognition reported on decision-level data fusion (e.g., [16], [56], [66], [97], [151], [153], [155], [138], [108]). In the decision-level data fusion, the input coming from each modality is modeled independently and these single-modal recognition results are combined at the end. Since humans display audio and visual expressions in a complementary and redundant manner, the assumption of conditional independence between audio and visual data streams in decision-level fusion is incorrect and results in the loss of information of mutual correlation between the two modalities. To address this problem, a number of model-level fusion methods have been proposed that aim at making use of the correlation between audio and visual data streams, and relax the requirement of synchronization of these streams (e.g., [17], [53], [122], [124], [150], [154]). Zeng et al. [154] presented Multi-stream Fused HMM to build an optimal connection among multiple streams from audio and visual channels according to maximum entropy and the maximum mutual information criterion. Zeng et al. [150] extended this fusion framework by introducing a middle-level training strategy under which a variety of learning schemes can be used to combine multiple component HMMs. Song et al. [124] presented tripled HMM to model correlation properties of three component HMMs that are based individually on upper face, lower face, and prosodic dynamic behaviors. Fragopanagos and Taylor [53] proposed an artificial neural network with a feedback loop called ANNA to integrate the information from face, prosody and lexical content. Caridakis et al. [17], Karpouzis et al. [74], and Petridis and Pantic [108] investigated combining the visual and audio data streams by using Neural Networks (NN). Sebe et al. [122] used Bayesian Network (BN) to fuse the facial expression and prosody expression.

Table 4 provides an overview of the currently existing, exemplar systems for audiovisual affect recognition with respect to the utilized auditory and visual features, classifier, and performance. As in Tables 2 and 3, we also specify a number of relevant issues in Table 4 to summarize the reported performance of surveyed systems.

In summary, the research on audiovisual affect recognition has witnessed significant progress in the past few years as follows:

- Efforts have been reported to detect and interpret non-basic genuine (spontaneous) affective displays in terms of coarse affective states such as positive and negative affective states (e.g., [151]), quadrants in evaluation-activation space (e.g., [17], [53], [74]), and application-dependent states (e.g., [122], [154], [97], [108]).
- Few studies have been reported on efforts to integrate other affective cues besides the face and the prosody, such as body and lexical features (e.g., [53], [74]).
- Few attempts have been made to recognize affective displays in specific naturalistic settings (e.g., in a car [66]) and in multiple languages (e.g., [138]).
- Various multimodal data fusion methods have been investigated. In particular, some advanced data fusion methods have been proposed such as HMM-based fusion (e.g., [124], [154], [150]), NN-based fusion (e.g., [53], [74]), and BN-based fusion (e.g., [122]).

4 CHALLENGES

The studies reviewed in the previous section indicate two new trends in the research on automatic human affect recognition: analysis of spontaneous affective behavior and multimodal analysis of human affective behavior including audiovisual analysis, combined linguistic and nonlinguistic analysis, and multi-cue visual analysis based on facial expressions, head movements, and/or body gestures. Several previously-recognized problems have been addressed including the development of more comprehensive datasets of training and testing material. At the same time, several new challenging issues have been recognized, including the necessity of studying the temporal correlations between the different modalities (audio and visual) as well as between various behavioral cues (e.g., facial, head, and body gestures). This section discusses these issues in detail.

4.1 Databases

Acquiring valuable spontaneous affective behavior data and the related ground truth is far from being solved. While it is relatively easy to elicit joyful laughter by showing clips from comedies to subjects, the majority of affective states are much more difficult (if possible at all) to elicit (e.g., fear, stress, sadness, or anger — which is particularly difficult to elicit in any laboratory setting, including face-to-face conversation [23]). Social psychology has provided a host of creative strategies for inducing emotion, which seem to be useful for collecting affective expressions that are difficult to elicit in the laboratory, and affective expressions that are contextually complex (such as embarrassment), or for research programs that emphasize the “mundane realism” of experimentally elicited emotions [23]. However, engineers, who are usually the designers of the databases of human behavior data, are often not even aware of these strategies, let alone putting them into the practice. This situation needs to be changed if the challenging and crucial issue of collecting valuable data on human spontaneous affective behavior is to be addressed.

Although many efforts have been done toward collec-

tion of databases of spontaneous human affective behavior, most of the data contained in the available databases currently lack labels. In other words, no metadata is available that could identify the affective state displayed in a video sample and the context in which this affective state was displayed. There are several related issues.

First, it is not clear which kind of metadata need to be provided. While data labeling is easy to accomplish in the case of prototypical expressions of emotions, it becomes a real challenge once we move beyond the six basic emotions. To reduce the subjectivity of data labeling, it is generally accepted that human facial expression data need to be FACS coded. The main reason is that FACS AUs are objective descriptors and independent of interpretation, and can be used for any high-level decision making process including recognition of affective states. However, while this solves the problem of attaining objective facial behavior coding, how to objectively code vocal behavior remains an open issue. Nonlinguistic vocalizations like laughter, coughs, cries, etc., can be labeled as such, but there is no set of interpretation-independent codes to label emotional speech. Another related issue is that of culture and context dependency. The metadata about the context in which the recordings were made such as the utilized stimuli, the environment, and the presence of other people, is needed since these contextual variables may influence masking of the emotional reactions.

Second, even if labeled data are available, engineers responsible to design an automated human affect analyzer usually assume that the data are accurately labeled. This assumption may or may not be accurate [26], [125]. The reliability of the coding can be ensured by asking several independent human observers to conduct the coding. If the inter-observer reliability is high, the reliability of the coding is assured. Inter-observer reliability can be improved by providing thorough training to observers on the utilized coding schemes such as FACS. When it comes to data coding in terms of affect labels, a possible method is to use multi-label multi-time-scale system in order to reduce the subjectivity of human judgment and to represent comprehensive properties of affect displays [37], [80].

Third, human labeling of affective behavior is very time consuming and expensive. In the case of facial expression data, it takes more than one hour to manually score 100 still images or a minute of video sequence in terms of AUs [43]. A remedy could be the semi-supervised active learning method [159] that is to combine semi-supervised learning [24] and active learning [51]. The semi-supervised learning mechanism aims at making use of the unlabeled data, and the active learning mechanism aims at enlarging the useful information conveyed by human feedback (annotation in this application), and provides the annotators the most ambiguous samples according to the current emotion classifier. More specifically, several promising prototype systems were reported in the past few years that can recognize deliberately produced AUs in either (near-) frontal view face images (e.g., [98], [130]) or profile-view face images (e.g., [99]). Although these systems will not be always able to be generalized to the subtlety and complexity of human

affective behaviour occurring in real-world settings, they can be used to attain an initial data labeling that can be subsequently controlled and corrected by human observers. However, as this has not been attempted in practice, there is no guarantee that such an approach will actually reduce the time needed for obtaining the ground truth. Future research is needed to determine whether this attempt is feasible.

Although much effort has been done toward collection of databases of spontaneous human affective behavior, many of these datasets are not publicly available (see Table 1). Some are still under construction, some are in the process of data publication, and some seem to have dim prospects of being published due to lack of appropriate agreement of subjects. More specifically, spontaneous displays of emotions, especially in multimedia format, reveal personal and intimate experience; privacy issues jeopardize the public accessibility of many databases.

Besides these problems concerned with acquiring valuable data, the related ground truth, and the agreement of subjects to make the data publicly available, another important issue is how one construct and administer such a large affective expression benchmark database. A noteworthy example is the MMI facial expression database [98], [106], which was built as a web-based direct-manipulation application, allowing easy access and easy search of the available images. In general, in the case of publicly available databases, once the permission for usage is issued, large, unstructured files of material are sent. Such unstructured data is difficult to explore and manage. Pantic et al. [102], [106] and Cowie et al. [32], emphasized a number of specific, research and development efforts needed to build a comprehensive, readily accessible reference set of affective displays that could provide a basis for benchmarks for all different efforts in the research on machine analysis of human affective behavior. Nonetheless, note that their list of suggestions and recommendations is not exhaustive of worthwhile contributions.

4.2 Vision-based Affect Recognition

Although several efforts discussed in section 3.2 were recently reported on machine analysis of spontaneous facial expressions, the problem of automatic analysis of facial behavior in unconstrained environments is still far from being solved.

Existing methods for machine analysis of facial affect typically assume that the input data are near frontal- or profile-view face image sequences showing non-occluded facial displays captured under constant lighting condition against a static background. In real interaction environment, such assumption is often invalid. Development of robust face detectors, head-, and facial feature trackers, which will be robust to arbitrary head movement, occlusions, and scene complexity like the presence of other people and dynamic background, forms the first step in the realization of facial affect analyzers capable of handling unconstrained environments. View-independent facial expression recognition based on 3D face model (e.g., [20], [148]) or multi-view face models (e.g., [160])

may be a (part of the) solution.

As mentioned already in section 2, a growing body of research in cognitive sciences argues that the dynamics of human behavior are crucial for its interpretation (e.g., [47], [113]). For instance, it has been shown that spontaneous smiles are longer in total duration, can have multiple apexes (multiple rises of the mouth corners), appear before or simultaneously with other facial actions such as the rise of the cheeks, and are slower in onset and offset time than the posed smiles (e.g., a polite smile) [28]. In spite of these findings, the vast majority of the past work in the field does not take dynamics of facial expressions into account when analyzing shown facial behavior. Some of the past work in the field has used aspects of temporal structure of facial expression such as the speed of a facial point displacement or the persistence of facial parameters over time (e.g., [87], [132], [158]). However, just few recent studies analyze explicitly the temporal structure of facial expressions (e.g., [98], [99], [135], [132], [134]). In addition, it remains unresolved how the grammar of facial behavior can be learned and how this information can be properly represented and used to handle ambiguities in the input data [100], [102].

Except for few studies (e.g., [27], [61]), the existing efforts toward machine analysis of human facial behavior focus only on the analysis of facial gestures without taking into consideration other visual cues like head movements, gaze direction, and body gestures. However, research in cognitive science reports that human judgments of behavioral cues are the most accurate when both of the face and the body are taken into account (e.g., [1], [117]). This seems to be of particular importance when judging certain complex mental states such as embarrassment [75]. However, integration, temporal structures and temporal correlations between different visual cues are virtually unexplored areas of research. One noteworthy study that investigated fully automatic coding of human behavior dynamics with respect to both temporal segments (onset, apex, offset, neutral) of various visual cues and temporal correlation between different visual cues (facial, head, and shoulder movements) is that by Valstar et al. [134], who investigated separating posed from genuine smiles in video sequences.

4.3 Audio-based Affect Recognition

One challenge in audio expression analysis is how to identify affect-related features in speech signals. When our aim is to detect spontaneous emotion expressions, we have to take into account both linguistic and paralinguistic cues that mingle together in audio channel. Although a number of linguistic and paralinguistic features (e.g. prosodic, dysfluency, lexicon, and discourse features) were proposed in the body of literature on vocal affect recognition, the optimal feature set has not yet been established.

Another challenge is how to reliably extract these linguistic and paralinguistic features from the audio signals in an automatic way. When prosody is analyzed in a naturalistic conversation, we have to consider the multiple functions of prosody that include information about the expressed affect as well as a variety of linguistic func-

tions [93]. A prosodic event model that could reflect both linguistic and paralinguistic (affective) functions simultaneously would be an ideal solution. Automatic extraction of spoken words from spontaneous emotional speech is still a difficult problem – the recognition rate of the exiting automatic speech recognition (ASR) systems drops significantly as soon as emotional speech needs to be processed. Some tentative studies on adapting an ASR system to emotional speech were reported in [5], [119]. We hope that in the future more such studies will be conducted. In addition, automatic extraction of high-level semantic linguistic information (e.g. dialogue act, repetitions, corrections, and syntactic information) is an even more challenging problem which remains open in the research field of natural language processing.

It is interesting to note that some mental states such as frustration and boredom seem to be identifiable from non-linguistic vocalizations like sighs and yawns [113]. Few efforts towards automatic recognition of non-linguistic vocalizations like laughers [133], [108], cries [97], and coughs [91] have been also recently reported. However, no effort towards human affect analysis based on vocal outbursts has been reported so far.

4.4 Audiovisual Affect Recognition

The research on audiovisual affect analysis in naturalistic data is still in its pioneering phase. While all agree that multisensory fusion including audiovisual data fusion, linguistic and paralinguistic data fusion, multi-visual-cue data fusion would be highly beneficial for machine analysis of human affect, it remains unclear how this should be accomplished. Studies in neurology on fusion of sensory neurons [126] are supportive of early data fusion (i.e., feature-level data fusion) rather than of late data fusion (i.e., decision-level fusion). However, it is an open issue how to construct suitable joint feature vectors composed of features from different modalities with different time scales, different metric levels and different temporal structures. Simply concatenating audio and video features into a single feature vector, as done in the current human affect analyzers that use feature level data fusion, is obviously not the solution to the problem.

Due to these difficulties, most researchers choose decision-level fusion in which the input coming from each modality is modeled independently and these single-modal recognition results are combined at the end. Decision-level fusion, also called classifier fusion, is now an active area in machine learning and pattern recognition field. Many studies have demonstrated the advantage of classifier fusion over the individual classifiers due to the uncorrelated errors from different classifiers (e.g., [78], [112]). Various classifier fusion methods (fixed rules and trained combiners) have been proposed in literature, but optimal design methods for classifier fusion are still not available. In addition, since humans simultaneously employ the tightly coupled audio and visual modalities, the multimodal signals cannot be considered mutually independent and should not be combined only at the end as is the case in decision-level fusion.

Model-level fusion or hybrid fusion that aims at com-

binning the benefits of both feature-level and decision-level fusion methods may be a good choice for this fusion problem. However, based on existing knowledge and methods, how to model multimodal fusion based on multi-label multi-time-scale labeling scheme mentioned above is largely unexplored. A number of issues relevant to fusion require further investigation, such as the optimal level of integrating these different streams, the optimal function for the integration, as well as inclusion of suitable estimations of reliability of each stream. In addition, how to build context-dependent multimodal fusion is an open and highly relevant issue.

Here we want to stress that temporal structures of the modalities (facial and vocal) and their temporal correlations play an extremely important role in interpretation of human naturalistic, audiovisual affective behavior (see section 2 for a discussion). Yet, these are virtually unexplored areas of research, due to the fact that facial expression and vocal expression of emotion are usually studied separately.

4.5 A Few Additional Related Issues

Context: An important related issue that should be addressed in all visual, vocal, and audiovisual affect recognition is how to make use of information about the context (environment, observed subject, his or her current task) in which the observed affective behavior was displayed. Affects are intimately related to a situation being experienced or imagined by human. Without context, human may misunderstand the observed person's emotion expressions. Yet, with the exception a few studies investigated the influence of context on affect recognition (e.g., [50], [52], [69], [72], [86], [104]), virtually all existing approaches to machine analysis of human affect are context insensitive. Building a context model that includes person ID, gender, age, conversation topic, and workload need the help from other research fields like face recognition, gender recognition, age recognition, topic detection, and task tracking. Since the problem of context sensing is very difficult to solve, pragmatic approaches (e.g. activity- and/ or subject-profiled approaches) should be taken.

Segmentation: Almost all of existing methods are tested just on pre-segmented emotion sequences or images, except few studies (e.g., [11], [25]) that use heuristic methods to segment the emotions from videos. Automatic continuous emotion recognition is a dynamic searching process that is to continuously make emotion inference in the presence of signal ambiguity and context. This is rather complicated, since the search algorithm has to consider the possibility of each emotion starting at any arbitrary time frame. Furthermore, the number of emotion changing in a video is not known, and the boundaries between different emotional expressions are full of ambiguity. It becomes more challenging in multimodal affect recognition because different modalities (e.g., face, body, vocal expressions) have difference temporal structures and often do not synchronize. If we aim at developing a practical affect recognizer, the emotion segmentation is one of the most important issues, but has not been largely unexplored so far.

Evaluation: Existing methods for machine analysis of human affect surveyed and discussed throughout this paper are difficult to compare because they are rarely (if ever) tested on a common dataset. United efforts of the relevant research communities are needed to specify evaluation procedures that could be used for establishing reliable measures of systems' performance based on a comprehensive, readily accessible benchmark database.

5 CONCLUSION

The research in machine analysis of human affect has witnessed a good deal of progress when compared to that described in the survey papers of Pantic and Rothkrantz from 2003 [102] and Cowie et al. in 2001 [31]. At that time, a few small-sized datasets of affective displays existed, and almost all methods for machine analysis of human affect were uni-modal, based on deliberate displays of either facial expressions or vocal expressions of six prototypical emotions. Available data was not shared among researchers, multimedia data and multimodal human affect analyzers were rare, and machine analysis of spontaneous displays of affective behavior seemed to be in a distant future. Today, several large collections of acted affective displays are shared by the researchers in the field and some datasets of spontaneously displayed expressions have been recently made available. A number of promising methods for vision-based, audio-based, and audiovisual analysis of human spontaneous behavior have been proposed. This paper focused on surveying and discussing these novel approaches to machine analysis of human affect as well as on summarizing the issues that have not received sufficient attention but are crucial for advancing machine interpretation of human behavior in naturalistic contexts. The most important of these issues yet to be addressed in the field include the following:

- Build a comprehensive, readily accessible reference set of affective displays that could provide a basis for benchmarks for all different efforts in the research on machine analysis of human affective behavior. Define the appropriate evaluation procedures.
- Develop methods for spontaneous affective behavior analysis that are robust to observed person's arbitrary movement, occlusion, complex and noisy background.
- Devise models and methods for human affect analysis that take into the consideration temporal structures of the modalities and temporal correlations between the modalities (and/or multiple cues), and context (subject, his or her task, environment).
- Develop better methods for multimodal fusion.

Since the complexity of these issues concerned with the interpretation of human behavior at a deeper level is tremendous and spans several different disciplines in computer and social sciences, we believe that a large, interdisciplinary, international program directed towards computer understanding of human behavioral patterns should be established if we are to experience true breakthroughs in this and the related research fields. The progress in research on machine analysis of human affect can aid in the creation of a new paradigm for HCI (affect-

sensitive interfaces, socially intelligent environments), and advance the research in several related fields including psychology, psychiatry, and education.

ACKNOWLEDGMENT

This paper is collaborative work. Thomas Huang is the leader of this team work but prefers to be the last in the author list. Zhihong Zeng wrote the first draft, Maja Pantic significantly improved it by rewriting it and offering important advice, and Glenn Roisman provided important comments and polished the whole paper. We would like to thank Qiang Ji and anonymous reviewers for encouragement and valuable comments. This work was supported in part by Beckman Postdoctoral Fellowship and NSF CCF 04-26627.

REFERENCES

- [1] Ambady, N., Rosenthal, R. (1992). Thin slices of expressive behavior as predictors of interpersonal consequences: A meta-analysis. *Psychological Bulletin*, Vol. 111, No. 2, 256-274
- [2] Anderson K and McOwan P W (2006). A real-time automated system for recognition of human facial expressions. *IEEE Trans. Systems, Man, and Cybernetics- Part B*, Vol. 36, No. 1, 96-105
- [3] Ang J, Dhillon R, Krupski A, et al. (2002). Prosody-based automatic detection of annoyance and frustration in human-computer dialog, *ICSLP*.
- [4] Ashraf, A.B., Lucey, S., Cohn, J.F., Chen, T., Ambadar, Z., Prkachin, K., Solomon, P. and Theobald, B.J. (2007). The painful face: pain expression recognition using active appearance models. *Int'l Conf. Multimodal Interfaces*, 9-14
- [5] Athanaselis T, Bakamidis S, Dologlou I, Cowie R, Douglas-Cowie E, Cox C (2005). ASR for emotional speech: Clarifying the issues and enhancing performance. *Neural Networks*, 18:437-444
- [6] Austermann, A. Esau, N. Kleinjohann, L. Kleinjohann, B. (2005). Prosody based emotion recognition for MEXI. *Int. Conf. Intelligent Robots and Systems*, 1138-1144
- [7] Balomenos, T., Raouzaoui, A., Ioannou, S., Drosopoulos, A., Karpouzis, K., Kollias, S. (2005). Emotion Analysis in Man-Machine Interaction Systems. *Lecture Notes in Computer Science*, vol. 3361, 318-328
- [8] Banse, R., Scherer, K.R. (1996). Acoustic profiles in Vocal emotion expression. *Journal Personality Social Psychology*, Vol. 70, No. 3, 614-636
- [9] Bartlett M S, Littlewort G, Braathen P, Sejnowski T J and Movellan J R (2003). A prototype for automatic recognition of spontaneous facial actions. *Advances in Neural Information Processing Systems*, Vol. 15, 1271-1278
- [10] Bartlett, M.S., Littlewort, G., Frank, M., Lainscsek, C., Fasel, I., and Movellan, J.(2005), Recognizing Facial Expression: Machine Learning and Application to Spontaneous Behavior, *IEEE International Conference on Computer Vision and Pattern Recognition*, 568-573
- [11] Bartlett M S, Littlewort G, Frank M G, Lainscsek C, Fasel I and Movellan J (2006). Fully automatic facial action recognition in spontaneous behavior. *Int. Conf. on Automatic Face and Gesture Recognition*, 223-230
- [12] Batliner A, Fischer K, Hubera R, Spilker J and Noth E. (2003). How to find trouble in communication. *Speech Communication*, Vol. 40, 117-143.
- [13] Batliner A, Hacker C, Steidl S, Noth E, D'Arcy S, et al. (2004). You stupid tin box—Children interacting with the AIBO robot: a cross-linguistic emotional speech. *Proceedings LREC*.
- [14] Blouin, C., and Maffiolo, V. (2005), "A study on the automatic detection and characterization of emotion in a voice service context", *Interspeech*, Lisbon, 469-472.
- [15] Burger S, MacLaren V and Yu H (2002). The ISL Meeting Corpus: The Impact of Meeting Type on Speech Style. In *Proceedings ICSLP*, Denver CO, USA.
- [16] Busso, C., Deng, Z., Yildirim, S., Bulut, M., Lee, C.M. et al. (2004), Analysis of Emotion Recognition using Facial Expressions, Speech and Multimodal Information. *Int. Conf. Multimodal Interfaces*. 205-211
- [17] Caridakis, G., Malatesta, L., Kessous, L., Amir, N., Paouzaoui, A. and Karpouzis, K.. (2006). Modeling Naturalistic Affective States via Facial and Vocal Expression Recognition. *Int. Conf. on Multimodal Interfaces*. 146-154
- [18] Chang Y, Hu C, Turk, M (2004). Probabilistic expression analysis on manifolds. *Proc. Computer Vision and Pattern Recognition*, 2:520-527
- [19] Chang Y, Hu C, Feris R and Turk M (2006). Manifold based analysis of facial expression. *J. Image and Vision Computing*, Vol. 24, No.6, 605-614
- [20] Chang Y, Vieira M, Turk M, and Velho L (2005). Automatic 3D facial expression analysis in videos. *Analysis and Modelling of Faces and Gestures, Proceedings*. 3723, pp. 293-307.
- [21] Chen, L.S (2000), Joint Processing of Audio-Visual Information for the Recognition of Emotional Expressions in Human-Computer Interaction, PhD thesis, UIUC
- [22] Chen, L, Huang, T. S., Miyasato, T., and Nakatsu, R. (1998). Multimodal human emotion/expression recognition. *Int. Conf. on Automatic Face and Gesture Recognition*. 396-401
- [23] Coan, J.A., Allen, J.J.B. (2007). *Handbook of Emotion Elicitation and Assessment*. Oxford University Press, New York, USA
- [24] Cohen, I., Cozman, F., Sebe, N., Cirelo, M., and Huang, T. S. (2004). Semi-supervised learning of classifiers: theory, algorithms, and their applications to human-computer interaction. *IEEE Trans on Pattern Analysis and Machine Intelligence*, Vol. 26, No. 12, 1553-1567
- [25] Cohen, L., Sebe, N., Garg, A., Chen, L., and Huang, T. (2003). Facial expression recognition from video sequences: Temporal and static modeling, *Computer Vision and Image Understanding*, 91(1-2):160-187
- [26] Cohn, J.F. (2006), Foundations of Human Computing: Facial Expression and Emotion, *Int. Conf. on Multimodal Interfaces*, 233-238
- [27] Cohn JF, Reed LI, Ambadar Z, Xiao J, and Moriyama T. (2004). Automatic Analysis and recognition of brow actions and head motion in spontaneous facial behavior. *Int. Conf. on Systems, Man & Cybernetics*, 1, 610-616
- [28] Cohn, J.F. and Schmidt, K.L.(2004). The timing of Facial Motion in Posed and Spontaneous Smiles, *International Journal of Wavelets, Multiresolution and Information Processing*, 2, 1-12
- [29] Cohn JF and Tronick EZ. (1988). Mother Infant Interaction: the sequence of dyadic states at three, six and nine months. *Development Psychology*, 23, 68-77
- [30] Cowie, R., Douglas-Cowie, E., Savvidou, S., McMahon, E., Sawey, M., & Schröder, M. (2000). 'Feeltrace': an instrument for recording perceived emotion in real time. *Proceedings of the ISCA Workshop on Speech and Emotion*, 19-24
- [31] Cowie, R., Douglas-Cowie, E., Tsapatsoulis, N., Votsis, G., Kollias, S., Fellenz, W., and Taylor, J.G. (2001). Emotion Recognition in Human-Computer Interaction, *IEEE Signal Processing Magazine*, January, 32-80
- [32] Cowie R, Douglas-Cowie E and Cox C (2005). Beyond emotion archetypes: databases for emotion modeling using neural networks. *Neural Networks*, 18: 371-388
- [33] Dellaert, F., Polzin, T. and Waibel, A. (1996). Recognizing emotion in speech. *Int. Conf. on Spoken Language Processing*, 1970-1973
- [34] Devillers L and Vasilescu I (2004). Reliability of lexical and prosodic cues in two real-life spoken dialog corpora. *Proceedings LREC*.
- [35] Devillers L, Vasilescu I. (2006). Real-Life Emotions Detection with Lexical and Paralinguistic Cues on Human-Human Call Center Dialogs. *Int. Conf. on Spoken Language Processing*
- [36] Vasilescu, I. and Devillers, L. (2005). Detection of real-life emotions in call centers. *Interspeech*.

- [37] Devillers L, Vidrascu L, and Lamel L (2005). Challenges in real-life emotion annotation and machine learning based detection. *Neural Networks*, 18: 407-422
- [38] Douglas-Cowie E, Campbell N, Cowie R and Roach P (2003). Emotional Speech: towards a new generation of database. *Speech Communication*, 40(1-2): 33-60
- [39] Duric, Z., Gray, W.D., Heishman, R., Li, F., Rosenfeld, A., Schoelles, M.J., Schunn, C., Wechsler, H. (2002). Integrating perceptual and cognitive modeling for adaptive and intelligent human-computer interaction. *Proceedings of the IEEE*, Vol. 90, No. 7, 1272-1289
- [40] Ekman P. (1972). Universals and cultural differences in facial expressions of emotion. *Nebr. Symp. Motiv.* 1971, 207-283
- [41] Ekman, P., editor (1982). *Emotion in the human face*. Cambridge University Press, New York, 2nd edition
- [42] Ekman, P. (1994). Strong Evidence for Universals in Facial Expressions: A Reply to Russell's Mistaken Critique, *Psychological Bulletin*, 115(2): 268-287
- [43] Ekman, P., Friesen, W.V., Hager, J.C. (2002). *Facial Action Coding System. A Human Face*, Salt Lake City, USA
- [44] Ekman, P., Huang, T.S., Sejnowski, T.J. and Hager, J.C., (Eds.), (1993). *NSF Understanding the Face, A Human Face eStore*, Salt Lake City, USA, (see Library).
- [45] Ekman P, Matsumoto D, and Friesen WV. (2005). Facial Expression in Affective Disorders. In *What the Face Reveals*. Edited by Ekman P and Rosenberg EL. 429-439
- [46] Ekman P. and Oster H. (1979). Facial expressions of emotion. *Ann. Rev. Psychol.* 1979, 30:527-554
- [47] Ekman P. and Rosenberg E.L. (2005). *What the face reveals: basic and applied studies of spontaneous expression using the facial action coding system*. 2nd edition, Oxford University Press.
- [48] El Kaliouby R and Robinson P (2004). Real-time Inference of complex mental states from facial expression and head gestures. *Computer Vision and Pattern Recognition Workshop*, Vol. 3, 154
- [49] Fasel, B. and Luttin, J. (2003). Automatic facial expression analysis: Survey. *Pattern Recognition*, 36(1): 259-275
- [50] Fasel B, Monay F and Gatica-Perez D (2004). Latent semantic analysis of facial action codes for automatic facial expression recognition. *ACM Int. Workshop on Multimedia Information Retrieval*, 181-188
- [51] Fiechter, C-N. (1994). Efficient reinforcement learning. *ACM Conf. on Computational Learning Theory*, 88-97
- [52] Forbes-Riley K and Litman D (2004). Predicting emotion in spoken dialogue from multiple knowledge sources. *Proc. Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics (HLT/NAACL)*
- [53] Fragopanagos, F. and Taylor, J.G. (2005). Emotion recognition in human-computer interaction, *Neural Networks*, 18: 389-405
- [54] Fried E. (1976). *The impact of nonverbal communication of facial affect on children's learning*. PhD thesis, Rutgers University, New Brunswick, NJ
- [55] Furnas, G., Landauer, T., Gomes, L., and Dumais, S. (1987). The vocabulary problem in human-system communication, *Communications of the ACM*, Vol. 30, No. 11, 964-972.
- [56] Go HJ, Kwak KC, Lee DJ, and Chun MG. (2003). Emotion recognition from facial image and speech signal. *Int. Conf. of the Society of Instrument and Control Engineers*. 2890-2895
- [57] Graciarena, M., Shriberg, E., Stolcke, A., Enos, F., Hirschberg, J., and Kajarekar, S. (2006). Combining prosodic, lexical and cepstral systems for deceptive speech detection. *Proc. Int. Conf. on Acoustics, Speech and Signal Processing*, I:1033-1036
- [58] Greenwald M, Cook E and Lang P. (1989). Affective judgment and psychophysiological response: dimensional covariation in the evaluation of pictorial stimuli. *J. Psychophysiol.* 3:51-64
- [59] Gross, R. (2005) Face databases. In *Handbook of Face Recognition*, Li S Z., Jain A.K., (Eds.), Springer, New York, USA, 301-328
- [60] Gu H, Ji Q (2004). An Automated Face Reader for Fatigue Detection. *Int. Conf. Automatic Face and Gesture Recognition*. 111-116
- [61] Gunes, H., Piccardi, M. (2005). Affect Recognition from Face and Body: Early Fusion vs. Late Fusion, In *Proc. Int'l Conf. Systems, Man and Cybernetics*, 3437- 3443
- [62] Gunes, H. and Piccardi, M. (2005). Fusing Face and Body Display for Bi-Modal Emotion Recognition: Single Frame Analysis and Multi-Frame Post Integration. *Int. Conf. on Affective Computing and Intelligent Interaction*, 102 – 111
- [63] Gunes, H. and Piccardi, M. (2006). A bimodal face and body gesture database for automatic analysis of human nonverbal affective behavior. *Int. Conf. on Pattern Recognition*, Vol. 1, 1148-1153
- [64] Guo G and Dyer C R (2005). Learning from examples in the small sample case – face expression recognition. *IEEE Trans. Systems, Man and Cybernetics – Part B*, Vol.35, No.3, 477-488
- [65] Hirschberg, J., Benus, S., Brenier, J.M., Enos, F., Friedman, S. (2005). Distinguishing Deceptive from Non-Deceptive Speech. *Interspeech*, 1833-1836
- [66] Hoch, S., Althoff, F., McGlaun, G., Rigoll, G. (2005). Bimodal fusion of emotional data in an automotive environment, *ICASSP*, Vol. II, 1085-1088, 2005
- [67] Ioannou, S., Raouzaoui, A., Tzouvaras, V., Mailis, T., Karpouzis, K., & Kollias, S. (2005). Emotion recognition through facial expression analysis based on a neurofuzzy method. *Neural Networks*: 18, 423-435.
- [68] Jaimes, A. and Sebe, N. (2005). Multimodal human computer interaction: a survey. *Workshop on Human Computer Interaction in conjunction with ICCV*.
- [69] Ji Q, Lan P and Looney C (2006). A probabilistic framework for modeling and real-time monitoring human fatigue. *IEEE SMC-Part A*, Vol. 36, No.5, 862-875
- [70] Juslin, P.N., Scherer, K.R. (2005). Vocal expression of affect. In *The New Handbook of Methods in Nonverbal Behavior Research*. Harrigan, J., Rosenthal, R., Scherer, K., Eds. Oxford University Press, Oxford, UK
- [71] Kanade, T., Cohn, J., and Tian, Y. (2000). Comprehensive Database for Facial Expression Analysis, In *Proceeding of International Conference on Face and Gesture Recognition*, 46-53
- [72] Kapoor, A., Burleson, W., and Picard, R. W. (2007). Automatic prediction of frustration. *Int. Journal of Human-Computer Studies*. Vol. 65(8), 724-736.
- [73] Kapoor, A. and Picard, R. W. (2005). Multimodal affect recognition in learning environment. *ACM Int'l Conf. on Multimedia*, 677-682
- [74] Karpouzis, K., Caridakis, G., Kessous, L., Amir, N., Raouzaoui, A., Malatesta, L., and Kollias, S. (2007). Modeling naturalistic affective states via facial, vocal, and bodily expression recognition, *Lecture Notes in Artificial Intelligence*, vol. 4451, 91-112.
- [75] Keltner D (1995). Signs of appeasement: evidence for the distinct displays of embarrassment, amusement and shame. *Journal of Personality and Social Psychology*, 68(3). 441-454
- [76] Kobayashi, H., and Hara, F. (1991). The recognition of basic facial expressions by neural network. *Proc. Int'l Joint Conf. on Neural Networks*, 460-466.
- [77] Kotsia, I. and Pitas, I. (2007). Facial expression recognition in image sequences using geometric deformation features and support vector machine. *IEEE Trans. On Image Processing*, 16(1): 172-187
- [78] Kuncheva, L.I. (2004). *Combining Pattern Classifier: Methods and Algorithms*, John Wiley and Sons, 2004
- [79] Kwon, O.W., Chan, K., Hao, J., Lee, T.W (2003). Emotion Recognition by Speech Signals, *EUROSPEECH*.
- [80] Laskowski, K. and Burger, S. (2006). Annotation and Analysis of Emotionally Relevant Behavior in the ISL Meeting Corpus, *LREC*, Genoa, Italy.
- [81] Lee, C. and Elgammal, A. (2005). Facial expression analysis using nonlinear decomposable generative models. In *IEEE International Workshop on Analysis and Modeling of Faces and Gestures*

- [82] Lee C and Narayanan (2003). Emotion recognition using a data-driven fuzzy inference system. In Proc. Eurospeech, 157-160
- [83] Lee C M Narayanan, S.S. (2005). Toward detecting emotions in spoken dialogs. *IEEE Tran. Speech and Audio Processing*, Vol. 13(2): 293-303
- [84] Liscombe, J., Hirschberg, J., Venditti, J.J. (2005). Detecting Certainty in Spoken Tutorial Dialogues. *Interspeech*.
- [85] Lisetti, C.L., Nasoz, F. (2002). MAUI: A multimodal affective user interface. *Proc. Int'l Conf. Multimedia*, 161-170
- [86] Litman, D.J. and Forbes-Riley, K. (2004). Predicting Student Emotions in Computer-Human Tutoring Dialogues. In Proc. of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL), July
- [87] Littlewort, G.C., Bartlett, M.S. and Lee, K. (2007). Faces of pain: Automated measurement of spontaneous facial expressions of genuine and posed pain. *Int'l Conf. Multimodal Interfaces*, 15-21
- [88] Lucey, S., Ashraf, A.B., and Cohn, J.F. (2007). Investigating Spontaneous Facial Action Recognition through AAM Representations of the Face. In *Face Recognition*, Delac, K. and Grgic, M., Eds. Vienna, Austria: I-Tech Education and Publishing, 275-286
- [89] Maat, L., and Pantic, M. (2006). Gaze-X: Adaptive affective multimodal interface for single-user office scenarios, *Proc. ACM Int'l Conf. Multimodal Interfaces*, 171-178
- [90] Mase, K. (1991). Recognition of facial expression from optical flow. *IEICE Trans. E74(10)*. 3474-3483
- [91] Matos, S., Birring, S.S., Pavord, I.D. and Evans, D.H. (2006). Detection of cough signals in continuous audio recordings using HMM. *IEEE Trans. Biomedical Engineering*, Vol. 53, No. 6, 1078-1083.
- [92] Mehrabian, A. (1968). Communication with words. *Psychology Today*, 2(4): 53-56
- [93] Mozziconacci, S. (2002). Prosody and Emotions. *Proc. Speech Prosody Aix-en-Provence*, 1-9.
- [94] Neiberg D, Elenius K, and Laskowski K. (2006). Emotion Recognition in Spontaneous Speech Using GMM. *Int. Conf. on Spoken Language Processing*, 809-812
- [95] O'Toole A J, Harms J, Snow S L, Hurst D R, Pappas M R, et al. (2006). A Video Database of Moving Faces and People. *IEEE PAMI*, VOL. 27, NO. 5, MAY 2005, 812-816
- [96] Oudeyer, P-Y. (2003). The production and recognition of emotions in speech: features and algorithms. *Int'l Journal Human-Computer Studies*, 59, 157-183.
- [97] Pal P, Iyer A N and Yantorno R E (2006). Emotion detection from infant facial expressions and cries. In Proc. Int'l Conf. Acoustics, Speech & Signal Processing, 2, pp. 721-724, 2006.
- [98] Pantic, M., and Bartlett, M.S. (2007). Machine analysis of facial expressions. In *Face Recognition*, Delac, K. and Grgic, M., Eds. Vienna, Austria: I-Tech Education and Publishing, 377-416
- [99] Pantic, M., and Patras, I. (2006). Dynamics of facial expression: recognition of facial actions and their temporal segments from face profile image sequences. *IEEE Trans. Systems, Man and Cybernetics - Part B*, Vol. 36, No.2, 433-449
- [100] Pantic, M., Pentland, A., Nijholt, A., and Huang, T.S. (2006). Human Computing and Machine Understanding of Human Behavior: A Survey, *Int. Conf. on Multimodal Interfaces*, 239-248
- [101] Pantic, M., and Rothkrantz, L.J.M. (2000). Automatic analysis of facial expressions—the state of the art. *IEEE PAMI*, Vol.22, No.12, 1424-1445
- [102] Pantic M., and Rothkrantz, L.J.M. (2003). Toward an affect-sensitive multimodal human-computer interaction, *Proceedings of the IEEE*, Vol. 91, No. 9, Sept., 1370-1390
- [103] Pantic, M., and Rothkrantz, L.J.M. (2004). Facial action recognition for facial expression analysis from static face images. *IEEE Trans. On Systems, Man and Cybernetics-Part B*, Vol. 34, No. 3, 1449-1461
- [104] Pantic, M., and Rothkrantz, L.J.M. (2004). Case-based reasoning for user-profiled recognition of emotions from face images. *Int. Conf. Multimedia & Expo*, 391-394
- [105] Pantic, M., Sebe, N., Cohn, J.F. and Huang, T. (2005), *Affective Multimodal Human-Computer Interaction*, in Proc. ACM Int'l Conf. on Multimedia, 669-676
- [106] Pantic, M., Valstar, M.F, Rademaker, R. and Maat, L. (2005), Web-based database for facial expression analysis, *Int. Conf. on Multimedia and Expo*, 317-321
- [107] Pentland, A. (2005). Socially aware, computation and communication, *IEEE Computer*, Vol.38, 33-40
- [108] Petridis, S. and Pantic, M. (2008). Audiovisual discrimination between laughter and speech. *Int'l Conf. Acoustics, Speech, and Signal Processing*.
- [109] Picard, R.W. (1997). *Affective Computing*, MIT Press, Cambridge.
- [110] Plutchik R. (1980). *Emotion: A psychoevolutionary synthesis*. New York: Harper and Row.
- [111] Roisman, G.I., Tsai, J.L., Chiang, K.S.(2004). The Emotional Integration of Childhood Experience: Physiological, Facial Expressive, and Self-reported Emotional Response During the Adult Attachment Interview, *Developmental Psychology*, Vol. 40, No. 5, 776-789
- [112] Roli, F., Kittler, J., et al., eds., (2001-2005). *Int. Workshop Multiple Classifier Systems (MCS)*.
- [113] Russell J.A., Bachorowski J. and Fernandez-Dols J. (2003). Facial and vocal expressions of emotion. *Ann. Rev. Psychol.* 54:329-349
- [114] Russell J and Mehrabian A. (1977). Evidence for a three-factor theory of emotions. *J. Res. Personality*, 11: 273-294
- [115] Samal A and Iyengar P A (1992). Automatic recognition and analysis of human faces and facial expressions: A survey. *Pattern Recognition*, Vol. 25, No.1, 65-77
- [116] Sander D. Grandjean D. and Scherer K.R. (2005). A system approach to appraisal mechanisms in emotion. *Neural Networks*. 18: 317-352
- [117] Scherer K.R. (1999). Appraisal theory. In Dalglish T and Power M J (Eds.), *Handbook of cognition and emotion*, New York: Wiley, 637-663
- [118] Schuller, B., Muller, R., Hornler, B., Hothker, A., Konosu, H. and Rigoll, G. (2007). Audiovisual recognition of spontaneous Interest within conversations. *Int. Conf. on Multimodal Interfaces*, 30-37
- [119] Schuller, B., Stadermann, J., Rigoll, G. (2006). Affect-Robust Speech Recognition by Dynamic Emotional Adaptation, *Proc. Speech Prosody 2006, Special Session on Prosody in Automatic Speech Recognition*.
- [120] Schuller, B., Villar, R. J., Rigoll, G., Lang, M. (2005). Meta-Classifiers in acoustic and linguistic feature fusion-based affect recognition. *Int. Conf. on Acoustics, Speech, and Signal Processing*, 325-328
- [121] Sebe, N., Cohen, I., and Huang, T.S. (2005). *Multimodal Emotion Recognition*, *Handbook of Pattern Recognition and Computer Vision*, World Scientific, 2005.
- [122] Sebe, N., Cohen, I., Gevers, T. and Huang, T.S. (2006). Emotion recognition based on joint visual and audio cues. *Int. Conf. on Pattern Recognition*, 1136-1139
- [123] Sebe, N., Lew, M.S., Cohen, I., Sun, Y., Gevers, T., Huang, T.S.(2004), *Authentic Facial Expression Analysis*, *Int. Conf. on Automatic Face and Gesture Recognition*
- [124] Song, M., Bu, J., Chen, C., and Li, N. (2004), Audio-visual based emotion recognition—A new approach, *Int. Conf. Computer Vision and Pattern Recognition*. 2004, 1020-1025
- [125] Steidl, S., Levit, M., Batliner, A, Noth, E., and Niemann, H. (2005), "Off all things the measure is man" Automatic classification of emotions and inter-labeler consistency, *ICASSP*, vol.1, 317-320
- [126] Stein, B., Meredith, M.A. (1993). *The Merging of Senses*. MIT Press, Cambridge, USA
- [127] Suwa, M., Sugie, N., Fujimora, K. (1978). A preliminary note on pattern recognition of human emotional expression. *Int. Joint Conference on Pattern Recognition*. 408-410
- [128] Tao, H. and Huang, T.S. (1999), Explanation-based facial motion tracking using a piecewise Bezier volume deformation mode, *IEEE CVPR*, vol.1, pp. 611-617,

- [129] Teeters, A., Kaliouby, R. E., and Picard, R. W. (2006). Self-Cam: Feedback From What Would Be Your Social Partner. ACM SIGGRAPH, Research Posters, p. 138
- [130] Tian Y L, Kanade T and Cohn J F (2005). Facial expression analysis. In: Handbook of Face Recognition, Li S Z and Jain A K (Eds.), Springer, New York, USA, 247-276
- [131] Tomkins SS. (1962). Affect, Imagery, Consciousness, Vol. 1. New York: Springer
- [132] Tong, Y., Liao, W. and Ji, Q. (2007). Facial action unit recognition by exploiting their dynamics and semantic relationships. IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 29, no. 10, 1683-1699
- [133] Truong K P and van Leeuwen D A (2007). Automatic discrimination between laughter and speech. Speech Communication, 49: 144-158.
- [134] Valstar, M.F., Gunes, H. and Pantic, M. (2007). How to distinguish posed from spontaneous smiles using geometric features. Int'l Conf. Multimodal Interfaces, 38-45.
- [135] Valstar, M., Pantic, M., Ambadar, Z., and Cohn, J.F. (2006). Spontaneous vs. Posed Facial Behavior: Automatic Analysis of Brow Actions. Int. Conf. on Multimedia Interfaces. 162-170
- [136] Valstar, M., Pantic, M., and Patras, I. (2004). Motion history for facial action detection from face video. Int. Conf. Systems, Man and Cybernetics, Vol.1, 635-640
- [137] Wang, H. and Aluja, N. (2003). Facial expression decomposition. IEEE International Conference on Computer Vision, p.958
- [138] Wang, Y. and Guan, L.(2005), Recognizing human emotion from audiovisual information, ICASSP, Vol. II, 1125-1128
- [139] Wang, J., Yin, L., Wei, X., and Sun, Y. (2006). 3D Facial Expression Recognition Based on Primitive Surface Feature Distribution. IEEE Conference on Computer Vision and Pattern Recognition, 2:1399-1406
- [140] Watson, D., Clark, L. A., and Tellegen, A. (1988). Development and validation of brief measures of positive and negative affect: The PANAS scales. Journal of Personality and Social Psychology, Vol. 54, 1063-1070.
- [141] Wen, Z. and Huang, T.S. (2003). Capturing Subtle Facial Motions in 3D Face Tracking. Int. Conf. on Computer Vision, 1343-1350
- [142] Whissell C M (1989). The dictionary of affect in language. In Plutchik R. and Kellerman H (Eds.). Emotion: Theory, research and experience. The measurement of emotions, Vol.4. 113-131. New York: Academic Press
- [143] Whitehill J. and Omlin, C. W. (2006). Haar features for FACS AU recognition. Int. Conf. on Automatic Face and Gesture Recognition, 217-222
- [144] Williams, A. C. de C. (2002) Facial expression of pain: An evolutionary account. Behavioral and Brain Sciences, Vol. 25, No. 4, 439-488
- [145] Williams, C. and Stevens, K. (1972). Emotions and speech: Some acoustic correlates. Journal of the Acoustic Society of America, 52(4), 1238-1250
- [146] Xiao J, Moriyama T, Kanade T and Cohn J F (2003). Robust full-motion recovery of head by dynamic templates and re-registration techniques. Int. J. Imaging Systems and Technology, Vol. 13, No.1, 85-94
- [147] Yeasin M., Bullot B. and Sharma R. (2006). Recognition of facial expressions and measurement of levels of interest from video, IEEE Trans. On Multimedia, Vol.8, No. 3, June, 500-507
- [148] Yin L, Wei X, Sun Y, Wang J, Rosato M J (2006). A 3D facial expression database for facial behavior research. Int. Conf. on Automatic Face and Gesture Recognition, 211-216
- [149] Zeng, Z., Fu, Y., Roisman, G.I., Wen, Z., Hu, Y., and Huang, T.S. (2006). Spontaneous Emotional Facial Expression Detection. Journal of Multimedia, 1(5): 1-8.
- [150] Zeng, Z., Hu, Y., Liu, M., Fu, Y. and Huang, T.S.(2006), Training Combination Strategy of Multi-stream Fused Hidden Markov Model for Audio-visual Affect Recognition, in Proc. ACM Int'l Conf. on Multimedia, 2006, 65-68
- [151] Zeng, Z., Hu, Y., Roisman, G.I., Wen, Z., Fu, Y., and Huang, T.S. (2007), Audio-visual Spontaneous Emotion Recognition, In Artificial Intelligence for Human Computing, Eds: Huang T.S., Nijholt A., Pantic M., and Pentland A., Springer, 72-90.
- [152] Zeng, Z., Pantic, M., Roisman, G.I. and Huang, T.S. (2007). A Survey of Affect Recognition Methods: Audio, Visual, and Spontaneous Expressions. Int'l Conf. Multimodal Interfaces, 126-133
- [153] Zeng, Z., Tu, J., Liu, M., Zhang, T., Rizzolo, N., Zhang, Z., Huang, T.S., Roth, D., and Levinson, S. (2004), Bimodal HCI-related Emotion Recognition, Int. Conf. on Multimodal Interfaces, 137-143.
- [154] Zeng, Z., Tu, J., Pianfetti, P., Liu, M., Zhang, T., Zhang Z., Huang T S and Levinson S (2005), Audio-visual Affect Recognition through Multi-stream Fused HMM for HCI, Int. Conf. Computer Vision and Pattern Recognition. 967-972
- [155] Zeng, Z., Tu, J., Liu, M., Huang, T.S., Pianfetti, B., Roth D. and Levinson, S. (2007), Audio-visual Affect Recognition, IEEE Transactions on Multimedia, Vol. 9, No. 2, February, 424-428
- [156] Zeng, Z., Zhang, Z., Pianfetti, B., Tu, J., and Huang, T.S. (2005), Audio-visual Affect Recognition in Activation-evaluation Space, Int. Conf. on Multimedia & Expo, 828-831.
- [157] Zhang T, Hasegawa-Johnson M and Levinson S E (2004). Children's Emotion Recognition in an Intelligent Tutoring Scenario, Interspeech 2004.
- [158] Zhang Y and Ji Q (2005). Active and Dynamic Information Fusion for Facial Expression Understanding from Image Sequences. IEEE Trans. Pattern Anal. Mach. Intell. 27(5): 699-714
- [159] Zhou, Z.-H., Chen, K.-J. and Dai, H.-B., Enhancing relevance feedback in image retrieval using unlabeled data. ACM Transactions on Information Systems, 2006, 24(2): 219-244
- [160] Zhu, Z. and Ji, Q. (2006). Robust Real-Time Face Pose and Facial Expression Recovery. IEEE Conference on Computer Vision and Pattern Recognition, 1: 681-688

Zhihong Zeng received his PhD in Institute of Automation, Chinese Academy of Sciences in 2002. He is currently Beckman Postdoctoral Fellow at Beckman Institute, UIUC. His research interests include multimodal affective computing, multimodal human computer interaction and computer vision. He is an IEEE member.

Maja Pantic received the MSc and PhD degrees in Computer Science from Delft University of Technology, The Netherlands, in 1997 and 2001. She is Reader in Multimodal HCI at Imperial College London, Computing Department, and Professor in Affective and Behavioural Computing at the University of Twente, Computer Science Department. Her research interests include computer vision and machine learning applied to face and body gesture recognition, multimodal human-computer interaction (HCI), context-sensitive HCI, affective computing, and e-learning tools. She is an IEEE senior member. She is an Associate Editor of IEEE Trans. on Systems, Man and Cybernetics - Part B, and of Image and Vision Computing Journal. She is a guest editor, organizer and committee member of over 10 major journals and conferences.

Glenn I. Roisman received his PhD from the University of Minnesota in 2002. He is currently assistant professor in the Department of Psychology at UIUC. His research interests concern social and emotional development across the lifespan. He has published over twenty-five scholarly journal articles and chapters, and received the Society for Research in Child Development's Award for Early Research Contributions in 2007.

Thomas S. Huang received his Sc.D. from MIT in 1963. He is William L. Everitt Distinguished Professor of Electrical and Computer Engineering, and Co-chair of Human Computer Intelligent Interaction Initiative (HCII) in the Beckman Institute, UIUC. His professional interests are computer vision, image compression and enhancement, pattern recognition, and multimodal signal processing. He has more than 80 honors, awards and outstanding achievements, including Member of National Academy of Engineering; Fellow of IEEE; Foreign Member of Chinese Academy of Sciences.