

Machine Analysis of Facial Behaviour: Naturalistic & Dynamic Behaviour

Journal:	<i>Philosophical Transactions B</i>
Manuscript ID:	RSTB-2009-0135
Article Type:	Research
Date Submitted by the Author:	27-Jun-2009
Complete List of Authors:	Pantic, Maja; Imperial College London, Computing Department; University of Twente, EEMCS
Issue Code: Click here to find the code for your issue.:	COMPUTATION
Subject:	Behaviour < BIOLOGY
Keywords:	facial expressions, dynamic behaviour, naturalistic behaviour

Machine Analysis of Facial Behaviour: Naturalistic & Dynamic Behaviour

Maja Pantic^{1,2}

¹ Imperial College London, Computing, London SW7 2AZ, UK

² University of Twente, EEMCS, 7500 AE Enschede, The Netherlands
m.pantic@imperial.ac.uk

Abstract

This article introduces recent advances in machine analysis of facial expressions. It describes the problem space, surveys the problem domain, and examines the state of the art. Two recent research topics are discussed with particular attention: analysis of facial dynamics and analysis of naturalistic (spontaneously displayed) facial behaviour. Scientific and engineering challenges in the field in general, and in these specific subproblem areas in particular, are discussed and recommendations for accomplishing a better facial expression measurement technology are outlined.

Keywords

Automatic facial expression analysis, dynamics of facial behaviour, naturalistic behaviour

1. Introduction

A widely accepted prediction is that computing will move to the background, weaving itself into the fabric of our everyday living and projecting the human user into the foreground. To realize this goal, next-generation computing (a.k.a. pervasive computing, ambient intelligence, and human computing) will need to develop human-centred user interfaces that respond readily to naturally occurring, multimodal, human communication [1]. These interfaces will need the capacity to perceive and understand intentions and emotions as communicated by social and affective signals. Motivated by this vision of the future, automated analysis of nonverbal behaviour, and especially of facial behaviour, has attracted increasing attention in computer vision, pattern recognition, and human-computer interaction [2], [3], [4], [5]. To wit, facial expression is one of the most cogent, naturally preeminent means for human beings to communicate emotions, to clarify and stress what is said, to signal comprehension, disagreement, and intentions, in brief, to regulate interactions with the environment and other persons in the vicinity [6], [7]. Automatic analysis of facial expressions forms, therefore, the essence of numerous next-generation-computing tools including affective computing technologies (proactive and affective user interfaces), learner-adaptive tutoring systems, patient-profiled personal wellness technologies, etc.

This article introduces recent advances in machine analysis of facial expressions. It describes the problem space, surveys the problem domain, and examines the state of the art. Two recent research topics will receive particular attention: analysis of facial dynamics, and analysis of naturalistic (spontaneously displayed) facial behaviour. Scientific and engineering challenges in the field in general, and in these specific subproblem areas in particular, will be discussed and recommendations for accomplishing a better facial expression measurement technology will be outlined.

2. The Process of Automatic Facial Behaviour Analysis

1
2
3 Facial expression recognition is a process performed by humans or computers, which consists of three
4 steps (Fig. 1):

- 5 1) locating faces in the scene (e.g., in an image; this step is also referred to as *face detection*),
6 2) extracting facial features from the detected face region (e.g., detecting the shape of facial components
7 or describing the texture of the skin in a facial area; this step is referred to as *facial feature extraction*), and
8 3) analysing the motion of facial features and/or the changes in the appearance of facial features and
9 classifying this information into some facial-expression-interpretative categories such as facial muscle
10 activations like smile or frown, emotion (affect) categories like happiness or anger, attitude categories like
11 (dis)liking or ambivalence, etc. (this step is also referred to as *facial expression interpretation*).
12
13

14 **Fig. 1.** Outline of an automated, geometric-features-based system for facial
15 expression recognition (for details of this system, see [55]).
16

17 The problem of *finding faces* can be viewed as a segmentation problem (in machine vision) or as a
18 detection problem (in pattern recognition). It refers to identification of all regions in the scene that contain
19 a human face. The problem of finding faces (*face localization, face detection*) should be solved regardless of
20 clutter, occlusions, and variations in head pose and lighting conditions. The presence of non-rigid
21 movements due to facial expression and a high degree of variability in facial size, colour and texture
22 make this problem even more difficult. Numerous techniques have been developed for face detection in
23 still images [8], [9]. However, most of them can detect only upright faces in frontal or near-frontal view.
24 Arguably the most commonly employed face detector in automatic facial expression analysis is the real-
25 time face detector proposed by Viola and Jones [10].
26

27 The problem of feature extraction can be viewed as a dimensionality reduction problem (in machine
28 vision and pattern recognition). It refers to transforming the input data into a reduced representation set
29 of features, which encode the relevant information from the input data. The problem of *facial feature*
30 *extraction* from input images may be divided into at least three dimensions [2], [3]: (a) Are the features
31 holistic (spanning the whole face) or atomistic (spanning subparts of the face)?; (b) Is temporal
32 information used?; (c) Are the features view- or volume based (2-D/3-D)? Given this glossary, most of
33 the proposed approaches to facial expression recognition are directed toward static, analytic, 2-D facial
34 feature extraction [3], [4]. The usually extracted facial features are either *geometric features* such as the
35 shapes of the facial components (eyes, mouth, etc.) and the locations of facial fiducial points (corners of
36 the eyes, mouth, etc.) or *appearance features* representing the texture of the facial skin in specific facial
37 areas including wrinkles, bulges, and furrows. Appearance-based features include learned image filters
38 from independent component analysis (ICA), principal component analysis (PCA), local feature analysis
39 (LFA), Gabor filters, integral image filters (also known as box-filters and Haar-like filters), features based
40 on edge-oriented histograms, etc. Several efforts have been also reported which use both geometric and
41 appearance features (e.g., [11]). These approaches to automatic facial expression analysis are referred to
42 as *hybrid* methods. Although it has been reported that methods based on geometric features are often
43 outperformed by those based on appearance features using, e.g., Gabor wavelets or eigenfaces, recent
44 studies show that in some cases geometric features can outperform appearance-based ones [12], [3]. Yet, it
45 seems that using both geometric and appearance features might be the best choice in the case of certain
46 facial expressions [12], [13].
47

48 Contractions of facial muscles, which produce facial expressions, induce movements of the facial skin
49 and changes in the location and/or appearance of facial features (e.g., contraction of the Corrugator
50 muscle induces a frown and causes the eyebrows to move towards each other, usually producing
51 wrinkles between the eyebrows; Fig. 2). Such *changes can be detected* by analyzing optical flow, facial-
52 point- or facial-component-contour-tracking results, or by using an ensemble of classifiers trained to
53 make decisions about the presence of certain changes (e.g., whether the nasolabial furrow is deepened or
54 not) based on the passed appearance features (see also section 3). The optical flow approach to describing
55 face motion has the advantage of not requiring a facial feature extraction stage of processing. Dense flow
56 information is available throughout the entire facial area, regardless of the existence of facial components,
57 even in the areas of smooth texture such as the cheeks and the forehead. Because optical flow is the
58
59
60

1
2
3 visible result of movement and is expressed in terms of velocity, it can be used to represent directly facial
4 expressions. Many researchers adopted this approach (e.g., [15], [44]; for comprehensive overviews see
5 [14], [2], [4]). Until recently, standard optical flow techniques were arguably most commonly used for
6 tracking facial characteristic points and contours as well. In order to address the limitations inherent in
7 optical flow techniques such as the accumulation of error and the sensitivity to noise, occlusion, clutter,
8 and changes in illumination, recent efforts in automatic facial expression recognition use sequential state
9

10
11
12 **Fig. 2.** Facial appearance of the Corrugator muscle contraction (coded as AU4
13 in the FACS system, [18]).
14

15
16 estimation techniques (such as Kalman filter and Particle filter) to track facial feature points in image
17 sequences (e.g., [11], [12], [3]).

18 Eventually, dense flow information, tracked movements of facial characteristic points, tracked
19 changes in contours of facial components, and/or extracted appearance features are translated into a
20 description of the displayed facial expression. This description (*facial expression interpretation*) is usually
21 given either in terms of shown affective states (emotions) or in terms of activated facial muscles
22 underlying the displayed facial expression (see section 3 for a detailed discussion and an overview of the
23 state of the art). Most facial expressions analyzers developed so far target human facial affect analysis and
24 attempt to recognize a small set of prototypic emotional facial expressions like happiness and anger [2],
25 [5]. However, several promising prototype systems were reported that can recognize deliberately
26 produced AUs in face images and even few attempts towards recognition of spontaneously displayed
27 AUs have been recently reported as well [4], (see also section 4). While the older methods employ simple
28 approaches including expert rules and machine learning methods such as neural networks to classify the
29 relevant information from the input data into some facial-expression-interpretative categories [14], [2],
30 the more recent (and often more advanced) methods employ probabilistic, statistical, and ensemble
31 learning techniques, which seem to be particularly suitable for automatic facial expression recognition
32 from face image sequence [3], [4].
33
34

35 **3. Facial Behaviour Interpretation: Emotions, Social Signals, and AUs**

36 Two main streams in the current research on automatic analysis of facial expressions consider facial affect
37 (emotion) detection and facial muscle action (action unit, AU) detection [14], [2], [4]. These streams stem
38 directly from two major approaches to facial expression measurement in psychological research [16]:
39 message and sign judgment. The aim of message judgment is to infer what underlies a displayed facial
40 expression, such as affect or personality, while the aim of sign judgment is to describe the “surface” of the
41 shown behaviour, such as facial movement or facial component shape. Thus, a brow frown (Fig. 2) can be
42 judged as “anger” in a message-judgment and as a facial movement that lowers and pulls the eyebrows
43 closer together in a sign-judgment approach. While message judgment is all about interpretation, sign
44 judgment attempts to be objective, leaving inference about the conveyed message to higher order
45 decision making.
46

47 Most commonly used facial expression descriptors in message judgment approaches are the six basic
48 emotions (fear, sadness, happiness, anger, disgust, surprise; see Fig. 3), proposed by Ekman and discrete
49 emotion theorists [17], who suggest that these emotions are universally displayed and recognized from
50 facial expressions. This trend can also be found in the field of automatic facial expression analysis. Most
51 facial expressions analyzers developed so far target human facial affect analysis and attempt to recognize
52 a small set of prototypic emotional facial expressions like happiness and anger [4]. Automatic detection of
53 the six basic emotions in posed, controlled displays can be done with reasonably high accuracy. More
54 specifically, recent studies report recognition accuracy of above 90% for prototypic facial expressions of
55 basic emotions displayed on command (e.g., [19]). Even though automatic recognition of acted
56 expressions of six basic emotions from face images and image sequences is considered largely solved,
57
58
59
60

1
2
3 reports on novel approaches are published even to date (e.g., [20]). Exceptions from this overall state of
4

5 **Fig. 3.** Prototypic facial expressions of six basic emotions (left-to-right from top
6 row): disgust, happiness, sadness, anger, fear, and surprise.
7

8
9 the art in machine analysis of human facial affect include few tentative efforts to detect acted expressions
10 of cognitive and psychological states like interest [21], fatigue [22], and pain [40]. However detecting
11 these facial expressions in the less constrained environments of real applications is a much more
12 challenging problem which is just beginning to be explored (see section 4).

13 Most commonly used facial action descriptors in sign judgment approaches are the Action Units
14 (AUs) defined in the Facial Action Coding System (FACS; [18]). FACS associates facial expression
15 changes with actions of the muscles that produce them. It defines 9 different action units (AUs) in the
16 upper face, 18 in the lower face, 5 miscellaneous AUs, 11 action descriptors (ADs) for head position, 9
17 ADs for eye position, and 14 additional descriptors for miscellaneous actions (for examples, see Fig. 4).
18 AUs are considered to be the smallest visually discernable facial movements. FACS also provides the
19 rules for recognition of AUs' temporal segments (onset, apex and offset) in a face video. Using FACS,
20 human coders can manually code nearly any anatomically possible facial expression, decomposing it into
21 the specific AUs and their temporal segments that produced the expression. As AUs are independent of
22 interpretation, they can be used for any higher order decision making process including recognition of
23 basic emotions (based on EMFACS rules, [18]), cognitive states like puzzlement [24], psychological states
24 like suicidal depression [7] or pain [25], social behaviours like accord and rapport [6], [24], personality
25 traits like extraversion and temperament [7], and social signals like emblems (i.e., culture-specific
26 interactive signals like wink), regulators (i.e., conversational mediators like nod and smile), and
27 illustrators (i.e., cues accompanying speech like raised eyebrows) [26], [6]. Because it is comprehensive,
28 FACS also allows for the discovery of new patterns related to emotional or situational states. For
29 example, what are the facial behaviours associated with cognitive states like comprehension, or social
30 behaviours like empathy or politeness? How do we build systems to detect comprehension, for example,
31 when we don't know for certain what faces display when students are comprehending? Having subjects
32 pose mental states such as comprehension and puzzlement is of limited use since there is a great deal of
33 evidence that people do different things with their faces when posing versus during a spontaneous
34 experience (see also section 4). Likewise, subjective labelling of expressions has also been shown to be less
35 reliable than objective coding for finding relationships between facial expression and other state
36 variables. An example where subjective judgments of expression failed to find relationships, which were
37 later found with FACS, is the failure of naive subjects to differentiate deception and intoxication from
38 facial display, whereas reliable differences were shown with FACS [27], [28]. Hence, AUs are very
39 suitable to be used as mid-level parameters in automatic facial behaviour analysis, as the thousands of
40 anatomically possible expressions [16], can be described as combinations of 32 AUs and can be mapped to
41 any higher order facial display interpretation including basic emotions, cognitive states, social signals
42 and behaviours, and complex mental states like depression.
43

44
45 It is not surprising, therefore, that automatic AU coding in face images and face image sequences
46 attracted the interest of computer vision researchers. Historically, the first attempts to encode AUs in
47 images of faces in an automatic way were reported by Bartlett et al. in 1996 [29], Lien et al. in 1998 [30],
48 and Pantic et al. in 1998 [31]. These three research groups are still the forerunners in this research field.
49 The focus of the research efforts in the field was first on automatic recognition of AUs in either static face
50 images or face image sequences picturing facial expressions produced on command [14]. Several
51 promising prototype systems were reported that can recognize deliberately produced AUs in either
52 (near-) frontal view face images (e.g., [32], [33], [34]) or profile view face images ([34], [12]). These systems
53 employ ranges of approaches including expert rules, machine learning methods such as neural networks,
54 feature-based image representations (i.e., use geometric features like facial points or shapes of facial
55 components; see also section 2) or appearance-based image representations (i.e., use texture of the facial
56 skin including wrinkles and furrows; see also section 2).
57
58
59
60

1
2
3 One of the main criticisms that these works received from both cognitive and computer scientists, is
4 that the methods are not applicable in real-life situations where subtle changes in facial expression typify
5

6
7 **Fig. 4:** Examples of AUs and their combinations defined in FACS.

8
9 the displayed facial behaviour rather than the exaggerated movements that typify posed expressions.
10 Hence, the focus of the research in the field started to shift to automatic AU recognition in spontaneous
11 facial expressions (produced in a reflex-like manner). Several works have recently emerged on machine
12 analysis of AUs in spontaneous facial expression data. Section 4 provides a detailed discussion of these
13 techniques and the challenges facing the researchers of vision-based analysis of human naturalistic facial
14 behaviour.
15

16 **4. Automatic Analysis of Facial Behaviour: Acted vs. Naturalistic Behaviour**

17
18 The importance of making a clear distinction between spontaneous and deliberately displayed facial
19 behavior for developing and testing computer vision systems becomes apparent when we examine the
20 neurological substrate for facial expression. There are two distinct neural pathways that mediate facial
21 expressions, each one originating in a different area of the brain. Volitional facial movements originate in
22 the cortical motor strip, whereas the more involuntary facial actions originate in the subcortical areas of
23 the brain. The facial expressions mediated by these two pathways have differences both in which facial
24 muscles are moved and in their dynamics [35], [7]. Subcortically initiated facial expressions (the
25 involuntary group) are characterized by synchronized, smooth, symmetrical, consistent, and reflex-like
26 facial muscle movements whereas cortically initiated facial expressions are subject to volitional real-time
27 control and tend to be less smooth, with more variable dynamics. For instance, it has been shown that
28 spontaneous smiles, in contrast to posed smiles (e.g., a polite smile), have smoother transitions between
29 onset, apex, and offset of movement [36], can have multiple AU12 apexes (multiple rises of the mouth
30 corners), and are accompanied by other AUs that appear either simultaneously with AU12 or follow
31 AU12 within 1s [37]. However, precise characterization of spontaneous expression dynamics has been
32 slowed down by the need to use non-invasive technologies (e.g. video), and the difficulty of manually
33 coding AUs, their temporal segments, and intensity frame-by-frame (manual coding of 1 minute of video
34 tape takes approximately 1 hour). Thus the importance of video based automatic coding systems.
35

36 Nonetheless, as already mentioned above, most of the existing works on automatic facial expression
37 recognition are based on deliberate and often exaggerated facial expressions (for survey papers on the
38 past work in the field, see [14], [2], [4]). Few works have been reported on machine analysis of
39 spontaneous facial expression data. When it comes to automatic recognition of affective and mental states
40 from naturalistic facial behaviour data, few tentative efforts have been reported to detect naturalistic
41 expressions of cognitive and psychological states like frustration [23], fatigue [38] and pain [39], [40].
42 Also, few studies investigating dimensional approach to automatic affect recognition have been reported.
43 For example, the study by Zeng et al. [41] investigates automatic discrimination between positive and
44 negative affect, while the study by Ioannou et al. [42] investigates classification of input facial expression
45 data into the quadrants in the evaluation-activation space. A small number of studies have been also
46 reported on automatic recognition of AUs in naturalistic facial behaviour data. These include studies on
47 upper-face AUs only [43], [44], [45], as well as on all AUs [46]. Finally, several recent studies explicitly
48 investigated the difference between spontaneous and deliberate facial behaviour. The work by Valstar et
49 al. from 2006, [45], concerns an automated system for distinguishing posed from spontaneous brow
50 actions (i.e. AU1, AU2, AU4, and their combinations). Conforming with the research findings in
51 psychology, the system was built around characteristics of temporal dynamics of brow actions and
52 employs parameters like speed, intensity, duration, and the occurrence order of brow actions to classify
53 brow actions present in a video as either deliberate or spontaneous facial actions. The work by Littlewort
54 et al. from 2007, [46], reports on an automated system for discriminating genuine from faked pain facial
55 expressions. The system was built around morphological (rather than temporal) characteristics of genuine
56 pain expression (i.e., presence of certain AUs and their intensity). The work by Valstar et al. from 2007,
57
58
59
60

[47], concerns an automated system for distinguishing acted from spontaneous smiles. The study shows that combining information from multiple visual cues (in this case, facial expressions, head movements, and shoulder movements) outperforms single-cue approaches to the target problem. It also clearly shows that the differences between spontaneous and deliberately displayed smiles are in the dynamics of shown behaviour (e.g., the amount of head and shoulder movement, the speed of onset and offset of the actions, and the order and the timing of actions' occurrences) rather than in the configuration of the displayed expression. These findings are in accordance with the research findings in psychology (e.g., [37], [48]). Most of the existing systems for facial expression analysis in naturalistic data are based on 2-D spatial or spatiotemporal facial features and employ advanced probabilistic (e.g., coupled and triple Hidden Markov Models (HMM)), statistical (e.g., Support Vector Machines (SVM) and Relevance Vector Machines (RVM)), and ensemble learning techniques (e.g., Adaboost and Gentleboost).

Although it is obvious that methods of automated facial behaviour analysis that have been trained on deliberate and often exaggerated behaviours may fail to generalize to the complexity of expressive behaviour found in real-world settings (and most probably will fail given the fact that deliberate behaviour differs in visual appearance and timing from spontaneously occurring behaviour), relatively few efforts have been reported so far towards development of systems trained and tested on naturalistic behaviour. There are at least two reasons for this. Firstly, automatic analysis of spontaneously occurring behaviour can hardly be done without analysing the dynamics of the displayed behaviour which, in turn, is a barely investigated research topic as explained in section 5. Secondly, to develop and evaluate facial behaviour analyzers capable of dealing with spontaneously occurring behaviour, large collections of suitable, annotated, publicly available training and test data are needed which, currently, is not the case (see section 6 for a further discussion on this topic).

5. Automatic Analysis of Facial Behaviour: Dynamics of Facial Behaviour

Automatic recognition of facial expression configuration (in terms of AUs constituting the observed expression) has been the main focus of the research efforts in the field. However, both the configuration and the dynamics of facial expressions (i.e., the timing, the duration, the speed of activation and deactivation of various AUs, etc.) are important for interpretation of human facial behaviour. The body of research in cognitive sciences, which argues that the dynamics of facial expressions are crucial for the interpretation of the observed behaviour, is ever growing [49], [7], [50]. Facial expression temporal dynamics are essential for categorization of complex psychological states like various types of pain and mood [25]. They improve judgment of observed facial behaviour (e.g., affect) by enhancing the perception of change and by facilitating processing of facial configuration [50]. They represent a critical factor for interpretation of social behaviours like social inhibition, embarrassment, amusement, and shame [51], [7]. They are also a key parameter in differentiation between posed and spontaneous facial displays [36], [37], [35], [28], as explained in section 4.

In spite of these findings, the vast majority of the past work on machine analysis of human facial behaviour does not take dynamics of facial expressions into account when analyzing shown facial behaviour. Some of the past work in the field has used aspects of temporal dynamics of facial expression such as the speed of a facial point displacement or the persistence of facial parameters over time. However, this was mainly done either in order to increase the performance of facial expression analyzers (e.g., [11], [52], [53]), or in order to report on the intensity of (a component of) the shown facial expression (e.g., [11], [19]), rather than to explicitly encode temporal dynamics of shown facial behaviour. Only few recent studies analyze explicitly the temporal dynamics of facial expressions. These studies explore feature-based approaches to automatic segmentation of AU activation into temporal segments (neutral, onset, apex, offset) in frontal-view [54], [55], and profile-view [12] face videos, appearance-based approaches to automatic coding of temporal segments of AUs [56], and approaches to modelling temporal relationships between AUs as present in expressions of basic emotions [53].

The works by Pantic & Patras from 2005 and 2006, [54], [12], employ rule-based reasoning to encode AUs and their temporal segments based on a set of spatiotemporal features extracted from the trajectories of tracked facial characteristic points. In contrast to biologically inspired learning techniques (such as

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
neural networks), which emulate human unconscious problem solving processes, rule-based techniques are inspired by human conscious problem solving processes. However, studies in cognitive sciences, like the one on “thin slices of behaviour” [6], suggest that facial displays are neither encoded nor decoded at an intentional, conscious level of awareness. They may be fleeting changes in facial appearance that we still accurately judge in terms of emotions or personality even from very brief observations. In turn, this finding suggests that learning techniques inspired by human unconscious problem solving may be more suitable for facial expression recognition than those inspired by human conscious problem solving [57]. Experimental evidence supporting this assumption for the case of prototypic emotional facial expressions was reported in [58]. Experimental evidence supporting this assumption for the case of expression configuration detection and its temporal activation model (neutral → onset → apex → offset) recognition has been recently reported as well [55]. In this latter work, a number of facial characteristic points is detected and tracked in an input face video (particle filtering framework has been used for tracking purposes), a set of spatiotemporal features is extracted from the trajectories of the tracked points, and a combination of statistical and probabilistic machine learning techniques (namely a combination of SVM and HMM) is used to detect AUs and their temporal segments (see Fig. 1 for the outline of the method). The reported experimental results clearly show that modelling facial expression temporal dynamics and analysing displayed facial expressions based on such models significantly improves the performance of the automated system (an increase of 6% in terms of F1 measure was reported).

33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
The work by Koelstra & Pantic from 2008, [56], proposes an appearance-based approach to automatic coding of AUs and their temporal segments. It presents a dynamic-texture-based approach based on non-rigid registration using Free-Form Deformations, in which the extracted facial motion representation is used to derive motion orientation histogram descriptors in both the spatial and temporal domain that, in turn, form further input to a set of AU classifiers based on ensemble and probabilistic machine learning techniques (more specifically, a combination of Gentleboost and HMM was used). This work represents the first appearance-based approach to explicit segmentation of AU activation into temporal segments, reconfirming the results reported in [55] -- modelling facial expression temporal dynamics and analysing displayed facial expressions based on such models significantly improves the performance of the automated system.

The only work reported so far on modelling temporal correlation of different AUs is that by Tong et al. [53]. It applies appearance-based approach to AU recognition, similar to that by Littlewort et al. [19], utilizing Gabor features and a set of Gentleboost classifiers, one for each target AU. It uses further a hierarchical probabilistic framework (more specifically, Dynamic Bayesian Networks) to model the relationships among different AUs as found in facial expressions of six basic emotions. The work reconfirms once again the results reported in [55] – the integration of AU relationships and AU dynamics with AU measurements yields significant improvement of AU recognition (an increase of 5% in correct recognition rate was reported).

Although these pioneering efforts towards automatic analysis of temporal structure of facial expressions are truly promising, many research issues are open and yet to be investigated. A crucial issue that remains unresolved is how the grammar of naturalistic facial behaviour can be learned and how this information can be properly represented and used to handle ambiguities in the input data. Another important issue relates to multi-cue visual analysis. Except for a few studies (e.g., [44], [47]), existing efforts toward the machine analysis of facial behaviour focus only on the analysis of facial gestures without taking into consideration other visual cues like head movements, gaze patterns, and body gestures like shoulder movements. However, research in cognitive science reports that human judgments of behavioural cues are the most accurate when both the face and the body are taken into account [6]. Experimental evidence supporting this finding for the case of automatic laughter analysis was reported in [47]. Taking into account both face and body movements seems to be of particular importance when judging certain complex mental states such as embarrassment [51]. However, integration, temporal structures, and temporal correlations between different visual cues are virtually unexplored areas of research.

6. Evaluating Performance of an Automated System for Facial Behaviour Analysis

The final step in the development of automated systems for facial behaviour analysis is the performance analysis of a developed system. The two crucial aspects of evaluating performance of a designed system are the utilized training/test dataset and the adopted evaluation strategy.

Having enough labelled data of the target human facial behaviour is a prerequisite in designing robust automatic facial expression recognizers. Explorations of this issue showed that, given accurate 3-D alignment of the face, at least 50 training examples are needed for moderate performance (in the 80% accuracy range) of a machine-learning approach to recognition of a specific facial expression [3]. Recordings of spontaneous facial behaviour are difficult to collect because they are difficult to elicit, short lived, and filled with subtle context-based changes. In addition, manual labelling of spontaneous facial behaviour for ground truth is very time consuming, error prone, and expensive. Due to these difficulties, most of the existing studies on automatic facial expression recognition are based on the “artificial” material of deliberately displayed facial behaviour (see also section 4), elicited by asking the subjects to perform a series of facial expressions in front of a camera. Most commonly used, publicly available, annotated datasets of posed facial expressions include the Cohn-Kanade facial expression database [59], and the MMI facial expression database [60]. Yet, as increasing evidence suggests that deliberate (posed) behaviour differs in appearance and timing from that which occurs in daily life (see section 4), it is not surprising that approaches that have been trained on deliberate and often exaggerated behaviours usually fail to generalize to the complexity of expressive behaviour found in real-world settings. To address the general lack of a reference set of (audio and/or) visual recordings of human spontaneous behaviour, several efforts aimed at development of such datasets have been recently reported. Most commonly used, publicly available, annotated datasets of spontaneous human behaviour recordings include SAL dataset, UT Dallas dataset, and MMI-Part2 database [3], [4].

In pattern recognition and machine learning, a common evaluation strategy is to consider correct classification rate (*classification accuracy*) or its complement error rate. However, this assumes that the natural distribution (prior probabilities) of each class are known and balanced. In an imbalanced setting, where the prior probability of the positive class is significantly less than the negative class (the ratio of these being defined as the *skew*), accuracy is inadequate as a performance measure since it becomes biased towards the majority class. That is, as the skew increases, accuracy tends towards majority class performance, effectively ignoring the recognition capability with respect to the minority class. This is a very common (if not the default) situation in facial expression recognition setting, where the prior probability of each target class (a certain facial expression) is significantly less than the negative class (all other facial expressions). Thus, when evaluating performance of an automatic facial expression recognizer, other performance measures such as *precision* (this indicates the probability of correctly detecting a positive test sample and it is independent of class priors), *recall* (this indicates the fraction of the positives detected that are actually correct and, as it combines results from both positive and negative samples, it is class prior dependent), *F1-measure* (this is calculated as $2 * recall * precision / (recall + precision)$), and ROC (this is calculated as $P(x|positive)/P(x|negative)$, where $P(x|C)$ denotes the conditional probability that a data entry has the class label C , and where a ROC curve plots the classification results from the most positive to the most negative classification) are more appropriate. However, as a confusion matrix shows all of the information about a classifier’s performance, it should be used whenever possible for presenting the performance of the evaluated facial expression recognizer.

7. Concluding Remark

Faces are tangible projector panels of the mechanisms which govern our emotional and social behaviours. The automation of the entire process of facial behaviour analysis is, therefore, a highly intriguing problem, the solution to which would be enormously beneficial for fields as diverse as medicine, law, communication, education, and computing. Although the research in the field has seen a lot of progress in the past few years, several issues remain unresolved. Arguably the most important unattended aspect of the problem is how the grammar of facial behaviour (i.e., temporal evolution of occurrences of visual

1
2
3 cues including facial gestures, gaze patterns, and body gestures like head and shoulder movements) can
4 be learned and how this information can be properly represented and used to handle ambiguities in the
5 observation data. This aspect of machine analysis of facial behaviour forms the main focus of the current
6 and future research in the field.
7

8 9 **Acknowledgments**

10 The work of Maja Pantic is funded in part by the European Research Council under the ERC Starting
11 Grant agreement no. ERC-2007-StG-203143 (MAHNOB) and in part by the European Community's 7th
12 Framework Programme [FP7/2007-2013] under the grant agreement no 231287 (SSPNet).
13

14 15 **References**

- 16 [1] Pantic, M., Pentland, A., Nijholt, A. & Huang, T.S., "Human-Centred Intelligent Human-Computer
17 Interaction (HCI²): How far are we from attaining it?", *J. Autonomous and Adaptive Communications*
18 *Systems*, Vol. 1, No. 2, pp. 168-187, 2008.
- 19 [2] Pantic, M. & Rothkrantz, L.J.M., "Toward an Affect-Sensitive Multimodal HCI", *Proc. of the IEEE*, Vol.
20 91, No. 9, pp. 1370-1390, 2003.
- 21 [3] Pantic, M. & Bartlett, M.S., "Machine Analysis of Facial Expressions", in *Face Recognition*, Delac, K. &
22 Grgic, M., Eds., pp. 377-416. I-Tech Education and Publishing, Vienna, Austria, 2007.
- 23 [4] Zeng, Z., Pantic, M., Roisman, G.I. & Huang, T.S., "A Survey of Affect Recognition Methods: Audio,
24 Visual, and Spontaneous Expressions", *IEEE Trans. Pattern Analysis & Machine Intelligence*, Vol. 31,
25 No.1, pp. 39-58, 2009.
- 26 [5] Vinciarelli, A., Pantic, M. & Bourlard, H., "Social Signal Processing: Survey of an Emerging Domain",
27 *J. Image & Vision Computing*, Vol. 27, 2009.
- 28 [6] Ambady, N. & Rosenthal, R. "Thin slices of expressive behavior as predictors of interpersonal
29 consequences: A meta-analysis", *Psychological Bulletin*, Vol. 111, No. 2, pp. 256-274, 1992.
- 30 [7] Ekman, P. & Rosenberg, E.L., Eds., *What the Face Reveals: Basic and Applied Studies of Spontaneous*
31 *Expression using the Facial Action Coding System*. Oxford University Press, Oxford, UK, 2005.
- 32 [8] Yang, M.H., Kriegman, D.J. & Ahuja, N., "Detecting faces in images: A survey", *IEEE Trans. Pattern*
33 *Analysis & Machine Intelligence*, Vol. 24, No. 1, pp. 34-58, 2002.
- 34 [9] Li, S.Z. & Jain, A.K., Eds., *Handbook of Face Recognition*. Springer, New York, USA, 2005.
- 35 [10] Viola, P. & Jones, M., "Robust real-time face detection", *J. Computer Vision*, Vol. 57, No. 2, pp. 137-154,
36 2004.
- 37 [11] Zhang, Y. & Ji, Q., "Active and dynamic information fusion for facial expression understanding from
38 image sequence", *IEEE Trans. Pattern Analysis & Machine Intelligence*, Vol. 27, No. 5, pp. 699-714.
- 39 [12] Pantic, M. & Patras, I., "Dynamics of facial expression: Recognition of facial actions and their
40 temporal segments from face profile image sequences", *IEEE Trans. Systems, Man and Cybernetics -*
41 *Part B*, Vol. 36, No. 2, pp. 433-449, 2006.
- 42 [13] Koelstra, S. & Pantic, M., "Non-rigid registration using free-form deformations for recognition of
43 facial actions and their temporal dynamics", *Proc. IEEE Int'l Conf. Automatic Face and Gesture*
44 *Recognition*, 2008.
- 45 [14] Pantic, M. & Rothkrantz, L.J.M., "Automatic Analysis of Facial Expressions: The State of the Art".
46 *IEEE Trans. Pattern Analysis & Machine Intelligence*, Vol. 22, No. 12, pp. 1424-1445, 2000.
- 47 [15] Gokturk, S.B., Bouguet, J.Y., Tomasi, C. & Girod, B. "Model-based face tracking for view independent
48 facial expression recognition", *Proc. IEEE Int'l Conf. Face and Gesture Recognition*, pp. 272-278, 2002.
- 49 [16] Cohn, J.F. & Ekman, P., "Measuring facial actions", in *The New Handbook of Methods in Nonverbal*
50 *Behavior Research*, Harrigan, J.A., Rosenthal, R. & Scherer, K., Eds., pp. 9-64. Oxford University Press,
51 New York, USA, 2005.
- 52 [17] Keltner, D. & Ekman, P., "Facial Expression of Emotion", in *Handbook of Emotions*, Lewis, M. &
53 Haviland-Jones, J.M., Eds., pp. 236-249. Guilford Press, New York, USA, 2000.
54
55
56
57
58
59
60

- 1
2
3
4 [18] Ekman, P., Friesen, W.V. & Hager, J.C., *Facial Action Coding System*. A Human Face, Salt Lake City, USA, 2002.
- 5
6 [19] Littlewort, G., Bartlett, M.S., Fasel, I., Susskind, J. & Movellan, J., "Dynamics of facial expression
7 extracted automatically from video", *J. Image & Vision Computing*, Vol. 24, No. 6, pp. 615-625, 2006.
- 8 [20] Kotsia, I. & Pitas, I., "Facial expression recognition in image sequences using geometric deformation
9 features and SVM", *IEEE Trans. Image Processing*, Vol. 16, No. 1, pp. 172-187, 2007.
- 10 [21] El Kaliouby, R. & Robinson, P., "Real-Time Inference of Complex Mental States from Facial
11 Expressions and Head Gestures", *Proc. IEEE Int'l Conf. Computer Vision & Pattern Recognition*, vol. 3, p.
12 154, 2004.
- 13 [22] Ji, Q., Lan, P. & Looney, C., "A Probabilistic Framework for Modeling and Real-Time Monitoring
14 Human Fatigue", *IEEE Trans. Systems, Man, and Cybernetics - Part A*, Vol. 36, No. 5, pp. 862-875, 2006.
- 15 [23] Kapoor, A., Bursleson, W. & Picard, R.W., "Automatic Prediction of Frustration," *J. Human-Computer
16 Studies*, Vol. 65, No. 8, pp. 724-736, 2007.
- 17 [24] Cunningham, D.W., Kleiner, M., Wallraven, C. & Bülthoff, H.H., "The components of conversational
18 facial expressions", *Proc. ACM Int'l Conf. Applied Perception in Graphics and Visualization*, pp. 143-149,
19 2004.
- 20 [25] Williams, A.C., "Facial expression of pain: An evolutionary account", *Behavioral & Brain Sciences*, Vol.
21 25, No. 4, pp. 439-488, 2002.
- 22 [26] Ekman, P. & Friesen, W.V., "The repertoire of nonverbal behavior", *Semiotica*, Vol. 1, pp. 49-98, 1969.
- 23 [27] Sayette, M.A., Smith, D.W., Breiner, M.J. & Wilson, G.T., "The effect of alcohol on emotional response
24 to a social stressor", *J. Studies on Alcohol*, Vol. 53, 541-545, 1992.
- 25 [28] Frank, M.G. & Ekman, P., "Appearing truthful generalizes across different deception situations", *J.
26 Personality & Social Psychology*, Vol. 86, 486-495, 2004.
- 27 [29] Bartlett, M.S., Viola, P.A., Sejnowski, T.J., Golomb, B.A., Larsen, J., Hager, J.C. & Ekman, P.,
28 "Classifying facial actions", *Advances in Neural Information Processing Systems 8*, pp. 823-829, 1996.
- 29 [30] Lien, J.J.J., Kanade, T., Cohn, J.F. & Li, C.C., "Subtly different facial expression recognition and
30 expression intensity estimation", *Proc. IEEE Int'l Conf. Computer Vision & Pattern Recognition*, pp. 853-
31 859, 1998.
- 32 [31] Pantic, M., Rothkrantz, L.J.M. & Koppelaar, H., "Automation of non-verbal communication of facial
33 expressions", *Proc. Conf. Euromedia*, pp. 86-93, 1998.
- 34 [32] Bartlett, M.S., Hager, J. C., Ekman, P. & Sejnowski, T.J., "Measuring facial expressions by computer
35 image analysis", *Psychophysiology*, Vol. 36, No. 2, pp. 253-263, 1999.
- 36 [33] Tian, Y.L., Kanade, T. & Cohn, J.F., "Recognizing action units for facial expression analysis", *IEEE
37 Trans. Pattern Analysis & Machine Intelligence*, Vol. 23, No. 2, pp. 97-115, 2001.
- 38 [34] Pantic, M. & Rothkrantz, L.J.M., "Facial action recognition for facial expression analysis from static
39 face images", *IEEE Trans. Systems, Man and Cybernetics - Part B*, Vol. 34, No. 3, pp. 1449-1461, 2004.
- 40 [35] Ekman, P., "Darwin, deception, and facial expression", *Annals New York Academy of sciences*, Vol.
41 1000, pp. 205-221, 2003.
- 42 [36] Frank, M.G., Ekman, P. & Friesen, W.V., "Behavioural markers and reognizability of the smile of
43 enjoyment", *J. Personality & Social Psychology*, Vol. 64, pp. 83-93, 1993.
- 44 [37] J.F. Cohn and K.L. Schmidt, "The timing of facial motion in posed and spontaneous smiles", *J.
45 Wavelets, Multi-resolution & Information Processing*, vol. 2, no. 2, pp. 121-132, 2004.
- 46 [38] Fan, X., Sun, Y. & Yin, B., "Multi-scale dynamic human fatigue detection with feature level fusion",
47 *Proc. IEEE Int'l Conf. Automatic Face & Gesture Recognition*, 2008.
- 48 [39] Ashraf, A.B., Lucey, S., Cohn, J.F., Chen, T., Ambadar, Z., Prkachin, K., Solomon, P. & Theobald, B.J.
49 "The Painful Face: Pain Expression Recognition Using Active Appearance Models," *Proc. ACM Int'l
50 Conf. Multimodal Interfaces*, pp. 9-14, 2007.
- 51 [40] Littlewort, G.C., Bartlett, M.S. & Lee, K., "Faces of Pain: Automated Measurement of Spontaneous
52 Facial Expressions of Genuine and Posed Pain", *Proc. ACM Int'l Conf. Multimodal Interfaces*, pp. 15-21,
53 2007.
- 54 [41] Zeng, Z., Fu, Y., Roisman, G.I., Wen, Z., Hu, Y. & Huang, T.S., "Spontaneous Emotional Facial
55 Expression Detection", *J. Multimedia*, vol. 1, no. 5, pp. 1-8, 2006.
- 56
57
58
59
60

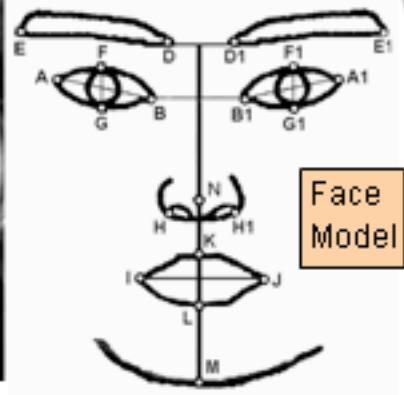
- 1
2
3 [42] Ioannou, S., Raouzaiou, A., Tzouvaras, V., Mailis, T., Karpouzis, K. & Kollias, S., "Emotion
4 Recognition through Facial Expression Analysis Based on a Neurofuzzy Method", *J. Neural Networks*,
5 Vol. 18, No. 4, pp. 423-435, 2005.
- 6 [43] Kapoor, A., Qi, Y. & Picard, R.W., "Fully Automatic Upper Facial Action Recognition", Proc. IEEE
7 Int'l Workshop on Analysis and Modeling of Faces and Gestures, 2003.
- 8 [44] Cohn, J.F., Reed, L.I., Ambadar, Z., Xiao, J. & Moriyama, T., "Automatic Analysis and Recognition of
9 Brow Actions and Head Motion in Spontaneous Facial Behavior", *Proc. IEEE Int'l Conf. Systems, Man,
10 and Cybernetics*, vol. 1, pp. 610-616, 2004.
- 11 [45] Valstar, M.F., Pantic, M., Ambadar, Z. & Cohn, J.F., "Spontaneous versus Posed Facial Behavior:
12 Automatic Analysis of Brow Actions", *Proc. ACM Int'l Conf. Multimodal Interfaces*, pp. 162-170, 2006.
- 13 [46] Bartlett, M.S., Littlewort, G., Frank, M., Lainscsek, C., Fasel, I. & Movellan, J. "Recognizing Facial
14 Expression: Machine Learning and Application to Spontaneous Behavior", *Proc. IEEE Int'l Conf.
15 Computer Vision and Pattern Recognition*, pp. 568-573, 2005.
- 16 [47] Valstar, M.F., Gunes, H. & Pantic, M., "How to Distinguish Posed from Spontaneous Smiles Using
17 Geometric Features", *Proc. ACM Int'l Conf. Multimodal Interfaces*, pp. 38-45, 2007.
- 18 [48] Krumhuber, E., Manstead, A.S.R. & Kappas, A., "Temporal aspects of facial displays in person and
19 expression perception: The effects of smile dynamics, head-tilt, and gender", *J. Nonverbal Behaviour*,
20 Vol. 31, pp. 39-56, 2007.
- 21 [49] Russell, J.A. & Fernandez-Dols, J.M., Eds. *The Psychology of Facial Expression*, Cambridge University
22 Press, New York, USA, 1997.
- 23 [50] Ambadar, Z., Schooler, J. & Cohn, J.F., "Deciphering the enigmatic face: The importance of facial
24 dynamics in interpreting subtle facial expressions", *Psychological Science*, Vol. 16, No. 5, pp. 403-410,
25 2005.
- 26 [51] Costa, M., Dinsbach, W., Manstead, A.S.R. & Bitti, P.E.R., "Social presence, embarrassment, and
27 nonverbal behaviour", *J. Nonverbal Behaviour*, Vol. 25., No. 4, pp. 225-240, 2001.
- 28 [52] Gralewski, L., Campbell, N. & Voak, I.P., "Using a tensor framework for the analysis of facial
29 dynamics", *Proc. IEEE Int'l Conf. Face & Gesture Recognition*, pp. 217-222, 2006.
- 30 [53] Tong, Y., Liao, W. & Ji, Q., "Facial action unit recognition by exploiting their dynamics and semantic
31 relationships", *IEEE Trans. Pattern Analysis & Machine Intelligence*, Vol. 29, No. 10, pp. 1683-1699, 2007.
- 32 [54] Pantic, M. & Patras, I., "Detecting facial actions and their temporal segments in nearly frontal-view
33 face image sequences", *Proc. IEEE Int'l Conf. Systems, Man & Cybernetics*, pp. 3358-3363, 2005.
- 34 [55] M.F. Valstar and M. Pantic, "Combined Support Vector Machines and Hidden Markov Models for
35 modeling facial action temporal dynamics", *LNCS*, Vol. 4796, pp. 118-127, 2007.
- 36 [56] Koelstra, S. & Pantic, M., "Non-rigid registration using free-form deformations for recognition of
37 facial actions and their temporal dynamics", *Proc. IEEE Int'l Conf. Automatic Face & Gesture
38 Recognition*, 2008.
- 39 [57] Pantic, M., Sebe, N., Cohn, J.F. & Huang, T., "Affective Multimodal Human-Computer Interaction",
40 *Proc. ACM Int'l Conf. Multimedia*, pp. 669-676, 2005.
- 41 [58] Valstar, M.F. & Pantic, M., "Biologically vs. logic inspired encoding of facial actions and emotions in
42 video", *Proc. IEEE Int'l Conf. Multimedia & Expo*, pp. 325-328, 2006.
- 43 [59] Kanade, T., Cohn, J.F. & Tian, Y., "Comprehensive database for facial expression analysis", *Proc. IEEE
44 Int'l Conf. Automatic Face & Gesture Recognition*, pp. 46-53, 2000.
- 45 [60] Pantic, M., Valstar, M.F., Rademaker, R. & Maat, L., "Web-based database for facial expression analysis", *Proc.
46 IEEE Int'l Conf. Multimedia & Expo*, pp. 317-321, 2005.
- 47
48
49
50
51
52
53
54
55
56
57
58
59
60

Input Video

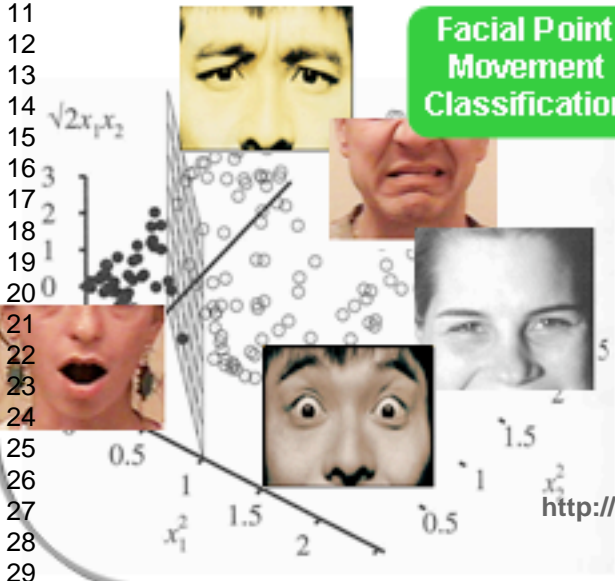
Submitted to Phil. Trans. R. Soc. B - Issue

Page 12 of 15

Facial Point Localization



Facial Point Movement Classification



Facial Point Tracking

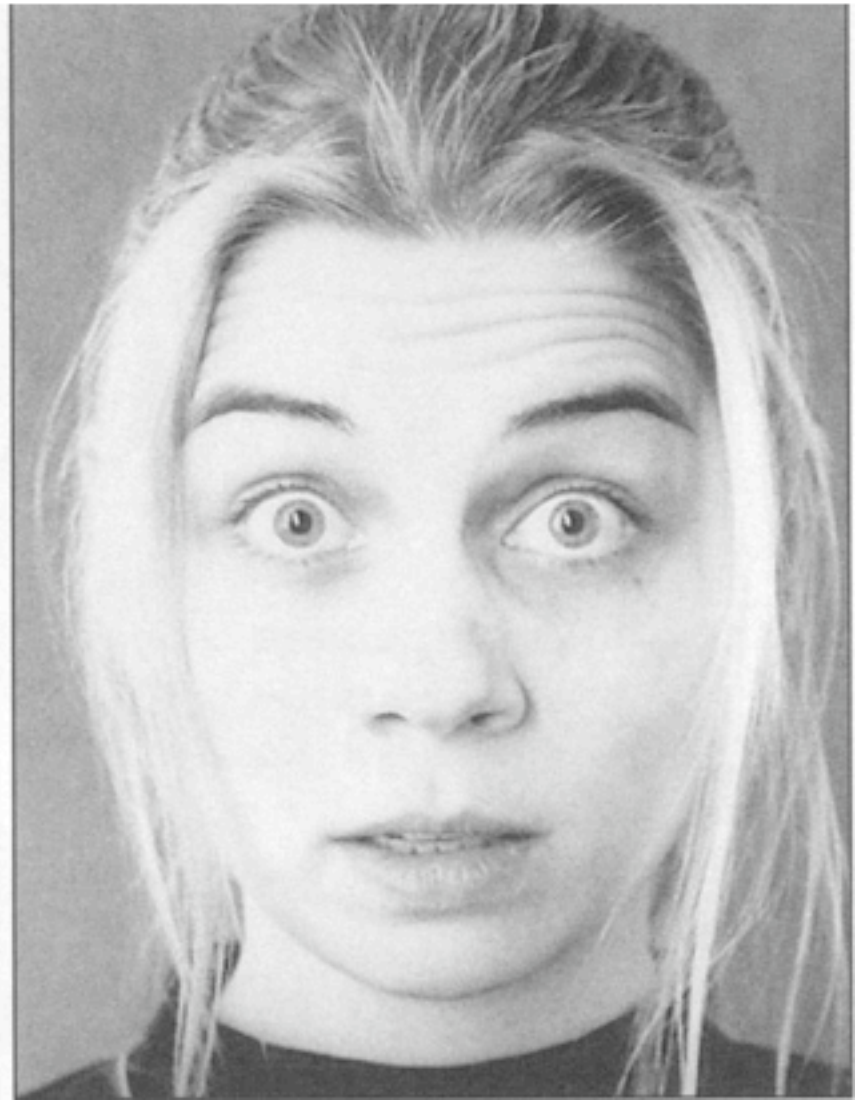
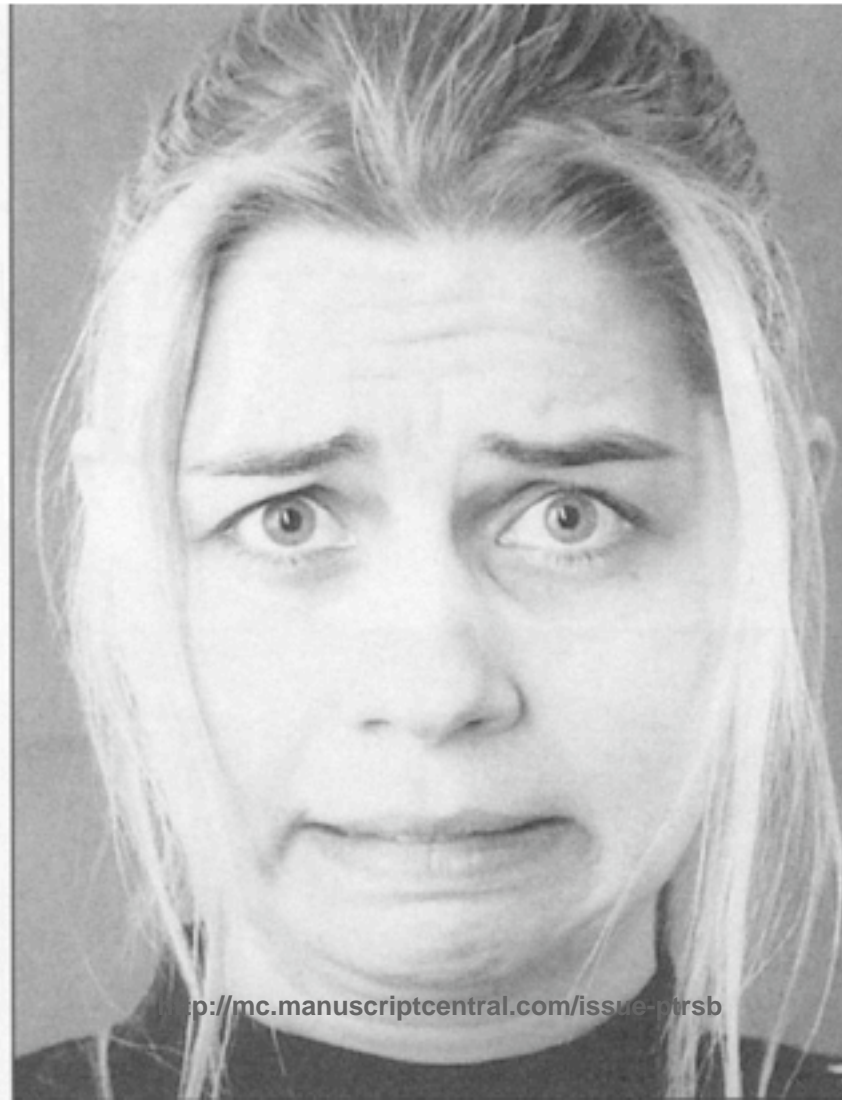
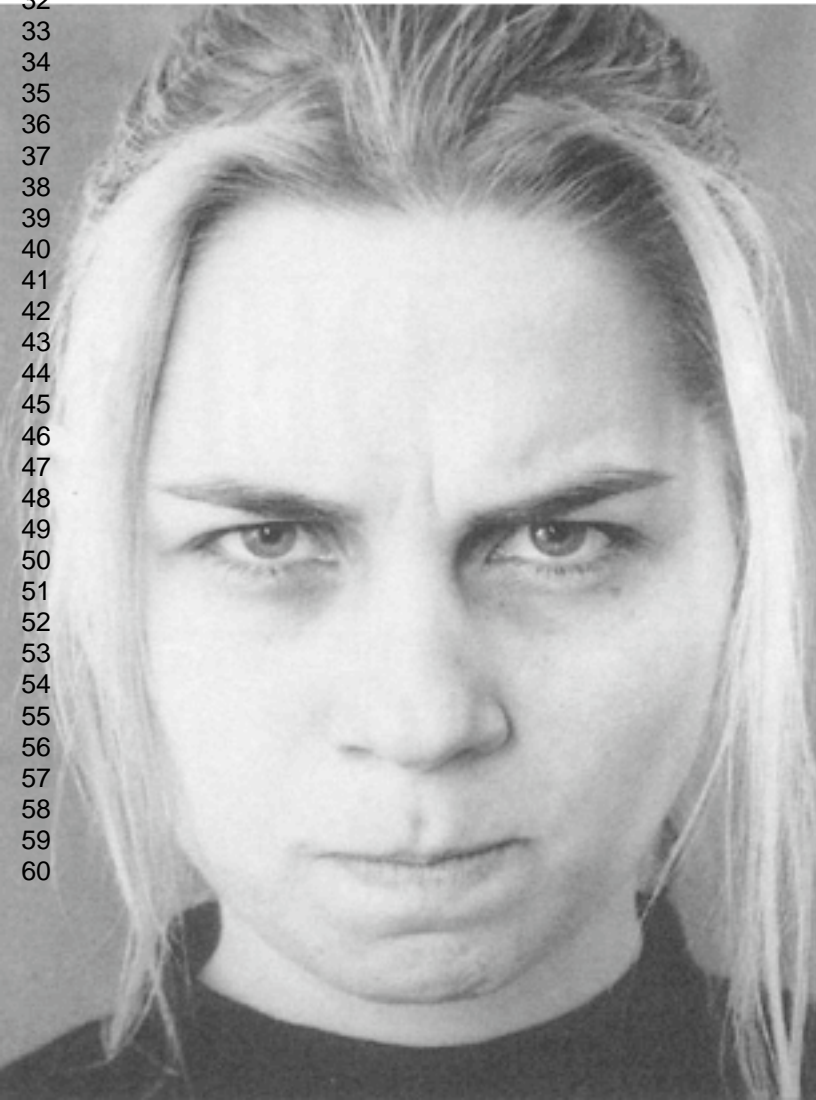
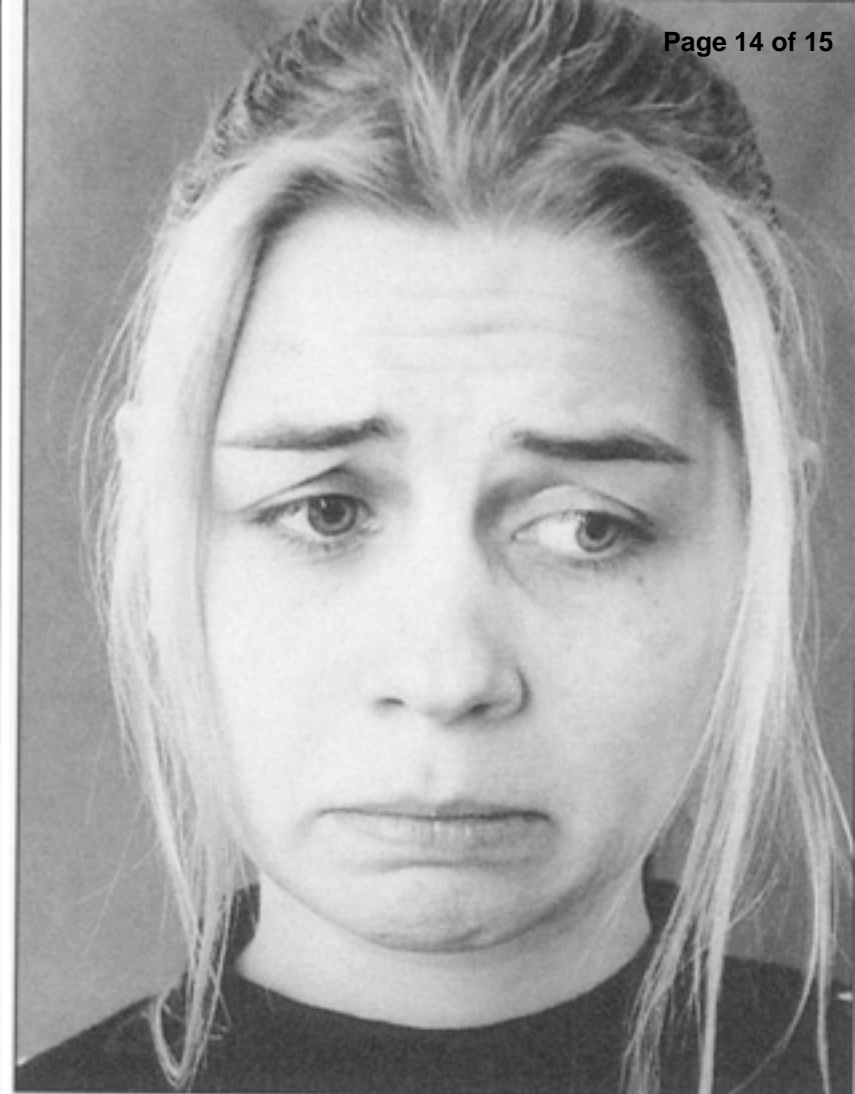
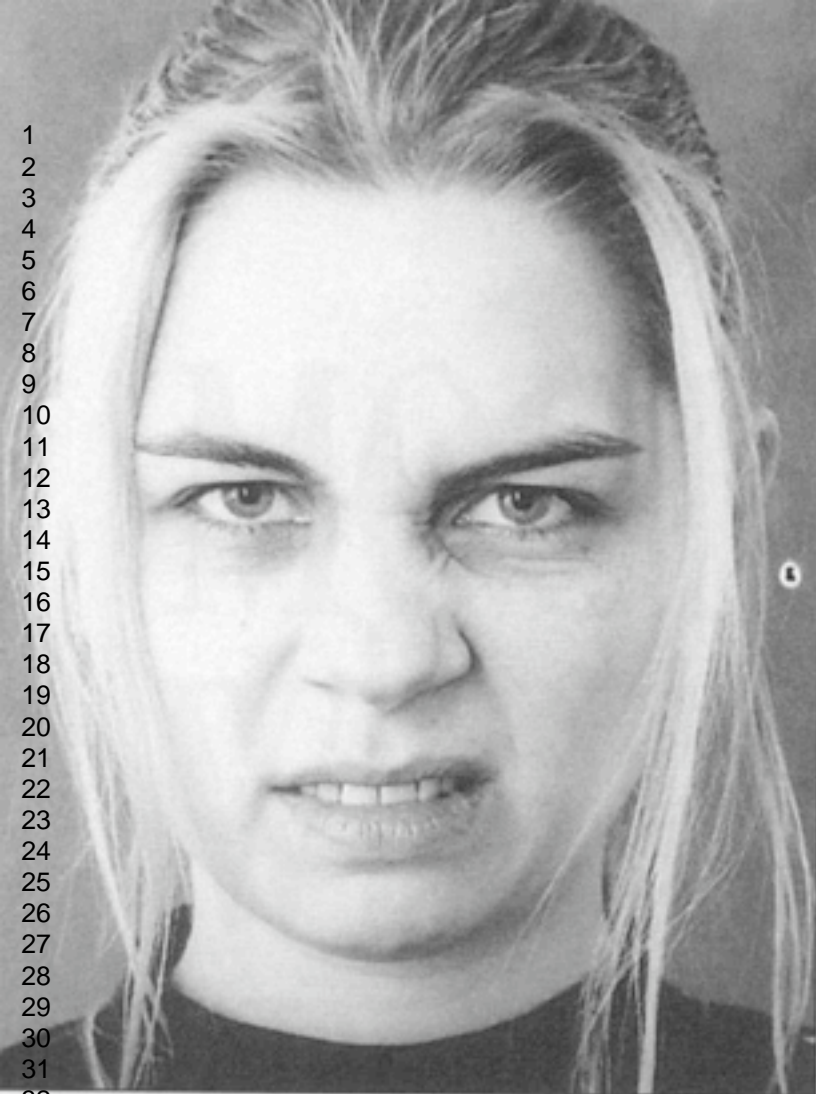


<http://mc.manuscriptcentral.com/issue-ptrsb>

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38

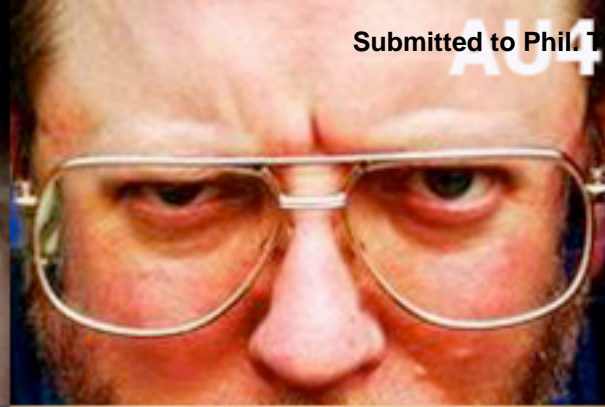


1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60





AU1



AU4



AU7



AU12



AU2



AU5



AU9



AU13



AU1+AU2



AU6



AU8+10+16



AU14

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37