

In A. Esposito, N. Bourbakis, N. Avouris, and I. Hatzilygeroudis. (Eds.) Lecture Notes in Computer Science, Vol 5042: Verbal and Nonverbal Features of Human-human and Human-machine Interaction, Springer Verlag, p. 1-21.

## Data mining spontaneous facial behavior with automatic expression coding

Marian Bartlett<sup>1</sup>, Gwen Littlewort<sup>1</sup>, Esra Vural<sup>1,3</sup>,  
Kang Lee<sup>2</sup>, Mujdat Cetin<sup>3</sup>, Aytul Ercil<sup>3</sup>, and Javier Movellan<sup>1</sup>

<sup>1</sup> Institute for Neural Computation, University of California, San Diego, La Jolla, CA 92093-0445, USA

<sup>2</sup> Human Development and Applied Psychology, University of Toronto, Ontario, Canada

<sup>3</sup> Engineering and Natural Science, Sabanci University, Istanbul, Turkey  
[mbartlett@ucsd.edu](mailto:mbartlett@ucsd.edu); [gwen@mpmlab.ucsd.edu](mailto:gwen@mpmlab.ucsd.edu), [movellan@mplab.ucsd.edu](mailto:movellan@mplab.ucsd.edu),  
[vesra@ucsd.edu](mailto:vesra@ucsd.edu), [kang.lee@utoronto.ca](mailto:kang.lee@utoronto.ca)

**Abstract.** The computer vision field has advanced to the point that we are now able to begin to apply automatic facial expression recognition systems to important research questions in behavioral science. The machine perception lab at UC San Diego has developed a system based on machine learning for fully automated detection of 30 actions from the facial action coding system (FACS). The system, called Computer Expression Recognition Toolbox (CERT), operates in real-time and is robust to the video conditions in real applications. This paper describes two experiments which are the first applications of this system to analyzing spontaneous human behavior: Automated discrimination of posed from genuine expressions of pain, and automated detection of driver drowsiness. The analysis revealed information about facial behavior during these conditions that were previously unknown, including the coupling of movements. Automated classifiers were able to differentiate real from fake pain significantly better than naïve human subjects, and to detect critical drowsiness above 98% accuracy. Issues for application of machine learning systems to facial expression analysis are discussed.

**Keywords:** Facial expression recognition, machine learning.

## 1 Introduction

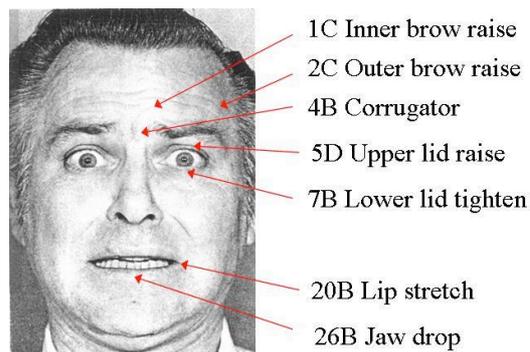
The computer vision field has advanced to the point that we are now able to begin to apply automatic facial expression recognition systems to important research questions in behavioral science. This paper is among the first applications of fully automated facial expression measurement to such research questions. It explores two applications of a machine learning system for automatic facial expression measurement to data mine spontaneous human behavior (1) differentiating fake from real expressions of pain, and (2) detecting driver drowsiness.

Based on the following two conference papers:

- (1) Littlewort, G., Bartlett, M.S. and Lee, K., (2007). Automated measurement of spontaneous facial expressions of genuine and posed pain. Proc. International Conference on Multimodal Interfaces, Nagoya, Japan. Copyright 2007 ACM 978-1-59593-817-6/07/0011..
- (2) Vural, E., Cetin, M., Ercil, A., Littlewort, G., Bartlett, M., and Movellan, J. (2007). Machine learning systems for detecting driver drowsiness. Proc. Digital Signal Processing for in-Vehicle and mobile systems, Istanbul, Turkey. Copyright 2007 IEEE.

## 1.1 The Facial Action Coding System

The facial action coding system (FACS) (Ekman and Friesen, 1978) is arguably the most widely used method for coding facial expressions in the behavioral sciences. The system describes facial expressions in terms of 46 component movements, which roughly correspond to the individual facial muscle movements. An example is shown in Figure 1. FACS provides an objective and comprehensive way to analyze expressions into elementary components, analogous to decomposition of speech into phonemes. Because it is comprehensive, FACS has proven useful for discovering facial movements that are indicative of cognitive and affective states. See Ekman and Rosenberg (2005) for a review of facial expression studies using FACS. The primary limitation to the widespread use of FACS is the time required to code. FACS was developed for coding by hand, using human experts. It takes over 100 hours of training to become proficient in FACS, and it takes approximately 2 hours for human experts to code each minute of video. The authors have been developing methods for fully automating the facial action coding system (e.g. Donato et al., 1999; Bartlett et al., 2006). In this paper we apply a computer vision system trained to automatically detect FACS to data mine facial behavior under two conditions: (1) real versus fake pain, and (2) driver fatigue.



**Fig. 1.** Example facial action decomposition from the facial action coding system. A prototypical expression of fear is decomposed into 7 component movements. Letters indicate intensity. A fear brow (1+2+4) is illustrated here.

## 1.2 Spontaneous Expressions

The machine learning system presented here was trained on spontaneous facial expressions. The importance of using spontaneous behavior for developing and testing computer vision systems becomes apparent when we examine the neurological substrate for facial expression. There are two distinct neural pathways that mediate facial expressions, each one originating in a different area of the brain. Volitional facial movements originate in the cortical motor strip, whereas spontaneous facial expressions originate in the subcortical areas of the brain (see Rinn, 1984, for a review). These two pathways have different patterns of innervation on the face, with the cortical system tending to give stronger innervation to certain muscles primarily

in the lower face, while the subcortical system tends to more strongly innervate certain muscles primarily in the upper face (e.g. Morecraft et al., 2001).

The facial expressions mediated by these two pathways have differences both in which facial muscles are moved and in their dynamics (Ekman, 2001; Ekman & Rosenberg, 2005). Subcortically initiated facial expressions (the spontaneous group) are characterized by synchronized, smooth, symmetrical, consistent, and reflex-like facial muscle movements whereas cortically initiated facial expressions (posed expressions) are subject to volitional real-time control and tend to be less smooth, with more variable dynamics (Rinn, 1984; Frank, Ekman, & Friesen, 1993; Schmidt, Cohn & Tian, 2003; Cohn & Schmidt, 2004). Given the two different neural pathways for facial expressions, it is reasonable to expect to find differences between genuine and posed expressions of states such as pain or drowsiness. Moreover, it is crucial that the computer vision model for detecting states such as genuine pain or driver drowsiness is based on machine learning of expression samples when the subject is actually experiencing the state in question.

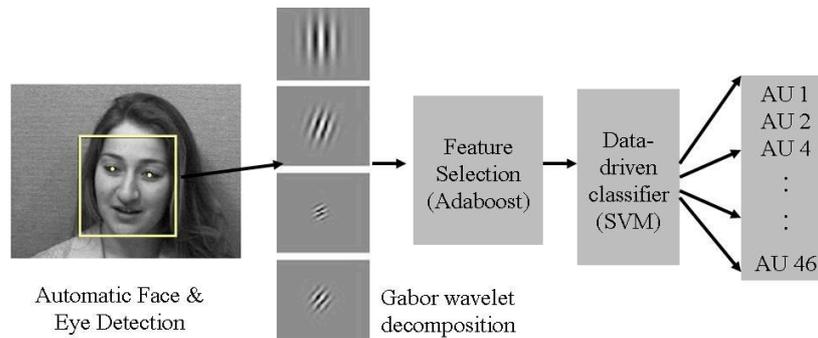


Fig. 2. Overview of the automated facial action recognition system.

## 2 The Computer Expression Recognition Toolbox (CERT)

Here we extend a system for fully automated facial action coding developed previously by the authors (Bartlett et al., 2006; Littlewort et al., 2006). It is a user independent fully automatic system for real time recognition of facial actions from the Facial Action Coding System (FACS). The system automatically detects frontal faces in the video stream and codes each frame with respect to 20 Action units. In previous work, we conducted empirical investigations of machine learning methods applied to the related problem of classifying expressions of basic emotions. We compared image features (e.g. Donato et al., 1999), classifiers such as AdaBoost, support vector machines, and linear discriminant analysis, as well as feature selection techniques (Littlewort et al., 2006). Best results were obtained by selecting a subset of Gabor filters using AdaBoost and then training Support Vector Machines on the outputs of the filters selected by AdaBoost. An overview of the system is shown in Figure 2.

## 2.1 Real Time Face and Feature Detection

We employed a real-time face detection system that uses boosting techniques in a generative framework (Fasel et al.) and extends work by Viola and Jones (2001). Enhancements to Viola and Jones include employing Gentleboost instead of AdaBoost, smart feature search, and a novel cascade training procedure, combined in a generative framework. Source code for the face detector is freely available at <http://kolmogorov.sourceforge.net>. Accuracy on the CMU-MIT dataset, a standard public data set for benchmarking frontal face detection systems (Schneiderman & Kanade, 1998), is 90% detections and 1/million false alarms, which is state-of-the-art accuracy. The CMU test set has unconstrained lighting and background. With controlled lighting and background, such as the facial expression data employed here, detection accuracy is much higher. All faces in the training datasets, for example, were successfully detected. The system presently operates at 24 frames/second on a 3 GHz Pentium IV for 320x240 images. The automatically located faces were rescaled to 96x96 pixels. The typical distance between the centers of the eyes was roughly 48 pixels. Automatic eye detection (Fasel et al., 2005) was employed to align the eyes in each image. The images were then passed through a bank of Gabor filters 8 orientations and 9 spatial frequencies (2:32 pixels per cycle at 1/2 octave steps). Output magnitudes were then passed to the action unit classifiers.

## 2.2 Automated Facial Action Classification

The approach presented here is a 2-stage system in which first an automated system CERT, is developed for detecting action units, and secondly CERT is applied to spontaneous examples of a state in question, and machine learning is applied to the CERT outputs. Here we describe the training of the facial action detectors in the first stage. The training data for the facial action classifiers came from three posed datasets and one dataset of spontaneous expressions. The facial expressions in each dataset were FACS coded by certified FACS coders. The first posed dataset was the Cohn-Kanade DFAT-504 dataset (Kanade, Cohn & Tian, 2000). This dataset consists of 100 university students who were instructed by an experimenter to perform a series of 23 facial displays, including expressions of seven basic emotions. The second posed dataset consisted of directed facial actions from 24 subjects collected by Ekman and Hager. Subjects were instructed by a FACS expert on the display of individual facial actions and action combinations, and they practiced with a mirror. The resulting video was verified for AU content by two certified FACS coders. The third posed dataset consisted of a subset of 50 videos from 20 subjects from the MMI database (Pantic et al., 2005). The spontaneous expression dataset consisted of the FACS-101 dataset collected by Mark Frank (Bartlett et. al. 2006). 33 subjects underwent an interview about political opinions on which they felt strongly. Two minutes of each subject were FACS coded. The total training set consisted of 5500 examples, 2500 from posed databases and 3000 from the spontaneous set.

Twenty linear Support Vector Machines were trained for each of 20 facial actions. Separate binary classifiers, one for each action, were trained to detect the presence of

the action in a one versus all manner. Positive examples consisted of the apex frame for the target AU. Negative examples consisted of all apex frames that did not contain the target AU plus neutral images obtained from the first frame of each sequence. Eighteen of the detectors were for individual action units, and two of the detectors were for specific brow region combinations: fear brow (1+2+4) and distress brow (1 alone or 1+4). All other detectors were trained to detect the presence of the target action regardless of co-occurring actions. A list is shown in Table 1A. Thirteen additional AU's were trained for the Driver Fatigue Study. These are shown in Table 1B.

**Table 1A.** AU detection performance on posed and spontaneous facial actions. Values are Area under the roc (A') for generalization to novel subjects.

<b>AU</b>	<b>Name</b>	<b>Posed</b>	<b>Spont</b>
1	Inner brow raise	.90	.88
2	Outer brow raise	.94	.81
4	Brow Lower	.98	.73
5	Upper Lid Raise	.98	.80
6	Cheek Raise	.85	.89
7	Lids tight	.96	.77
9	Nose wrinkle	.99	.88
10	Upper lip raise	.98	.78
12	Lip corner pull	.97	.92
14	Dimpler	.90	.77
15	Lip corner Depress	.80	.83
17	Chin Raise	.92	.80
18	Lip Pucker	.87	.70
20	Lip stretch	.98	.60
23	Lip tighten	.89	.63
24	Lip press	.84	.80
25	Lips part	.98	.71
26	Jaw drop	.98	.71
1,1+4	Distress brow	.94	.70
1+2+4	Fear brow	.95	.63
Mean:		.93	.77

**Table 1B:** Additional 13 AU's trained for the driver fatigue study.

AU	Name
8	Lip Toward Each Other
11	Nasolabial Furrow Deepener
13	Sharp Lip Puller
16	Lower Lip Depress
19	Tongue Show
22	Lip Funneller
27	Mouth Stretch
28	Lips Suck
30	Jaw Sideways
32	Bite
38	Nostril Dilate
39	Nostril Compress
45	Blink

The output of the system was a real valued number indicating the distance to the separating hyperplane for each classifier. Previous work showed that the distance to the separating hyperplane (the margin) contained information about action unit intensity (e.g. Bartlett et al., 2006).

In this paper, area under the ROC ( $A'$ ) is used to assess performance rather than overall percent correct, since percent correct can be an unreliable measure of performance, as it depends on the proportion of targets to non-targets, and also on the decision threshold. Similarly, other statistics such as true positive and false positive rates depend on decision threshold, which can complicate comparisons across systems.  $A'$  is a measure derived from signal detection theory and characterizes the discriminative capacity of the signal, independent of decision threshold. The ROC curve is obtained by plotting true positives against false positives as the decision threshold shifts from 0 to 100% detections. The area under the ROC ( $A'$ ) ranges from 0.5 (chance) to 1 (perfect discrimination).  $A'$  can also be interpreted in terms of percent correct.  $A'$  is equivalent to the theoretical maximum percent correct achievable with the information provided by the system when using a 2-Alternative Forced Choice testing paradigm.

Table 1 shows performance for detecting facial actions in posed and spontaneous facial actions. Generalization to novel subjects was tested using 3-fold cross-validation on the images in the training set. Performance was separated into the posed set, which was 2,500 images, and a spontaneous set, which was 1100 images from the FACS-101 database which includes speech.

The overall CERT system gives a frame-by-frame output with N channels, consisting of N facial actions. This system can be applied to data mine human behavior. By applying CERT to face video while subjects experience spontaneous expressions of a given state, we can learn new things about the facial behaviors associated with that state. Also, by passing the N channel output to a machine learning system, we can directly train detectors for the specific state in question. (See Figure 3.) In Sections 3 and 4, two implementations of this idea are described.

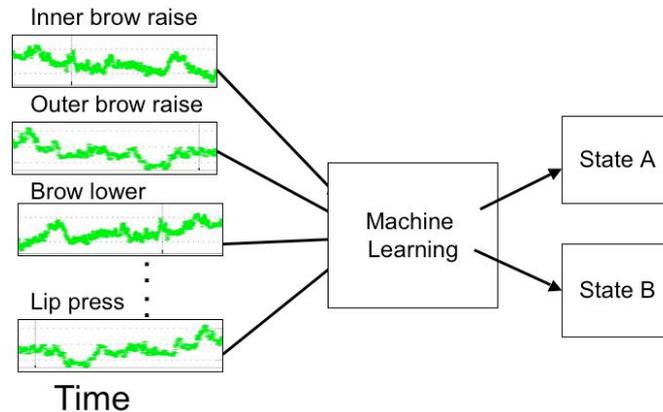


Figure 3. Data mining human behavior. CERT is applied to face videos containing spontaneous expressions of states in question. Machine learning is applied to the outputs of CERT to learn a classifier to automatically discriminate state A from State B.

### 3 Classification of Real versus Faked Pain Expressions

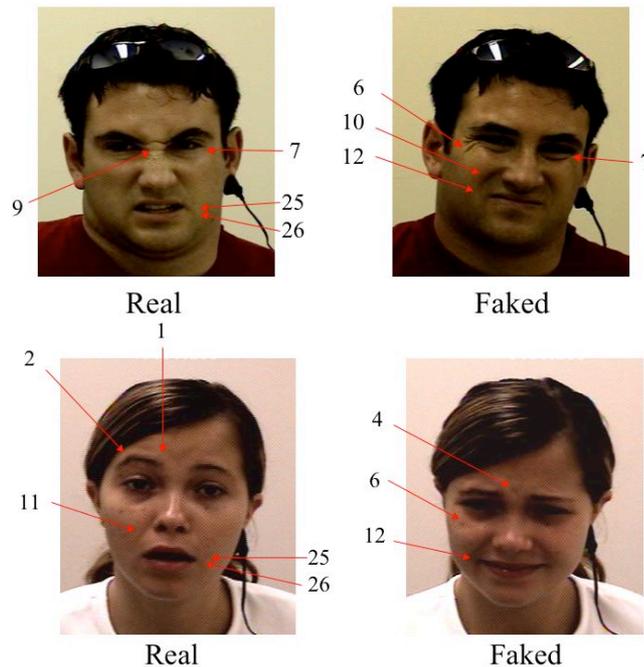
An important issue in medicine is the ability to distinguish real pain from faked pain, (malingering). Some studies suggest that malingering rates are as high as 10% in chronic pain patients (Fishbain et al., 1999), and much higher in litigation contexts (Schmand et al., 1998). Even more important is to recognize when patients are experiencing genuine pain so that their pain is taken seriously. There is presently no reliable method for physicians to differentiate faked from real pain (Fishbain, 2006). Naïve human subjects are near chance for differentiating real from fake pain from observing facial expression (e.g. Hadjistavropoulos et al., 1996). In the absence of direct training in facial expressions, clinicians are also poor at assessing pain from the face (e.g. Prkachin et al. 2002; Grossman, 1991). However a number of studies using the Facial Action Coding System (FACS) (Ekman & Friesen, 1978) have shown that information exists in the face for differentiating real from posed pain (e.g. Hill and Craig, 2002; Craig et al., 1991; Prkachin 1992).

In previous studies using manual FACS coding by human experts, at least 12 facial actions showed significant relationships with pain across multiple studies and pain modalities. Of these, the ones specifically associated with cold pressor pain were 4, 6, 7, 9, 10, 12, 25, 26 (Craig & Patrick, 1985; Prkachin, 1992). See Table 1 and Figure 2 for names and examples of these AU's. A previous study compared faked to real pain, but in a different pain modality (lower back pain). This study found that when faking subjects tended to display the following AU's: 4, 6, 7, 10, 12, 25. When faked pain expressions were compared to real pain expressions, the faked pain expressions contained significantly more brow lower (AU 4), cheek raise (AU 6), and lip corner pull (AU 12) (Craig, Hyde & Patrick, 1991). These studies also reported substantial individual differences in the expressions of both real pain and faked pain.

Recent advances in automated facial expression measurement, such as the CERT system described above, open up the possibility of automatically differentiating posed from real pain using computer vision systems (e.g. Bartlett et al., 2006; Littlewort et al., 2006; Cohn & Schmidt, 2004; Pantic et al., 2006). This section explores the application of CERT to this problem.

### 3.1 Human subject methods

Video data was collected of 26 human subjects during real pain, faked pain, and baseline conditions. Human subjects were university students consisting of 6 men and 20 women. The pain condition consisted of cold pressor pain induced by immersing the arm in cold water at 5<sup>o</sup> Celsius. For the baseline and faked pain conditions, the water was 20<sup>o</sup> Celsius. Subjects were instructed to immerse their forearm into the water up to the elbow, and hold it there for 60 seconds in each of the three conditions. The order of the conditions was baseline, faked pain, and then real pain. For the faked pain condition, subjects were asked to manipulate their facial expressions so that an “expert would be convinced they were in actual pain.” Participants facial expressions were recorded using a digital video camera during each condition. Examples are shown in Figure 4.



**Figure 4.** Sample facial behavior and facial action codes from the real and faked pain conditions.

A second subject group underwent the conditions in the counterbalanced order, with real pain followed by faked pain. This ordering involves immediate motor

memory, which is a fundamentally different task. The present paper therefore analyzes only the first subject group. The second group will be analyzed separately in a future paper, and compared to the first group.

After the videos were collected, a set of 170 naïve observers were shown the videos and asked to guess whether each video contained faked or real pain. Subjects were undergraduates with no explicit training in facial expression measurement. They were primarily Introductory Psychology students at UCSD. Mean accuracy of naïve human subjects for discriminating fake from real pain in these videos was at chance at 49.1% (standard deviation 13.7%). These observers had no specific training in facial expression and were not clinicians. One might suppose that clinicians would be more accurate. However previous studies suggest that clinicians judgments of pain from the face are similarly unreliable (e.g. Grossman, 1991). Facial signals do appear to exist however (Hill & Craig, 2002, Craig et al., 1991; Prkachin 1992), and immediate corrective feedback has been shown to improve observer accuracy (Hill & Craig, 2004).

### **3.2 Human expert FACS coding**

In order to assess the validity of the automated system, we first obtained FACS codes for a portion of the video from a human FACS expert certified in the Facial Action Coding System. For each subject, the last 500 frames of the fake pain and real pain conditions were FACS coded (about 15 seconds each). It took 60 man hours to collect the human codes, over the course of more than 3 months, since human coders can only code up to 2 hours per day before having negative repercussions in accuracy and coder burn-out.

The sum of the frames containing each action unit were collected for each subject condition, as well as a weighted sum, multiplied by the intensity of the action on a 1-5 scale. To investigate whether any action units successfully differentiated real from faked pain, paired t-tests were computed on each individual action unit. (Tests on specific brow region combinations 1+2+4 and 1,1+4 have not yet been conducted.) The one action unit that significantly differentiated the two conditions was AU 4, brow lower, ( $p < .01$ ) for both the sum and weighted sum measures. This finding is consistent with the analysis of the automated system, which also found action unit 4 most discriminative.

### **3.3 Automated coding**

#### **3.3.1 Characterizing the differences between real and faked pain**

Applying CERT to the pain video data produced a 20 channel output stream, consisting of one real value for each learned AU, for each frame of the video. This data was further analyzed to predict the difference between baseline and pained faces, and the difference between expressions of real pain and fake pain.

We first examined which facial action detectors were elevated in real pain compared to the baseline condition. Z-scores for each subject and each AU detector were

computed as  $Z=(x-\mu)/\sigma$ , where  $(\mu,\sigma)$  are the mean and variance for the output of frames 100-1100 in the baseline condition (warm water, no faked expressions). The mean difference in Z-score between the baseline and pain conditions was computed across the 26 subjects. Table 2 shows the action detectors with the largest difference in Z-scores. We observed that the actions with the largest Z-scores for genuine pain were Mouth opening and jaw drop (25 and 26), lip corner puller by zygomatic (12), nose wrinkle (9), and to a lesser extent, lip raise (10) and cheek raise (6). These facial actions have been previously associated with cold pressor pain (e.g. Prkachin, 1992; Craig & Patrick 1985).

The Z-score analysis was next repeated for faked versus baseline. We observed that in faked pain there was relatively more facial activity than in real pain. The facial action outputs with the highest z-scores for faked pain relative to baseline were brow lower (4), distress brow (1 or 1+4), inner brow raise (1), mouth open and jaw drop (25 and 26), cheek raise (6), lip raise (10), fear brow (1+2+4), nose wrinkle (9), mouth stretch (20), and lower lid raise (7).

Differences between real and faked pain were examined by computing the difference of the two z-scores. Differences were observed primarily in the outputs of action unit 4 (brow lower), as well as distress brow (1 or 1+4) and inner brow raise (1 in any combination).

**Table 2.** Z-score differences of the three pain conditions, averaged across subjects. FB: Fear brow 1+2+4. DB: Distress brow (1,1+4).

**A. Real Pain vs baseline:**

<u>Action Unit</u>	25	12	9	26	10	6
Z-score	1.4	1.4	1.3	1.2	0.9	0.9

**B. Faked Pain vs Baseline:**

<u>Action Unit</u>	4	DB	1	25	12	6	26	10	FB	9	20	7
Z-score	2.7	2.1	1.7	1.5	1.4	1.4	1.3	1.3	1.2	1.1	1.0	0.9

**C. Real Pain vs Faked Pain:**

<u>Action Unit</u>	4	DB	1
Z-score difference	1.8	1.7	1.0

**Table 3.** Individual subject differences between faked and genuine pain. Differences greater than 2 standard deviations are shown. F>P: Number of subjects in which the output for the given AU was greater in faked than genuine pain. P>F: Number of subjects for which the output was greater in genuine than faked pain. FB: Fear brow 1+2+4. DB: Distress brow (1,1+4).

<u>AU</u>	1	2	4	5	6	7	9	10	12	14	15	17	18	20	23	24	25	26	FB	DB
F>P	6	4	9	1	7	4	3	6	5	3	5	5	1	4	3	4	4	4	6	5
P>F	3	3	0	0	4	0	4	4	4	2	3	1	3	1	1	1	2	4	2	0

Individual subject differences between faked and real pain are shown in Table 3. Difference-of-Z-scores between the genuine and faked pain conditions were computed

for each subject and each AU. There was considerable variation among subjects in the difference between their faked and real pain expressions. However the most consistent finding is that 9 of the 26 subjects showed significantly more brow lowering activity (AU4) during the faked pain condition, whereas none of the subjects showed significantly more AU4 activity during the real pain condition. Also 7 subjects showed more cheek raise (AU6), and 6 subjects showed more inner brow raise (AU1), and the fear brow combination (1+2+4). The next most common differences were to show more 12, 15, 17, and distress brow (1 alone or 1+4) during faked pain.

Paired t-tests were conducted for each AU to assess whether it was a reliable indicator of genuine versus faked pain in a within-subjects design. Of the 20 actions tested, the difference was statistically significant for three actions. It was significant for AU 4 at  $p < .001$ , and marginally significant for AU 7 and distress brow at  $p < .05$ .

In order to characterize action unit combinations that relate to the difference between fake and real pain expressions, principal component analysis was conducted on the difference-of-Z-scores. The first eigenvector had the largest loading on distress brow and inner brow raise (AU 1). The second eigenvector had the largest loading on lip corner puller (12) and cheek raise (6) and was *lower* for fake pain expressions. The third eigenvector had the largest loading on brow lower (AU 4). Thus when analyzed singly, the action unit channel with the most information for discriminating fake from real pain was brow lower (AU 4). However when correlations were assessed through PCA, the largest variance was attributed to two combinations, and AU 4 accounted for the third most variance.

Overall, the outputs of the automated system showed similar patterns to previous studies of real and faked pain using manual FACS coding by human experts. Exaggerated activity of the brow lower (AU 4) during faked pain is consistent with previous studies in which the real pain condition was exacerbated lower back pain (Craig et al. 1991, Hill & Craig, 2002). Another study performed a FACS analysis of fake and real pain expressions with cold pressor pain, but with children ages 8-12 (LaRochette et al., 2006). This study observed significant elevation in the following AUs for fake pain relative to baseline: 1 4 6 7 10 12 20 23 25 26. This closely matches the AUs with the highest z-scores in the automated system output of the present study (Table 2B). LaRochette et al. did not measure AU 9 or the brow combinations. When faked pain expressions were compared with real cold pressor pain in children, LaRochette et al found significant differences in AU's 1 4 7 10. Again the findings of the present study using the automated system are similar, as the AU channels with the highest z-scores were 1, 4, and 1+4 (Table 2C), and the t-tests were significant for 4, 1+4 and 7.

### **3.3.2 Automatic Discrimination of Real from Fake Pain**

The above analysis examined which AU outputs contained information about genuine versus faked pain. We next turned to the problem of discriminating genuine from faked pain expressions in a subject-independent manner. If the task were to

simply detect the presence of a red-flag set of facial actions, then differentiating fake from real pain expressions would be relatively simple. However, it should be noted that subjects display actions such as AU 4, for example, in both real and fake pain, and the distinction is in the quantity of AU 4. Also, there is inter-subject variation in expressions of both real and fake pain, there may be combinatorial differences in the sets of actions displayed during real and fake pain, and the subjects may cluster. We therefore applied machine learning to the task of discriminating real from faked pain expressions in a subject-independent manner.

A second-layer classifier was trained to discriminate genuine pain from faked pain based on the 20-channel output stream. For this second-layer classification step, we explored SVMs, Adaboost, and linear discriminant analysis. Nonlinear SVMs with radial basis function kernels gave the best performance. System performance for generalization to novel subjects was assessed using leave-one-out cross-validation, in which all the data from one subject was reserved for testing on each trial.

Prior to learning, the system performed an automatic reliability estimate based on the smoothness of the eye positions, and those frames with low reliability were automatically excluded from training and testing the real pain / fake pain classifier. Those frames with abrupt shifts of 2 or more pixels in the returned eye positions were automatically detected and labeled unreliable. This tends to occur during eyeblinks with the current eye detector. However future versions of the eye detector will correct that issue. The reliability filter had a relatively small effect on performance. The analysis of Table 2 was repeated under this criterion, and the Z-scores improved by about 0.1. Note also that the reliability filter on the frames is not to be confused with dropping the difficult trials since a real pain / fake pain decision was always made for each subject.

The 60 second video from each condition was broken up into 6 overlapping segments of 500 frames. For each segment, the following 5 statistics were measured for each of the 20 AU's: median, maximum, range, first to third quartile difference, 90 to 100 percentile difference. Thus the input to the SVM for each segment contained 100 dimensions. Each cross-validation trial contained 300 training samples (25 subjects x 2 conditions x 6 segments).

A nonlinear SVM trained to discriminate posed from real facial expressions of pain obtained an area under the ROC of .72 for generalization to novel subjects. This was significantly higher than performance of naïve human subjects, who obtained a mean accuracy of 49% correct for discriminating faked from real pain on the same set of videos.

### **3.4 Discussion of pain study**

The field of automatic facial expression analysis has advanced to the point that we can begin to apply it to address research questions in behavioral science. Here we describe a pioneering effort to apply fully automated facial action coding to the problem of differentiating fake from real expressions of pain. While naïve human subjects were only at 49% accuracy for distinguishing fake from real pain, the automated system obtained .72 area under the ROC, which is equivalent to 72%

correct on a 2-alternative forced choice. Moreover, the pattern of results in terms of which facial actions may be involved in real pain, fake pain, and differentiating real from fake pain is similar to previous findings in the psychology literature using manual FACS coding.

Here we applied machine learning on a 20-channel output stream of facial action detectors. The machine learning was applied to samples of spontaneous expressions during the subject state in question. Here the state in question was fake versus real pain. The same approach can be applied to learn about other subject states, given a set of spontaneous expression samples. Section 4 develops another example in which this approach is applied to the detection of driver drowsiness from facial expression.

#### **4 Automatic Detection of Driver Fatigue**

The US National Highway Traffic Safety Administration estimates that in the US alone approximately 100,000 crashes each year are caused primarily by driver drowsiness or fatigue (Department of Transportation, 2001). Thus incorporating automatic driver fatigue detection mechanism into vehicles may help prevent many accidents.

One can use a number of different techniques for analyzing driver exhaustion. One set of techniques places sensors on standard vehicle components, e.g., steering wheel, gas pedal, and analyzes the signals sent by these sensors to detect drowsiness (Takei & Furukawa, 2005). It is important for such techniques to be adapted to the driver, since Abut and his colleagues note that there are noticeable differences among drivers in the way they use the gas pedal (Igarashi, et al., 2005).

A second set of techniques focuses on measurement of physiological signals such as heart rate, pulse rate, and Electroencephalography (EEG) (e.g. Cobb, 1983). It has been reported by re-searchers that as the alertness level decreases EEG power of the alpha and theta bands increases (Hung & Chung, 2005). Hence providing indicators of drowsiness. However this method has draw-backs in terms of practicality since it requires a person to wear an EEG cap while driving.

A third set of solutions focuses on computer vision systems that can detect and recognize the facial motion and appearance changes occurring during drowsiness (Gu & Ji, 2004; Zhang & Zhang, 2006). The advantage of computer vision techniques is that they are non-invasive, and thus are more amenable to use by the general public. There are some significant previous studies about drowsiness detection using computer vision techniques. Most of the published research on computer vision approaches to detection of fatigue has focused on the analysis of blinks and head movements. However the effect of drowsiness on other facial expressions have not been studied thoroughly. Recently Gu & Ji presented one of the first fatigue studies that incorporates certain facial expressions other than blinks. Their study feeds action unit information as an input to a dynamic Bayesian network. The network was trained on subjects posing a state of fatigue (Gu, Zhang & Ji, 2005). The video segments were classified into three stages: inattention, yawn, or falling asleep. For predicting falling asleep, head nods, blinks, nose wrinkles and eyelid tighteners were used.

Previous approaches to drowsiness detection primarily make pre-assumptions about the relevant behavior, focusing on blink rate, eye closure, and yawning. Here

we employ machine learning methods to data mine actual human behavior during drowsiness episodes. The objective of this study is to discover what facial configurations are predictors of fatigue. In this study, facial motion was analyzed automatically from video using a fully automated facial expression analysis system based on the Facial Action Coding System (FACS) (Ekman & Friesen, 1978). In addition to the output of the automatic FACS recognition system we also collected head motion data using an accelerometer placed on the subject's head, as well as steering wheel data.

#### 4.1 Driving task

Subjects played a driving video game on a windows machine using a steering wheel<sup>1</sup> and an open source multi- platform video game<sup>2</sup> (See Figure 5). At random times, a wind effect was applied that dragged the car to the right or left, forcing the subject to correct the position of the car. This type of manipulation had been found in the past to increase fatigue (Orden, Jung & Makeig, 2000). Driving speed was held constant. Four subjects performed the driving task over a three hour period beginning at midnight. During this time subjects fell asleep multiple times thus crashing their vehicles. Episodes in which the car left the road (crash) were recorded. Video of the subjects face was recorded using a DV camera for the entire 3 hour session.

In addition to measuring facial expressions with CERT, head movement was measured using an accelerometer that has 3 degrees of freedom. This three dimensional accelerometer<sup>3</sup> has three one dimensional accelerometers mounted at right angles measuring accelerations in the range of 5g to +5g where g represents earth gravitational force.



**Fig. 5.** Driving Simulation Task.

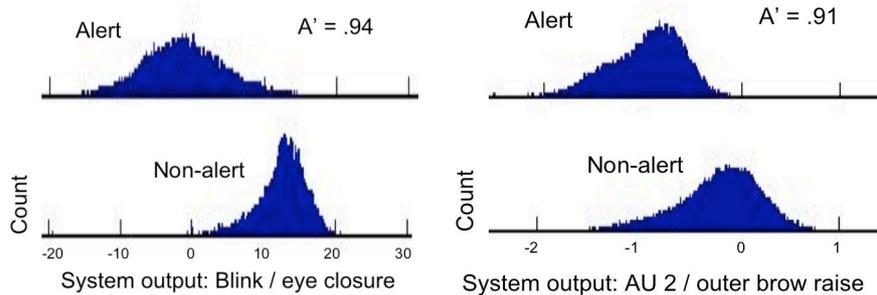
---

<sup>1</sup>Thrustmaster<sup>®</sup>Ferrari Racing Wheel

<sup>2</sup>The Open Racing Car Simulator (TORCS)

## 4.2 Facial Actions Associated with Driver Fatigue

Subject data was partitioned into drowsy (non-alert) and alert states as follows. The one minute preceding a sleep episode or a crash was identified as a non-alert state. There was a mean of 24 non-alert episodes with a minimum of 9 and a maximum of 35. Fourteen alert segments for each subject were collected from the first 20 minutes of the driving task.



**Fig. 6.** Example histograms for blink and Action Unit 2 in alert and non-alert states for one subject. A' is area under the ROC.

The output of the facial action detector consisted of a continuous value for each facial action and each video frame which was the distance to the separating hyperplane, i.e., the margin. Histograms for two of the action units in alert and non-alert states are shown in Figure 6. The area under the ROC (A') was computed for the outputs of each facial action detector to see to what degree the alert and non-alert output distributions were separated.

In order to understand how each action unit is associated with drowsiness across different subjects, Multinomial Logistic Ridge Regression (MLR) was trained on each facial action individually. Examination of the A' for each action unit reveals the degree to which each facial movement was able to predict drowsiness in this study. The A's for the drowsy and alert states are shown in Table 4. The five facial actions that were the most predictive of drowsiness by increasing in drowsy states were 45, 2 (outer brow raise), 15 (frown), 17 (chin raise), and 9 (nose wrinkle). The five actions that were the most predictive of drowsiness by decreasing in drowsy states were 12 (smile), 7 (lid tighten), 39 (nostril compress), 4 (brow lower), and 26 (jaw drop). The high predictive ability of the blink/eye closure measure was expected. However the predictability of the outer brow raise (AU 2) was previously unknown.

We observed during this study that many subjects raised their eyebrows in an attempt to keep their eyes open, and the strong association of the AU 2 detector is consistent with that observation. Also of note is that action 26, jaw drop, which occurs during yawning, actually occurred less often in the critical 60 seconds prior to a crash. This is consistent with the prediction that yawning does not tend to occur in the final moments before falling asleep.

**Table 4.** MLR model for predicting drowsiness across subjects. Predictive performance of each facial action individually is shown.

**More when critically drowsy**

AU	Name	A'
45	Blink/Eye Closure	0.94
2	Outer Brow Raise	0.81
15	Lip Corner Depressor	0.80
17	Chin Raiser	0.79
9	Nose Wrinkle	0.78
30	Jaw Sideways	0.76
20	Lip stretch	0.74
11	Nasolabial Furrow	0.71
14	Dimpler	0.71
1	Inner Brow Raise	0.68
10	Upper Lip Raise	0.67
27	Mouth Stretch	0.66
18	Lip Pucker	0.66
22	Lip funneler	0.64
24	Lip presser	0.64
19	Tongue show	0.61

**Less when critically drowsy**

AU	Name	A'
12	Smile	0.87
7	Lid tighten	0.86
39	Nostril Compress	0.79
4	Brow lower	0.79
26	Jaw Drop	0.77
6	Cheek Raise	0.73
38	Nostril Dilate	0.72
23	Lip tighten	0.67
8	Lips toward	0.67
5	Upper lid raise	0.65
16	Upper lip depress	0.64
32	Bite	0.63

**4.3 Automatic Detection of Driver Fatigue**

The ability to predict drowsiness in novel subjects from the facial action code was then tested by running MLR on the full set of facial action outputs. Prediction performance was tested by using a leave-one-out cross validation procedure, in which one subjects' data was withheld from the MLR training and retained for testing, and the test was repeated for each subject. The data for each subject by facial action was first normalized to zero-mean and unit standard deviation. The MLR output for each AU feature was summed over a temporal window of 12 seconds (360 frames) before computing A'.

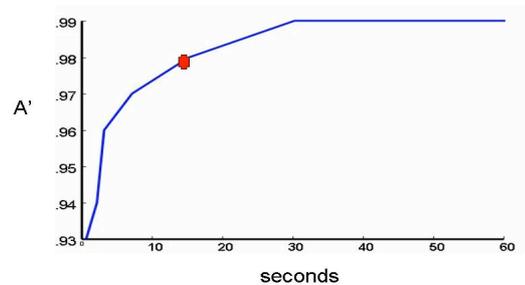
MLR trained on all AU features obtained an A' of .90 for predicting drowsiness in novel subjects. Because prediction accuracy may be enhanced by feature selection, in which only the AU's with the most information for discriminating drowsiness are included in the regression, a second MLR was trained by contingent feature selection, starting with the most discriminative feature (AU 45), and then iteratively adding the next most discriminative feature given the features already selected. These features are shown on Table 5. Best performance of .98 was obtained with five features: 45, 2, 19 (tongue show), 26 (jaw drop), and 15. This five feature model outperformed the MLR trained on all features.

**Effect of Temporal Window Length.** We next examined the effect of the size of the temporal window on performance. The five feature model was employed for this analysis. The performances shown in Table 5 employed a temporal window of 12 seconds. Here, the MLR output in the 5 feature model was summed over windows of N seconds, where N ranged from 0.5 to 60 seconds. Figure 7 shows the area under the

ROC for drowsiness detection in novel subjects over time periods. Performance saturates at about 0.99 as the window size exceeds 30 seconds. In other words, given a 30 second video segment the system can discriminate sleepy versus non-sleepy segments with 0.99 accuracy across subjects.

**Table 5.** Drowsiness detection performance for novel subjects, using an MLR classifier with different feature combinations. The weighted features are summed over 12 seconds before computing A'.

Feature	A'
AU45,AU2,AU19,AU26,AU15	.9792
All AU features	.8954



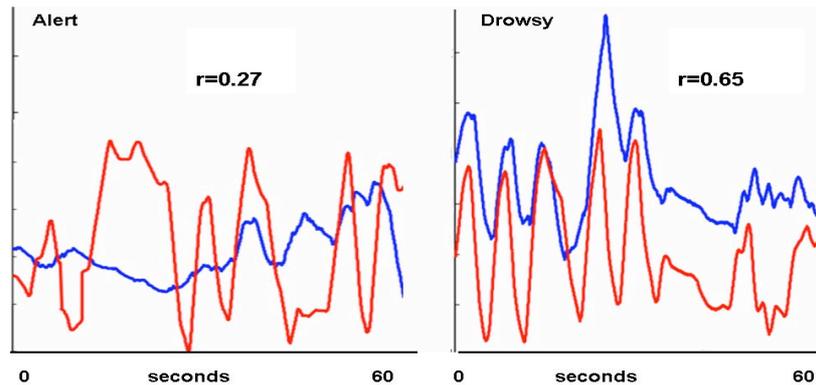
**Fig. 7.** Performance for drowsiness detection in novel subjects over temporal window sizes. Red dot indicates the priorly obtained performance for a temporal window of 12 seconds.

#### 4.4 Coupling of Behaviors

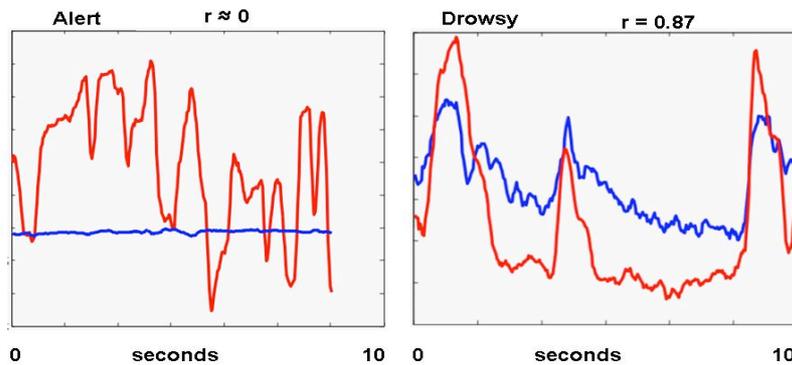
**Coupling of steering and head motion.** Observation of the subjects during drowsy and nondrowsy states indicated that the subjects head motion differed substantially when alert versus when the driver was about to fall asleep. Surprisingly, head motion increased as the driver became drowsy, with large roll motion coupled with the steering motion as the driver became drowsy. Just before falling asleep, the head would become still.

We also investigated the coupling of the head and arm motions. Correlations between head motion as measured by the roll dimension of the accelerometer output and the steering wheel motion are shown in Figure 8. For this subject (subject 2), the correlation between head motion and steering increased from 0.33 in the alert state to 0.71 in the non-alert state. For subject 1, the correlation between head motion and steering similarly increased from 0.24 in the alert state to 0.43 in the non-alert state. The other two subjects showed a smaller coupling effect. Future work includes combining the head motion measures and steering correlations with the facial movement measures in the predictive model.

**Coupling of eye openness and eyebrow raise.** We observed that for some of the subjects coupling between eye brow up's and eye openness increased in the drowsy state. In other words subjects tried to open their eyes using their eyebrows in an attempt to keep awake. See Figure 9.



**Fig. 8.** Head motion (blue/gray) and steering position (red/black) for 60 seconds in an alert state (left) and 60 seconds prior to a crash (right). Head motion is the output of the roll dimension of the accelerometer.



**Fig. 9.** Eye Openness (red/black) and Eye Brow Raises (AU2) (Blue/gray) for 10 seconds in an alert state (left) and 10 seconds prior to a crash (right).

#### 4.4 Conclusions of Driver Fatigue Study

This chapter presented a system for automatic detection of driver drowsiness from video. Previous approaches focused on assumptions about behaviors that might be predictive of drowsiness. Here, a system for automatically measuring facial expressions was employed to data mine spontaneous behavior during real drowsiness episodes. This is the first work to our knowledge to reveal significant associations between facial expression and fatigue beyond eyeblinks. The project also revealed a

potential association between head roll and driver drowsiness, and the coupling of head roll with steering motion during drowsiness. Of note is that a behavior that is often assumed to be predictive of drowsiness, yawn, was in fact a negative predictor of the 60-second window prior to a crash. It appears that in the moments before falling asleep, drivers yawn less, not more, often. This highlights the importance of using examples of fatigue and drowsiness conditions in which subjects actually fall sleep.

The computer vision field has advanced to the point that we are now able to begin to apply automatic facial expression recognition systems to important research questions in behavioral science. This chapter explored two such applications, in which the automated measurement system revealed information about facial expression that was previously unknown. While the accuracy of individual facial action detectors is still below that of human experts, automated systems can be applied to large quantities of video data. Statistical pattern recognition on this large quantity of data can reveal emergent behavioral patterns that previously would have required hundreds of coding hours by human experts, and would be unattainable by the non-expert. Moreover, automated facial expression analysis will enable investigations into facial expression dynamics that were previously intractable by human coding because of the time required to code intensity changes. Future work will explore facial expression dynamics.

### **Acknowledgements**

Support for this work was provided in part by NSF grants CNS-0454233, SBE-0542013, and NSF ADVANCE award 0340851, and by a grant from Turkish State Planning Organization. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation. Portions of the research in this paper use the MMI Facial Expression Database collected by M. Pantic & M.F. Valstar.

### **References**

1. Bartlett M.S., Littlewort G.C., Frank M.G., Lainscsek C., Fasel I., and Movellan J.R., (2006). Automatic recognition of facial actions in spontaneous expressions., *Journal of Multimedia.*, 1(6) p. 22-35.
2. Cobb. W. (1983). Recommendations for the practice of clinical neurophysiology. Elsevier.
3. Cohn, J.F. & Schmidt, K.L. (2004). The timing of facial motion in posed and spontaneous smiles. *J. Wavelets, Multi-resolution & Information Processing*, Vol. 2, No. 2, pp. 121-132.
4. Craig KD, Hyde S, Patrick CJ. (1991). Genuine, suppressed, and faked facial behaviour during exacerbation of chronic low back pain. *Pain* 46:161– 172.
5. Craig KD, Patrick CJ. (1985). Facial expression during induced pain. *J Pers Soc Psychol.* 48(4):1080-91.
6. Donato, G., Bartlett, M.S., Hager, J.C., Ekman, P. & Sejnowski, T.J. (1999). Classifying facial actions. *IEEE Trans. Pattern Analysis and Machine Intelligence*, Vol. 21, No. 10, pp. 974-989.

7. DOT (2001). Saving lives through advanced vehicle safety technology. USA Department of Transportation. <http://www.its.dot.gov/ivi/docs/AR2001.pdf>.
8. Ekman P. and Friesen, W. *Facial Action Coding System: A Technique for the Measurement of Facial Movement*, Consulting Psychologists Press, Palo Alto, CA, 1978.
9. Ekman, P. (2001). *Telling Lies: Clues to Deceit in the Marketplace, Politics, and Marriage*. W.W. Norton, New York, USA.
10. Ekman, P. & Rosenberg, E.L., (Eds.), (2005). *What the face reveals: Basic and applied studies of spontaneous expression using the FACS*, Oxford University Press, Oxford, UK.
11. Fasel I., Fortenberry B., Movellan J.R. "A generative framework for real-time object detection and classification.," *Computer Vision and Image Understanding* 98, 2005.
12. Fishbain DA, Cutler R, Rosomoff HL, Rosomoff RS. (1999). Chronic pain disability exaggeration/malingering and submaximal effort research. *Clin J Pain*. 15(4):244-74.
13. Fishbain DA, Cutler R, Rosomoff HL, Rosomoff RS. (2006). Accuracy of deception judgments. *Pers Soc Psychol Rev*. 10(3):214-34.
14. Frank MG, Ekman P, Friesen WV. (1993). Behavioral markers and recognizability of the smile of enjoyment. *J Pers Soc Psychol*. 64(1):83-93.
15. Grossman, S., Shielder, V., Swedeen, K., Mucenski, J. (1991). Correlation of patient and caregiver ratings of cancer pain. *Journal of Pain and Symptom Management* 6(2), p. 53-57.
16. Gu, H., Ji, Q. (2004). An automated face reader for fatigue detection. In: *FGR*. (2004) 111–116.
17. Gu, H., Zhang, Y., Ji, Q. (2005). Task oriented facial behavior recognition with selective sensing. *Comput. Vis. Image Underst.* 100(3) p. 385–415.
18. Hadjistavropoulos HD, Craig KD, Hadjistavropoulos T, Poole GD. (1996). Subjective judgments of deception in pain expression: accuracy and errors. *Pain*. 65(2-3):251-8.
19. Hill ML, Craig KD (2002) Detecting deception in pain expressions: the structure of genuine and deceptive facial displays. *Pain*. 98(1-2):135-44.
20. Hong, Chung, K. (2005). Electroencephalographic study of drowsiness in simulated driving with sleep deprivation. *International Journal of Industrial Ergonomics*. 35(4), Pages 307-320.
21. Igarashi, K., Takeda, K., Itakura, F., Abut, H.: (2005). *DSP for In-Vehicle and Mobile Systems*. Springer US
22. Kanade, T., Cohn, J.F. and Tian, Y., "Comprehensive database for facial expression analysis," in *Proceedings of the fourth IEEE International conference on automatic face and gesture recognition (FG'00)*, Grenoble, France, 2000, pp. 46–53.
23. Larochette AC, Chambers CT, Craig KD (2006). Genuine, suppressed and faked facial expressions of pain in children. *Pain*. 2006 Dec 15;126(1-3):64-71.
24. Littlewort, G., Bartlett, M.S., Fasel, I., Susskind, J. & Movellan, J. (2006). Dynamics of facial expression extracted automatically from video. *J. Image & Vision Computing*, Vol. 24, No. 6, pp. 615-625.
25. Morecraft RJ, Louie JL, Herrick JL, Stilwell-Morecraft KS. (2001). Cortical innervation of the facial nucleus in the non-human primate: a new interpretation of the effects of stroke and related subtotal brain trauma on the muscles of facial expression. *Brain* 124(Pt 1):176-208.
26. Orden, K.F.V., Jung, T.P., Makeig, S. (2000). Combined eye activity measures accurately estimate changes in sustained visual task performance. *Biological Psychology* 52(3):221-40.
27. Pantic, M., Pentland, A., Nijholt, A. & Huang, T. (2006). Human Computing and machine understanding of human behaviour: A Survey, *Proc. ACM Int'l Conf. Multimodal Interfaces*, pp. 239-248.
28. Pantic, M.F. Valstar, R. Rademaker and L. Maat, "Web-based Database for Facial Expression Analysis", *Proc. IEEE Int'l Conf. Multimedia and Expo (ICME'05)*, Amsterdam, The Netherlands, July 2005.

29. Prkachin KM. (1992). The consistency of facial expressions of pain: a comparison across modalities. *Pain*. 51(3):297-306.
30. Prkachin KM, Schultz I, Berkowitz J, Hughes E, Hunt D. Assessing pain behaviour of low-back pain patients in real time: concurrent validity and examiner sensitivity. *Behav Res Ther*. 40(5):595-607.
31. Rinn WE. The neuropsychology of facial expression: a review of the neurological and psychological mechanisms for producing facial expression. *Psychol Bull* 95:52-77.
32. Schmand B, Lindeboom J, Schagen S, Heijt R, Koene T, Hamburger HL. Cognitive complaints in patients after whiplash injury: the impact of malingering. *J Neurol Neurosurg Psychiatry*. 64(3):339-43.
33. Schmidt KL, Cohn JF, Tian Y. (2003). Signal characteristics of spontaneous facial expressions: automatic movement in solitary and social smiles. *Biol Psychol*. 65(1):49-66.
34. Schneiderman, H. and Kanade, T. (1998). Probabilistic Modeling of Local Appearance and Spatial Relationships for Object Recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 45-51.
35. Takei, Y. Furukawa, Y. (2005): Estimate of driver's fatigue through steering motion. In: *Man and Cybernetics, 2005 IEEE International Conference*. Volume: 2, pg. 1765- 1770.
36. Viola, P. & Jones, M. (2004). Robust real-time face detection. *J. Computer Vision*, Vol. 57, No. 2, pp. 137-154.
37. Vural, E., Ercil, A., Littlewort, G.C., Bartlett, M.S., and Movellan, J.R. (2007). Machine learning systems for detecting driver drowsiness. *Proceedings of the Biennial Conference on Digital Signal Processing for in-Vehicle and Mobile Systems*.
38. Zhang, Z., shu Zhang, J. (2006). Driver fatigue detection based intelligent vehicle control. In: *Proceedings of the 18th International Conference on Pattern Recognition*, Washington, DC, USA, IEEE Computer Society p. 1262-1265.