

Spontaneous vs. Posed smiles - can we tell the difference?

Bappaditya Mandal and Nizar Ouarti

Visual Computing Department, Institute for Infocomm Research, Singapore
Email address: bmandal@i2r.a-star.edu.sg (Bappaditya Mandal);
nizarouarti@gmail.com (Nizar Ouarti)

Abstract. Smile is an irrefutable expression that shows the physical state of the mind in both true and deceptive ways. Generally, it shows happy state of the mind, however, ‘smiles’ can be deceptive, for example people can give a smile when they feel happy and sometimes they might also give a smile (in a different way) when they feel pity for others. This work aims to distinguish spontaneous (felt) smile expressions from posed (deliberate) smiles by extracting and analyzing both global (macro) motion of the face and subtle (micro) changes in the facial expression features through both tracking a series of facial fiducial markers as well as using dense optical flow. Specifically the eyes and lips features are captured and used for analysis. It aims to automatically classify all smiles into either ‘spontaneous’ or ‘posed’ categories, by using support vector machines (SVM). Experimental results on large UvA-NEMO smile database show promising results as compared to other relevant methods.

Keywords: Posed, spontaneous smiles, feature extraction, face analysis.

1 Introduction

People believe that human face is the mirror/screen showing internal emotional state of the human body as and when it responds to the external world. This means that, what an individual thinks, feels or understands, etc, deep inside the brain, get imitated into the outside world through its face [7]. Facial smile expression undeniably plays a huge and pivotal role [25, 1, 11] in understanding social interactions within a community. People often give smile imitating the internal state of the body. For example, generally, people smile when they are happy or when sudden humorous things happen/appear in front of them. However, people are sometimes forced to pose smile because of the outside pressure or external factors. For example, people would pose a smile even when they don’t understand the joke or the humor. Sometimes people would also pose a smile even when they are reluctantly or unwillingly do or perform something in front of their bosses/peers [6].

Therefore being able to identify the type of smiles of individuals would give affective computing a deeper understanding of the human interactions. A large amount of research in psychology and neuroscience studying facial behavior

demonstrate that spontaneous deliberately displayed facial behavior has differences both in utilized facial muscles and their dynamics as compared to posed ones [8]. For example, spontaneous smiles are smaller in amplitude, longer in duration, slower in onset and offset times than posed smiles [3, 8, 22]. For humans, capturing such subtle facial movements is difficult and we often fail to distinguish between them. It is not surprising that in computer vision, algorithms developed for classifying such smiles usually fail to generalize to the subtlety and complexity of human posed and spontaneous affective behaviors [25, 15].

Numerous researchers asserted that dynamic features such as duration and speed of the smile play a part in differentiating the nature of the smile [11]. A spontaneous smile usually takes longer time to reach from onset to apex and then offset as compared to a posed smile [5]. As for non-dynamic features, the aperture size of the eyes is found to be a useful clue and is generally of a higher value when extracted from a spontaneous smile as compared to a posed one. On the other hand, the symmetry in (or the lack of) movement of spontaneous and posed smiles do not produce significant distinction in identifying them and is therefore not much useful [21]. In [22] a multi-modal system using geometric features such as shoulder, head and inner facial movements are fused together and GentleSVM-sigmoid is used to classify the posed and spontaneous smiles. He *et al.* in [10] proposed a technique for feature extraction and compared the performance using geometric and facial appearance features. Appearance based features are computed by recording statistics of overall pixel values of the image, or even using edge detection algorithm such as Gabor Wavelet Filter. Their comprehensive study shows that geometric features are generally more effective in detecting posed from spontaneous expressions [10].

A spatiotemporal method involving both natural and infrared face videos to distinguish posed and spontaneous expressions is proposed in [20]. Using temporal space and image sequences as volume, they extended the complete local binary patterns texture based descriptor into the spatiotemporal features to classify posed and spontaneous smiles. Dibeklioglu *et al.* in [4] used the dynamics of eyelid movements, distance measures and angular features in the changes of the eye aperture. Using several classifiers they have shown the superiority of eyelid movements over the eyebrows, cheek and lip movements for smile classification. Later in [5], they used dynamic characteristics of eyelid, cheek and lip corner movements for classifying posed and spontaneous smiles. Temporal facial information is obtained in [13] through segmenting the facial expression into onset, apex and offset which cover the entire duration of the smile. They reported good classification performance by using a combination of features extracted from the different phases.

The block diagram of our proposed method is shown in Fig. 1. Given smile video sequences of various subjects, we apply the facial features detection and tracking of the fiducial points over the entire smile video clip. Using D-markers, 25 important parameters (like duration, amplitude, speed acceleration, etc) are extracted from two important regions of the face: eyes and lips. Smile discriminative features are extracted using dense optical flow along the temporal domain

from the global (macro) motion and local (micro) motion of the face. All these information are fused and support vector machine (SVM) is then used as a classifier on these parameters to distinguish posed and spontaneous smiles.

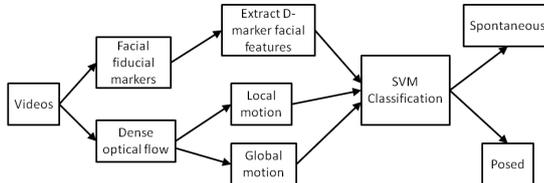


Fig. 1: Block diagram of the proposed system.

2 Feature Extraction from Various Face Components

We use the facial tracking algorithm developed by Nguyen *et al.* in [17] to obtain the fiducial points on the face. The 21 tracking markers each are labeled and placed following the convention as shown in Fig. 2 (a). The markers are manually annotated in the first frame of each video by user input and thereafter it automatically tracks the remaining frames of the smile video, it is of good accuracy and precision as compared to other facial tracking software [2]. The markers are placed on important facial feature points such as eyelids and corner of the lips for each subject. The convention followed in our approach for selecting fiducial markers are shown in Fig. 2 (a).

2.1 Face Normalization

To reduce inaccuracy due to the subject's head motion in the video that can cause change in angle with respect to roll, yaw and pitch rotations, we use the face normalization procedure described in [5]. Let l_i represents each of the feature points used to align the faces as shown in Fig. 2. Three non-collinear points (eye centers and nose tip) are used to form a plane ρ . Eye centers are defined as $c_1 = \frac{l_1+l_3}{2}$ and $c_2 = \frac{l_4+l_6}{2}$. Angles between the positive normal vector N_ρ of ρ and unit vectors U on X (horizontal), Y (vertical), and Z (perpendicular) axes give the relative head pose as follows:

$$\theta = \arccos \frac{U \cdot N_\rho}{\|U\| \|N_\rho\|}, \text{ where } N = \overrightarrow{l_g c_2} \times \overrightarrow{l_g c_1}. \quad (1)$$

$\overrightarrow{l_g c_2}$ and $\overrightarrow{l_g c_1}$ denote the vectors from point l_g to points c_2 and c_1 , respectively. $\|U\|$ and $\|N_\rho\|$ represents the magnitudes of U and N_ρ vectors respectively. Using the human face configuration, (1) can estimate the exact roll (θ_z) and yaw (θ_y) angles of the face with respect to the camera. If we start with the frontal face, the pitch angles (θ'_x) can be computed by subtracting the initial value. Using the estimated head pose, tracked fiducial points are normalized with respect to rotation, scale and translation as follows:

$$l'_i = [l_i - \frac{c_1 + c_2}{2}] R_x(-\theta'_x) R_y(-\theta_y) R_z(-\theta_z) \frac{100}{\epsilon(c_1 + c_2)}, \quad (2)$$

where l'_i is the aligned point. R_x , R_y and R_z denote the 3D rotation matrices for the given angles. $\epsilon()$ is the Euclidean distance measure. Essentially (1) constructs a normal vector perpendicular to the plane of the face using three points (nose tip and eye centers), then calculate the angle formed between X , Y and Z axis with regards to the normal vector of face plane. Thereafter, (2) process and normalize each and every point of the frame accordingly and set the interocular distance to 100 pixels with the middle point acting as the new origin of the face center.

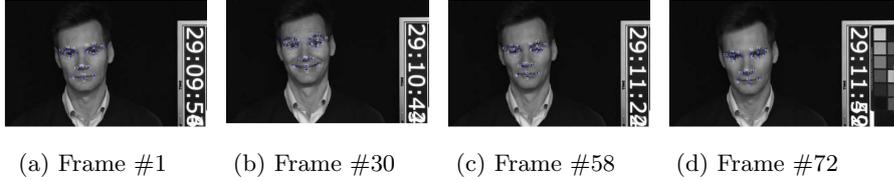


Fig. 2: (a) Shows the tracked points on the 1st frame, (b) shows the tracked points on 30th frame, (c) shows the tracked points on 58th frame and (d) shows the tracked points on 72nd frame on one subject. (Best viewed when zoomed in.)

2.2 D-Marker Facial Features

In the first part of our strategy, we focus on extracting the subject's eyelid and lips features. We first construct a amplitude signal variable based on the facial feature markers on the eyelid regions. We compute the amplitude of eyelid and lip end movements during a smile using the procedure described in [21]. Eyelid amplitude signals are computed using the eyelid aperture (D_{eyelid}) displacement at time t , given by:

$$D_{eyelid}(t) = \frac{\kappa(\frac{l_1^t+l_3^t}{2}, l_2^t)\epsilon(\frac{l_1^t+l_3^t}{2}, l_2^t) + \kappa(\frac{l_4^t+l_6^t}{2}, l_5^t)\epsilon(\frac{l_4^t+l_6^t}{2}, l_5^t)}{2\epsilon(l_1^t, l_3^t)} \quad (3)$$

where $\kappa(l_i, l_j)$ denotes the relative vertical location function, which equals to -1 if l_j is located (vertically) below l_i on the face, and 1 otherwise. The equation above uses the markers for eyelids namely 1-6 as shown in Fig. 2, to construct the amplitude signal that calculate the eyelid aperture size in each frame t . The amplitude signal D_{eyelid} is then further computed to obtain a series of features. In addition to the amplitudes, speed and acceleration signal are also extracted by computing the second derivatives of the amplitudes.

Smile amplitude is estimated as the mean amplitude of right and left lip corners, normalized by the length of the lip. Let $D_{lip}(t)$ be the value of the mean amplitude signal of the lip corners in the frame t . It is estimated as

$$D_{lip}(t) = \frac{\epsilon(\frac{l_{10}^t+l_{11}^t}{2}, l_{10}^t) + \epsilon(\frac{l_{10}^t+l_{11}^t}{2}, l_{11}^t)}{2\epsilon(l_{10}^t, l_{11}^t)} \quad (4)$$

where l_i^t denotes the 2D location of the i^{th} point in frame t . For each video of our subject we are able to acquire a 25-dimensional feature vectors based on

the eyelids markers and lip corner points. Onset phase is defined as the longest continuous increase in D_{lip} . Similarly, the offset phase is detected as the longest continuous decrease in D_{lip} . Apex is defined as the phase between the last frame of the onset and the first frame of the offset. The displacement signals of eyelids and lip corners could then be calculated using the tracked points. Onset, apex and offset phases of the smile are estimated using the maximum continuous increase and decrease of the mean displacement of the eyelids and lip corners. The D-Marker is then able to extract 25 descriptive features each for eyelids and lip corner, so a vector of 50 features are obtained from each frame (using two frames at a time). The features are then concatenated and passed through SVM for training and classification.

2.3 Features from Dense Optical Flow

In the second phase of the feature extraction, we use our own proposed dense optical flow [19] for capturing both global and local motions appearing in the smile videos. Our approach is divided into four distinct stages that are fully automatic and does not require any human intervention. The first step is to detect each frame in which the face is present. We use our previously developed face, integration of sketch and graph patterns (ISG) eyes and mouth detectors for face recognition on wearable devices and human-robot-interaction [14, 23]. So we get the region of interest (ROI) for the face (as shown in Fig. 3, left, yellow ROI) with 100% accuracy on the entire UvA-NEMO smile database [5]. In the second step, we determine the area corresponding to the right eye, left eye in red ROI and mouth in blue ROI for which we get 96.9% accuracy on the entire database.

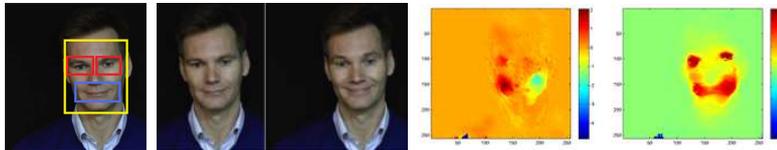


Fig. 3: Left: Face, eyes and mouth detections. Yellow ROI for face detection, red ROI for eyes detection and blue ROI for mouth detection. Middle: Two consecutive frames of a subject’s smile video and Right: their optical flows in x- and y-directions. (Best viewed in color and zoomed in.)

In the third step, the optical flow is computed between the image at time t and at time $t+1$ of the video sequence (see Fig. 3, middle). The two components of the optical flow are illustrated in Fig. 3, right, which shows the optical flow along the x-axis and the optical flow along the y-axis. Because we are using a dense optical flow algorithm, the time to process one picture is relatively important. To speed up the processing, we computed the optical flow only in the three ROI regions: right eye, left eye and mouth. The optical flow computed in our approach is a pyramidal differential dense algorithm that is based on the following constraint:

$$F = F_{smooth} + \beta F_{attach}, \quad (5)$$

where the *attach* term is based on thresholding method [24] and the regularization term (*smooth*) is based on the method developed by Meyer in [16], β is a weight controlling the ratio between the end attachment and the term control. Ouarti *et al.* in [19] proposed to use a regularization that do not use an usual wavelet but a non-stationary wavelet packet [18], which generalize the concept of wavelet for extracting optical flow information. We extend this idea for extracting fine grained information for both micro and macro motion variations in smile videos as shown in Fig. 4. Fig. 5 shows the dense optical flows with spontaneous and posed smiles variations. In the fourth step, for each of the three ROIs, the



Fig. 4: Original images and their dense optical flows with their corresponding micro and macro motion variations of a subject. (Best viewed in color and zoomed in.)

median of the optical flow is determined that give a cue to the global motion of the area. An histogram is computed based on the optical flow that has 10 bins. The top three bins in term of cardinality are kept among all the bins. A linear regression is then applied to find the major axis of the point group for each of the three bins determined. In the end, for each ROI we obtain: the median value of the bin 1, the value of the bin 2 and the value of the bin 3. It also calculates the intercept and slope for points of bins 1, 2 and 3. These result in 60 features for each frame (using two consecutive frames in a smile video). SVM is then used on these features to classify the posed and spontaneous smiles.



Fig. 5: Original images and their dense optical flows with their corresponding spontaneous and posed smiles variations of a subject. (Best viewed in color and zoomed in.)

The major advantage of this approach is that we can obtain useful smile discriminative features using a fully automatic analysis of videos, no marker are needed to be annotated by an operator/user. Moreover, rather than attempting to classify raw optical flow we design some processing to obtain a sparse representation of the optical flow signal. This representation helps in classification by extracting only the useful information in low dimensions and speeds up the calculation of the SVM. Finally, information is not completely connected to the

positioning of the different ROI knowing that this positioning may vary from one frame to another, it is dependent on the depth and highly variable depending on the individuals. Therefore a treatment which would be too closely related to the choice of the ROI would lead to non-consistent results.

3 Experimental Results

We test our proposed algorithm on UvA-NEMO Smile Database [5], it is the largest and most extensive smile (both posed and spontaneous) database with videos from a total of 400 subjects, (185 female, 215 male) aged between 8 to 76 years old, giving us a total of 1240 individual videos. Each video consists of a short segment of 3-8 seconds. The videos are extracted into frames at 50 frames per second. The extracted frames are also converted to gray scale and downsized to 480×270 . In all the experiments, we split the database, in which 80% is used as training samples and the remaining 20% is used as testing samples. Binary classifier SVM with radial basis function as the kernel and default parameters as in LIBSVM [12], is used to form a hyperplane based on the training samples. When a new testing sample is passed into the SVM it uses the hyperplane to determine which class the new sample falls under. This process is then repeated 5 times using a 5-fold cross validation method. To measure the subtle differences in the spontaneous and posed smiles we compute the confusion matrices between the two smiles so as to find out how much accuracy we can obtain in using each of them in the actual and classified separately. The results from all 5 processes are averaged and shown in Tables 1-5 and compared with other methods in Table 6.

3.1 Results using parameters from the facial components

Tables 1 and in bracket (\cdot) show the accuracy rates in distinguishing spontaneous smiles from the posed ones using eyes and lips features respectively. The results show that the eye features play very crucial role in finding the posed smiles where as the lips features are important for spontaneous smiles. Overall we could obtain an accuracy of 71.14% and 73.44% using eyes and lips features respectively. Table 2 shows the classification performance using combined features from eyes and lips. It is evident from the table that using these facial component features, pose smile can be classified better as compared to the spontaneous ones.

Actual \ Classified	Classified	
	Spontaneous	Posed
Spontaneous	60.1 (67.5)	39.9 (32.5)
Posed	17.5 (20.4)	82.5 (79.6)

Table 1: The overall accuracy (%) in classifying spontaneous and posed smiles using only the eyes features is 71.14%. In bracket (\cdot) shows accuracy using only the lips features as 73.44%.

Actual \ Classified	Classified	
	Spontaneous	Posed
Spontaneous	65.3	34.7
Posed	16.3	83.7

Table 2: The overall accuracy (%) in classifying spontaneous and posed smiles using the combined features from eyes and lips is 74.68%. (rows are gallery, columns are testing)

3.2 Results using Dense Optical flow

We use the features using dense optical flow as described in Section 2.3, the movement in both X- and Y-directions are recorded between every consecutive frames of each video. The confusion matrices are shown in Tables 3, in bracket (·) and 4. It can be see from the tables that the performance of optical flow is lower as compared to the component based approach. However, the facial component based feature extraction method requires user initialization to find and track fiducial points, whereas the dense optical flow features are fully automatic. It does not require any user intervention, so it is more useful for practical applications like first-person-views (FPV) or egocentric views on wearable devices like Google Glass for improving real-time social interactions [14, 9].

Actual \ Classified	Spontaneous	Posed
Spontaneous	57.8 (58.3)	42.2 (41.7)
Posed	39.8 (30.8)	60.2 (69.2)

Table 3: The accuracy (%) in classifying spontaneous and posed smiles using our proposed X-directions dense optical flow is 59%. In bracket (·) the accuracy using our proposed Y-directions is 63.8%.

Actual \ Classified	Spontaneous	Posed
Spontaneous	58.0	42.0
Posed	45.1	54.9

Table 4: The accuracy (%) in classifying spontaneous and posed smiles using our proposed fully automatic system using X- and Y-directions of dense optical flow is 56.6%.

3.3 Results using both Component based features and Dense Optical Flow

We combine all the features obtained from facial component based parameters and dense optical flow in to a single vector and apply SVM. Table 5 shows the confusion matrix using spontaneous and posed smiles. It can be seen that the performance of spontaneous smiles classification improved using features from dense optical flow. The experimental results in Table 5 show that both features from facial components and dense optical flows are important for improving the overall accuracy. Features from facial components (as shown in Table 2) are useful for encoding information arising from the muscle artifacts within a face, however, the regularized dense optical flow features helps in encoding fine grained information for both micro and macro motion variations in face smile videos. So combining them the overall accuracy has been improved.

Actual \ Classified	Spontaneous	Posed
Spontaneous	83.6	16.4
Posed	22.9	77.1

Table 5: The accuracy (%) in classifying spontaneous and posed smiles using our proposed fused approach comprising of both features from facial components and dense optical flow is 80.4%.

3.4 Comparison with Other Methods

Correct classification rates (%) using various methods on UvA-NEMO are shown in Table 6. It is evident from the table that our proposed approach is quite competitive as compared to the other state-of-the-arts methodologies.

Method	Correct Classification Rate (%)
Pfister <i>et al.</i> [20]	73.1
Dibeklioglu <i>et al.</i> [4]	71.1
Cohn & Schmidt [21]	77.3
Eyelid Features [5]	85.7
Mid-level Fusion (voting) [5]	87.0
<i>Ours Eye+Lips+dense optical flow</i>	80.4

Table 6: Correct classification rates (%) on UvA-NEMO database.

4 Conclusions

Differentiating spontaneous smiles from the posed ones is a challenging problem as it involves extracting subtle minute facial features and learning them. In this work we have analysed features extracted from facial component based parameters using fiducial points markers and tracking them. We have also obtained fully automatic features from dense optical flow on both eyes and mouth patches. It has been shown that the facial component based parameters give higher accuracy as compared to dense optical flow features for smile classification. However, the former requires initialization of the fiducial markers on the first frame and hence, it is not fully automatic. Dense optical flow has advantage that the features can be obtained without any manual intervention. Combining the facial components parameters and dense optical flow gives us highest accuracy for classifying the spontaneous and posed smiles. Experimental results on the largest UvA-NEMO smile database shows the efficacy of our proposed approach as compared to other state-of-the-arts methods.

References

1. Ambadar, Z., Cohn, J., Reed, L.: All smiles are not created equal: Morphology and timing of smiles perceived as amused, polite, and embarrassed/nervous. *Journal of Nonverbal Behavior* 33, 17–34 (2009)
2. Asthana, A., Zafeiriou, S., Cheng, S., Pantic, M.: Incremental face alignment in the wild. In: CVPR. Columbus, Ohio, USA (2014)
3. Cohn, J., Schmidt, K.: The timing of facial motion in posed and spontaneous smiles. *Intl J. Wavelets, Multiresolution and Information Processing* 2, 1–12 (2004)
4. Dibeklioglu, H., Valenti, R., Salah, A., Gevers, T.: Eyes do not lie: Spontaneous versus posed smiles. In: ACM Multimedia. pp. 703–706 (2010)
5. Dibeklioglu, H., Salah, A.A., Gevers, T.: Are you really smiling at me? spontaneous versus posed enjoyment smiles. In: IEEE ECCV. pp. 525–538 (2012)
6. Ekman, P.: Telling lies: Cues to deceit in the marketplace, politics, and marriage. WW. Norton & Company, New York (1992)

7. Ekman, P., Hager, J., Friesen, W.: The symmetry of emotional and deliberate facial actions. *Psychophysiology* 18, 101–106 (1981)
8. Ekman, P., Rosenberg, E.: *What the Face Reveals: Basic and Applied Studies of Spontaneous Expression Using the Facial Action Coding System*. Second ed. Oxford Univ. Press (2005)
9. Gan, T., Wong, Y., Mandal, B., Chandrasekhar, V., Kankanhalli, M.: Multi-sensor self-quantification of presentations. In: *ACM Multimedia*. pp. 601–610. Brisbane, Australia (Oct 2015)
10. He, M., Wang, S., Liu, Z., Chen, X.: Analyses of the differences between posed and spontaneous facial expressions. *Humaine Association Conference on Affective Computing and Intelligent Interaction* pp. 79–84 (2013)
11. Hoque, M., McDuff, D., Picard, R.: Exploring temporal patterns in classifying frustrated and delighted smiles. *IEEE Trans. Affective Computing* 3, 323–334 (2012)
12. Hsu, C., Chang, C., Lin, C.: *A practical guide to support vector classification* (2010)
13. Huijser, M., Gevers, T.: The influence of temporal facial information on the classification of posed and spontaneous enjoyment smiles. Tech. rep., Univ. of Amsterdam (2014)
14. Mandal, B., Ching, S., Li, L., Chandrasekha, V., Tan, C., Lim, J.H.: A wearable face recognition system on google glass for assisting social interactions. In: *3rd International Workshop on Intelligent Mobile and Egocentric Vision, ACCV*. pp. 419–433 (Nov 2014)
15. Mandal, B., Eng, H.L.: Regularized discriminant analysis for holistic human activity recognition. *IEEE Intelligent Systems* 27(1), 21–31 (2012)
16. Meyer, Y.: Oscillating patterns in image processing and in some nonlinear evolution equations. *The Fifteenth Dean Jacqueline B. Lewis Memorial Lectures*, American Mathematical Society (2001)
17. Nguyen, T., Ranganath, S.: Tracking facial features under occlusions and recognizing facial expressions in sign language. In: *International Conference on Automatic Face & Gesture Recognition*. vol. 6, pp. 1–7 (2008)
18. Ouarti, N., Peyre, G.: Best basis denoising with non-stationary wavelet packets. In: *International Conference on Image Processing*. vol. 6, pp. 3825–3828 (2009)
19. Ouarti, N., SAFRAN, A., LE, B., PINEAU, S.: Method for highlighting at least one moving element in a scene, and portable augmented reality (Aug 22 2013), <http://www.google.com/patents/WO2013121052A1?cl=en>, wO Patent App. PCT/EP2013/053,216
20. Pfister, T., Li, X., Zhao, G., Pietikainen, M.: Differentiating spontaneous from posed facial expressions within a generic facial expression recognition framework. In: *ICCV Workshop*. pp. 868–875 (2011)
21. Schmidt, K., Bhattacharya, S., Denlinger, R.: Comparison of deliberate and spontaneous facial movement in smiles and eyebrow raises. *Journal of Nonverbal Behavior* 33, 35–45 (2009)
22. Valstar, M., Pantic, M.: How to distinguish posed from spontaneous smiles using geometric features. In: *In Proceedings of ACM ICML*. pp. 38–45 (2007)
23. Yu, X., Han, W., Li, L., Shi, J., Wang, G.: An eye detection and localization system for natural human and robot interaction without face detection. *TAROS* pp. 54–65 (2011)
24. Zach, C., Pock, T., Bischof, H.: A duality based approach for realtime tv-l1 optical flow. In: *In Ann. Symp. German Association Patt. Recogn.* pp. 214–223 (2007)
25. Zeng, Z., Pantic, M., Roisman, G.I., Huang, T.S.: A survey of affect recognition methods: Audio, visual, and spontaneous expressions. *PAMI* 31(1), 39–58 (2009)