# An integrated approach for visual tracking of hands, faces and facial features

Maria Pateraki and Haris Baltzakis and Panos Trahanias

*Abstract*— **This paper presents an integrated approach for tracking hands, faces and specific facial features (eyes, nose, and mouth) in image sequences. For hand and face tracking, we employ a state-of-the-art blob tracker which is specifically trained to track skin-colored regions. The skin-color tracker is extended by incorporating an incremental probabilistic classifier, which is used to maintain and continuously update the belief about the class of each tracked blob, which can be left-hand, right hand or face as well as to associate hand blobs with their corresponding faces. Then, in order to detect and track specific facial features within each detected facial blob, a hybrid method consisting of an appearance-based detector and a feature based tracker is employed. The proposed approach is intended to provide input for the analysis of hand gestures and facial expressions that humans utilize while engaged in various conversational states with robots that operate autonomously in public places. It has been integrated into a system which runs in real time on a conventional personal computer which is located on the mobile robot itself. Experimental results confirm its effectiveness for the specific task at hand.**

## I. INTRODUCTION

In this paper, we propose an integrated approach to identify and track human hands, human faces and specific facial features in image sequences. The proposed approach is mainly intended to support natural interaction with autonomously navigating robots that guide visitors in museums and exhibition centers and, more specifically, to provide input for the analysis of hand gestures and facial expressions that humans utilize while engaged in various conversational states with a robot. The operational requirements of such an application challenge existing approaches in that the visual perception system should operate effectively under difficult conditions regarding occlusions, variable illumination, moving cameras, and varying background. The proposed approach combines and integrates a set of state-of-the-art techniques to solve three different but closely related problems: (a) identification and tracking of human hands and human faces which are detected as skin-colored blobs, (b) robust classification of the identified tracks to faces and hands, and, finally, (c) identification and tracking of specific facial features (eyes, nose and mouth) within each recognized facial blob.

For the first of the above defined problems (identification and tracking of human hands and faces) a variety of approaches have been reported in the literature [1]. Several of them rely on the detection of skin-colored areas [2]. The idea behind this family of approaches is to build appropriate color models of human skin and then classify image pixels based on how well they fit to these color models. On top of that, various segmentation techniques are used to cluster skin-colored pixels together into solid blobs that correspond to human hands and/or human faces.

In contrast to blob tracking approaches, model based ones [3] do not track objects on the image plane but, rather, in a hidden model-space. This is commonly facilitated by means of sequential Bayesian filters such as Kalman or particle filters. The state of each object is assumed to be an unobserved Markov process which evolves according to specific dynamics and which generates measurement predictions that can be evaluated by comparing them with the actual image measurements. Model based approaches are computationally more expensive and often require the adoption of additional constraints for the dynamics of the system and for the plausibility of each pose but they inherently provide richer information regarding the actual pose of the tracked human as well as the correspondence of specific body parts with the observed image.

In this work, we employ and extend a blob-tracking approach which is based on our previous work [4]. The blob-tracking approach, described in section III, has been extended by incorporating an incremental classifier which is used to maintain and continuously update a belief about whether a tracked hypothesis corresponds to a facial region, a left hand or a right hand (see section IV).

In the field of facial feature detection and tracking a number of approaches have already been presented in the existing literature [1]. Still, complexities arising from inter-personal variation (i.e. gender, race), intra-personal changes (i.e. pose, expression) and inconsistency of acquisition conditions render the task difficult and challenging. Related methods can be categorized on the basis of their inherent techniques. Facial feature localization based on attributes of geometrical shapes has been adopted in several works, e.g. [5], but the methods fail to show good performance in face images with large pose and expression variation. A variety of shape-based approaches tries to overcome this limitation by employing deformable templates [6], graph matching [7], active contours [8] or Hough transformation [9]. Color-based approaches were exploited by face detection systems to verify that a candidate blob is a face, by observing the darker appearance of facial elements in relation to their surroundings or local context [10]. Such approaches although may succeed to perform fast detection, they usually encounter difficulties in robustly detecting the skin color in the presence of different illuminations. Approaches based
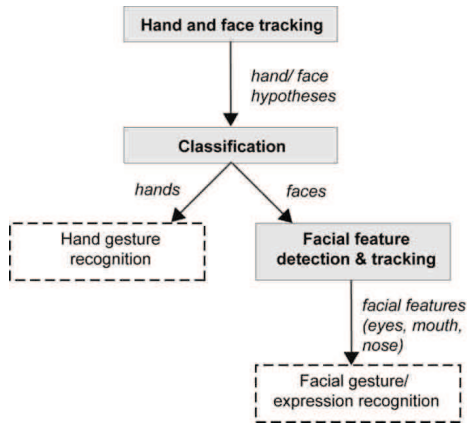
Fig. 1. Block diagram of the proposed system for hands and face tracking.

on machine learning techniques, like Principal Components Analysis [11], Neural Networks [12] and Adaboost Classifiers [13] require a large number of images for training and are computationally less efficient in the case of high resolution video sequences.

For detecting and tracking the facial features within the detected facial blobs, we propose an approach which combines the boosted cascade detector of Viola and Jones [14] with a feature based tracker and is described in section V. The resulting, combined detector and tracker extends our previous work on facial feature localization [15] in that specific anthropometric constraints are imposed after the initial detection step in order to enforce the elimination of false positives and provide reliable initial values for tracking.

The purpose of the above-described approach for hand, face and facial features tracking is to support recognition of hand gestures and facial expressions for rich interaction with an autonomous mobile robot. It has been integrated into a system which runs in real time on a conventional personal computer which is located on the mobile robot itself. Experimental results presented in section VI, confirm its effectiveness for this demanding task.

## II. METHODOLOGY

A block diagram of the components that comprise the proposed approach is depicted in Figure 1. The first block in Figure 1 is the hand and face tracker. This component is responsible for identifying and tracking hand and face blobs based on their color and on the information of whether they lay in the image foreground or not. The second step of the proposed system involves the classification of the resulting tracks into tracks that belong to facial blobs and tracks that belong to hands; left and right hands are also classified separately in this step.

Hand trajectories are forwarded to the hand-gesture recognition system (not described in this paper) while facial regions are further analyzed in order to detect and track specific facial features (eyes, nose and mouth) and to facilitate facial gestures and expression recognition at a later processing stage of the system (also not part of this paper).

In the following sections we describe each of the above mentioned components in detail.

## III. HAND AND FACE TRACKING

In this work, hand and face regions are detected as solid blobs of skin-colored, foreground pixels and they are tracked over time using the propagated pixel hypotheses algorithm [4]. This specific tracking algorithm allows the tracked regions to move in complex trajectories, change their shape, occlude each other in the field of view of the camera and vary in number over time.

Initially, the foreground area of the image is extracted by the use of a background subtraction algorithm [16]. Then, foreground pixels are characterized according to their probability to depict human skin and then grouped together into solid skin color blobs using hysteresis thresholding and connected components labeling. The location and the speed of each blob is modelled as a discrete time, linear dynamical system which is tracked using the Kalman filter equations. Information about the spatial distribution of the pixels of each tracked object (i.e. its shape) is passed on from frame to frame by propagating a set of pixel hypotheses, uniformly sampled from the original object's projection, to the target frame using the object's current dynamics, as estimated by the Kalman filter. The density of the propagated pixel hypotheses provides the metric which is used in order to associate observed skin-colored pixels with existing object tracks in a way that is aware of each object's shape and the uncertainty associated with its track.
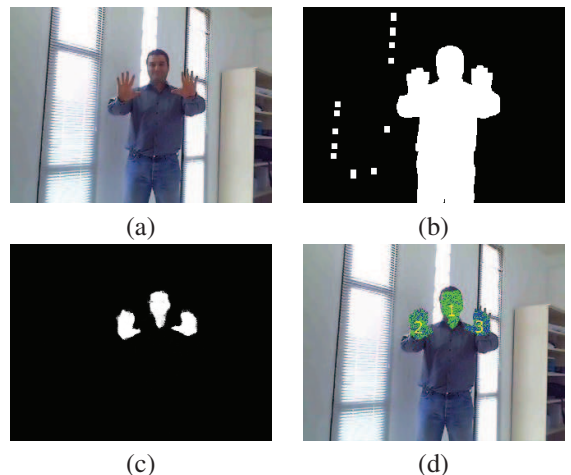


Fig. 2. The tracking approach. (a) Initial image, (b) background subtraction result, (c) pixel probabilities, (d) hand and face hypotheses.

Figures 2 and 3 demonstrate the operation of the employed hand and face tracker on a test sequence which involves a man performing hand gestures in an office environment. Figure 2(a) shows a single frame from this sequence. Figures 2(b) and 2(c) depict foreground pixels and skin-colored pixels, respectively. White pixels are pixels with high probability to be foreground/skin-colored pixels and black pixels are non-skin pixels. Finally, Fig. 2(d) depicts the hand and face hypotheses as tracked by the proposed tracker.

The output of the tracking algorithm in a number of frames from the same sequence is demonstrated in Fig. 3. As can be easily observed, this specific tracker succeeds in keeping track of all the three hypotheses, despite the occlusions and the blob merging events introduced at various fragments of the sequence.
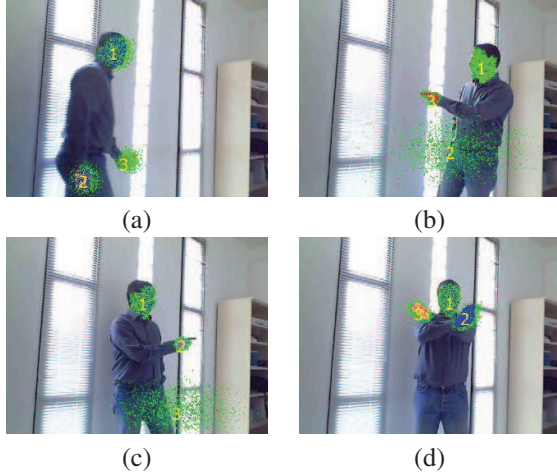


(a)  (b)

(c)  (d)

Fig. 3. Indicative tracking results for four segments of the office image sequence used in the previous example. In all cases the algorithm succeeds in correctly tracking the three skin-colored regions.

A more detailed explanation of the above-described hand and face tracker is given in [4].

## IV. CLASSIFYING BETWEEN HANDS AND FACES

To proceed with higher level tasks, like hand gestures and facial expressions recognition, one has to distinguish between tracks that belong to hands and tracks that belong to faces. Moreover, for hand tracks, one has to know which tracks belong to left hands and which tracks belong to right hands. Towards this goal, we have developed a technique that incrementally classifies a track into one of three classes: faces, left hands and right hands.

The input of the technique is a feature vector $O_t$ which is extracted at each time instant $t$ and is used to update the belief of the robot $B_t$ regarding the class $F$ of each track. The feature vector $O_t$ consists of the following components:

- The periphery-to-area ratio $r_t$ of the current track's blob. The ratio $r_t$ is normalized to the corresponding ratio of a circle and provides a measure of the complexity of the blob's contour. It is expected that hands will generally have more complex contours than faces, i.e. larger values for $r_t$.
- The vertical and the horizontal components $u_t$ and $v_t$ of the speed of a tracked skin-colored blob. The intuition behind this choice is that hands are generally expected to move faster than faces. Moreover, faces are not expected to have large vertical components in their motion.
- The orientation $\theta_t$ of the blob. It is expected that faces will tend to have orientations close to $\pi/2$.
- The location $l_t$ of the blob within the image. This location is relative to the location of each possible head

hypothesis and it is normalized according to the radius of this head, as it will be explained later in this section.

We define the belief $B_t$ of the robot at time instant $t$ to be the probability that the track belongs to class $f$, given all observations $O_i$ up to time instant $t$. That is:

$$B_t = P(F = f|O_1, \ldots, O_{t-1}, O_t) \tag{1}$$

$$= \frac{P(O_t|F = f, O_1, \ldots, O_{t-1})P(F = f|O_1, \ldots, O_{t-1})}{P(O_t|O_1, \ldots, O_{t-1})} \tag{2}$$

Since the denominator $P(O_n|O_1, \ldots, O_{t-1})$ is independent of $F$, we can substitute it with $1/\alpha$ and we obtain

$$B_t = \alpha P(O_n|F = f, O_1, \ldots, O_{t-1})P(F = f|O_1, \ldots, O_{t-1}) \tag{3}$$

$$= \alpha P(O_n|F = f, O_1, \ldots, O_{t-1})B_{t-1} \tag{4}$$

The above equation defines an incremental way to compute $B_t$, i.e. to classify the track by incrementally improving the belief $B_t$ based on the previous belief $B_{t-1}$ and the current observations.

Taking into account the Markov property and the independence assumptions indicated by Figure 4(a), we can further simplify the above equation:

$$B_t = \alpha P(O_t|F = f)B_{t-1} \tag{5}$$

In order to compute the term $P(O_t|F = f)$ in the right hand of Equation (5), we assume the naive Bayes classifier depicted in the graph of Figure 4 (b). According to this graph, we have at time instant $t$:

$$P(O_t|F) = \frac{P(F, O_t)}{P(F)} \tag{6}$$

$$= P(r_t|F)P(u_t|F)P(v_t|F)P(\theta_t|F)P(l_t|F) \tag{7}$$



(a)  (b)

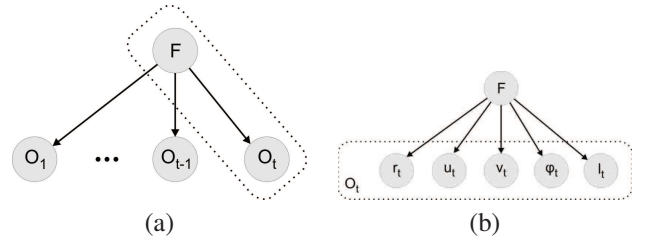Fig. 4. (a) Bayes graph encoding the independence assumptions of our approach, (b) The naive Bayes classifier used to compute the $P(O_t|F = f)$.

All the probabilities in the right side of Equation (7) can be estimated according to training data. Hence they are encoded and stored in appropriate look-up tables, thus permitting real-time computations.

The lookup tables for $P(r_t|F)$, $P(u_t|F)$, $P(v_t|F)$ and $P(\theta_t|F)$ are depicted in Figure 5. They are 1D lookup tables encoding the relevant quantity ($r$, $u$, $v$, or $\theta$) with the probability of appearance of this quantity in the training set. These lookup tables are identical for left hands and right hands but they are different in the case of faces. This is because, the relevant quantities are not expected to vary significantly between left and right hands but, as can be

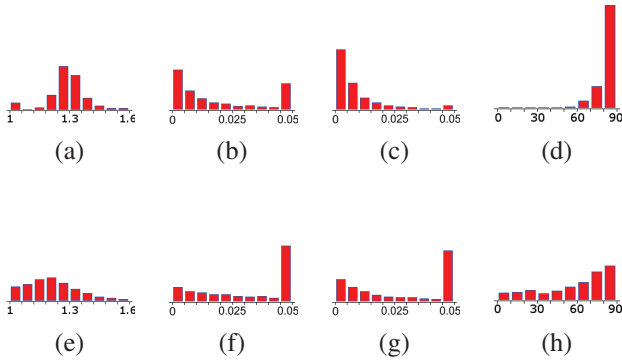easily observed in Figure 5, they differ significantly in the case of faces.



Fig. 5. 1D Look-up tables used for the computation of Equation (7). (a): $P(r_t|F=face)$, (b): $P(u_t|F=face)$, (c): $P(v_t|F=face)$, (d): $P(\theta_t|F=face)$, (e): $P(r_t|F=hand)$, (f): $P(u_t|F=hand)$, (g): $P(v_t|F=hand)$, (h): $P(\theta_t|F=hand)$.

$P(l_t|F)$, which is the probability of a blob being observed at location $l_t$ given its class $F$, is computed and stored differently for faces and differently for hands.

For faces, $P(l_t|face)$ is retrieved as the probability for a facial blob to be centered at this specific image location $l_t$. Obviously, the 2D lookup table for $P(l_t|face)$ depends on the actual application at hand and involves assumptions about the pose of the camera and the relative location of the human(s) with respect to the camera. In our case, which involves a human-robot interaction application, we assumed a camera placement such that the field of view of the camera includes the upper body part of one or more humans standing at a convenient distance between 0.5m and 2m in front of the robot. The actual lookup table that we compiled and used in our experiments is depicted in Figure 6(a).
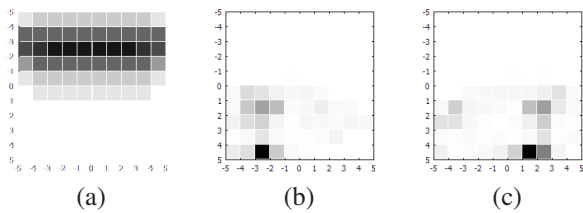


Fig. 6. 2D Look-up tables used for the computation of $P(l_t|F)$ in Equation (7). (a) for faces, (b) for left hands, (c) for right hands.

For hands, $P(l_t|left\,hand)$ and $P(l_t|right\,hand)$ are computed relatively to the location of the corresponding person's face. Since we don't know which is the corresponding person's face, we marginalize over all possible face hypotheses.

That is, for $P(l_t|right\,hand)$ we have:

$$P(l_t|right\,hand) = \sum_h P(l_t|right\,hand, h=face)P(h=face) \tag{8}$$

and similarly for the left hand:

$$P(l_t|left\,hand) = \sum_h P(l_t|left\,hand, h=face)P(h=face) \tag{9}$$

Figures 6(b) and 6(c) depict the resulting lookup tables for $P(l_t|right\,hand, h=face)$ and $P(l_t|left\,hand, h=face)$.

## V. DETECTION AND TRACKING OF FACIAL FEATURES

For tracking individual facial features within the detected facial blob, we utilize a hybrid approach by integrating an appearance-based detector and a feature-based tracker for the eyes, the nose and and mouth. The facial feature detector and tracker combines the advantages of appearance-based methods in detection (i.e. robustness in various lighting conditions), and feature-based methods in tracking (i.e. computational speed and high accuracy when initial estimation is close to the real solution) and permits robust identification of the facial features as well as real-time computations.

The overview of the implemented approach is illustrated in Figure 7 and is based on three steps: (a) initial detection of facial features using an appearance-based detector, (b) elimination of false positive detections via the application of anthropometric constraints, and, (c) real time tracking of the detected and filtered facial features using a feature-based method.
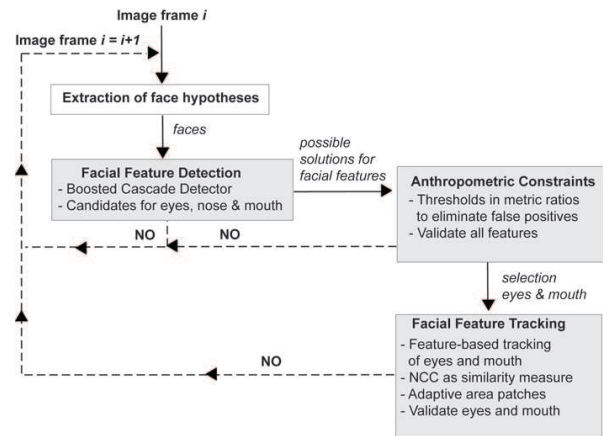


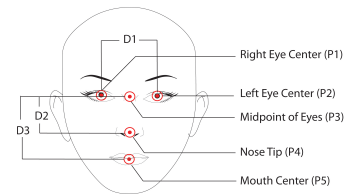Fig. 7. Diagram of the proposed approach for detection and tracking of facial features.



Fig. 8. Landmarks in the Anthropometric Face Model.

For the initial detection of facial features we use the Boosted Cascade Detector of Viola and Jones [14]. The detector has been designed for general object detection and
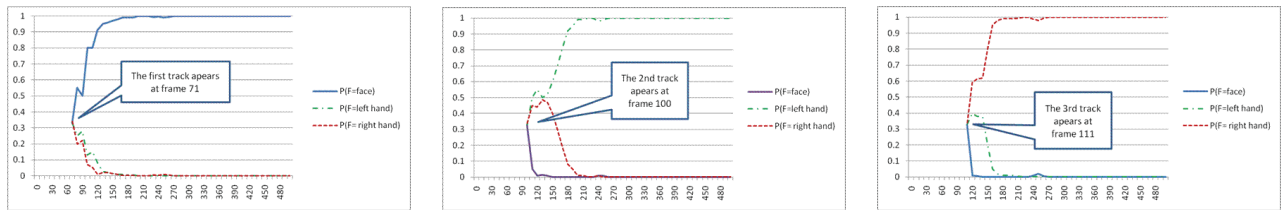
Fig. 9. The belief of each of the three tracks of the office sequence, as it evolves over the first 500 frames. The solid blue lines correspond to the probability of each blob being a face blob, the dot-dashed green lines correspond to left hands and the dashed red lines corresponds to right hands.

has gained widespread acceptance due to the availability of an implementation in an open source library [17]. In our case, for the detection of the features within each face blob, individual sets of Haar-like features for eyes, nose and mouth are utilized and the method is initialized with frontal-view faces.

An important factor which affects both the reliability of detection and the tracking accuracy of facial features is the size of the detected face blob. According to Tian [18], facial features become hard to detect when the face region is smaller than a threshold of approximately $70 \times 90$ pixels. Therefore, the procedure of facial feature detection and tracking is only activated when the face blob satisfies the above size requirements.

After all features have been detected, specific anthropometric constraints are applied in order to cast out false positives. Motivated by the work of Sohail and Bhattacharya [19], we have collected a large set of measurements from images depicting faces in frontal view. The collected measurements were used to built an anthropometric model of the human face and to define the necessary thresholds and validation gates used to filter out false positive detections. The selected validation criteria involve the location and the size of the eyes, the nose and the mouth. Landmarks on other regions such as the eyebrows, used, for example, in [19], were not selected because they often proved to be occluded by hair, eyeglasses or, in some cases, they were entirely non-existent.

More specifically, we define the following criteria:

- All four selected features (eyes, nose, mouth) should be detected.
- The normalized sizes of the two eyes and mouth should be within certain bounds.
- The normalized distance between the midpoint of the eye centers and nose tip should be approximately 0.6. That is $D_2/D_1 \simeq 0.6$, where $D2$ is the distance between points $P_3$ and $P_4$ (see Figure 8).
- The normalized distance between the midpoint of eye centers and mouth center should be approximately $\simeq$ 1.2. That is $D_3/D_1 \simeq 1.2$, where $D_3$ is the distance between points ($P_3$ and $P_5$).

If the above criteria are not met by the system and a new re-initialization is attempted (by repeating the facial feature detection step) in the next frame, otherwise the tracking procedure is invoked. It is to be noted that the nose region is not tracked because it's actual location is not considered important for our target application, which is expression

recognition and visual speech detection.

Our tracking approach is based on template matching which is implemented using the normalized cross-correlation (NCC) measure as matching score/quality measure. The selection of NCC as quality measure is justified as only small deviations in the relative positions of the feature areas with respect to the position of the face blob in the image are expected. The detected eye and mouth regions from each face are used as templates in the matching process, updated in every consecutive frame and the matching score is used to block results of low reliability.

## VI. EXPERIMENTAL RESULTS WITH REAL WORLD DATA

Figure 10 presents hand and face classification results for various frames of the office sequence of Fig. 3. Blobs classified as faces are marked with an "F", left hands are marked with an "L", and right hands are marked with an "R". The proposed approach has been successful in classifying the three observed tracks and it also managed to maintain its belief over the whole sequence.
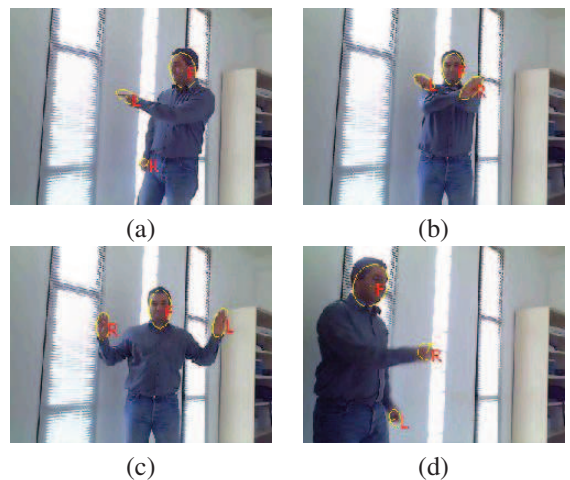


(a)  (b)

(c)  (d)

Fig. 10. Four frames of a sequence depicting a person performing various hand gestures in an office environment.

Figure 9 depicts the belief of each of the three tracks of the office sequence, as it evolves over the first 500 frames of this sequence. As can be easily observed, the belief of each track is initially uncertain but very soon it stabilizes to the correct class. The belief stays stable to the correct classes thoughout the whole sequence consisting of a total of 2600

frames (for clarity of presentation, only the first 500 frames are shown in these graphs).
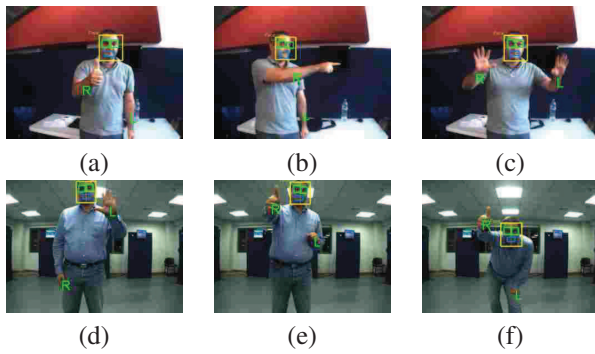


Fig. 11. Frames of two different sequences captured in an exhibition center that show results from hand, face and facial feature tracking.
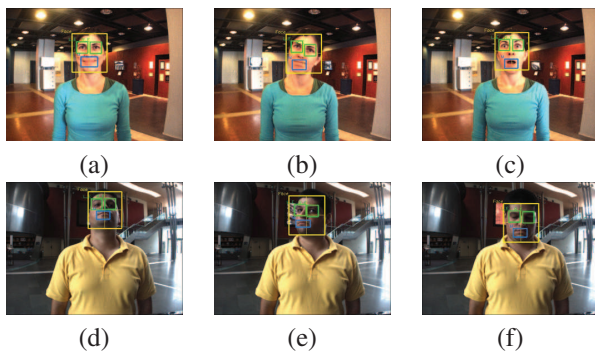


Fig. 12. Frames of different sequences captured in an exhibition center that show results from facial feature tracking of a person.

Figure 11 depicts some frames from two additional sequences captured by the robot's camera in two different application environments within an exhibition center. In all our experiments the algorithm successfully tracked the skin-colored blobs and very fast converged to the correct class for each track (i.e. left hands, right-hands and faces), following convergence curves which were very similar to the ones depicted in Fig. 9. Eyes, nose and mouth regions were also correctly localized and tracked, even in cases of usual off-plane head rotations and different facial expressions.

Figure 12 depicts some additional facial feature tracking results from two additional, close-up, sequences captured at the same exhibition center. As with the previous figure, facial features were correctly localized and tracked. More experimental results from different application environments (office, exhibition center) can be found in http://www.ics.forth.gr/~pateraki/handfacetracking.html.

## VII. CONCLUSIONS

In this paper we have presented an integrated approach for tracking of hands, faces and facial features in image sequences, intended to support natural interaction with autonomously navigating robots in public spaces and, more specifically, to provide input for the analysis of hand gestures and facial expressions that humans utilize while engaged in various conversational states with the robot.

For hand and face tracking, we employ a blob tracker which is specifically trained to track skin-colored regions. The skin-color tracker is extended by incorporating an incremental probabilistic classifier which is used to maintain and continuously update the belief about the class of each tracked blob which can be a left-hand, a right hand or a face. Facial feature detection and tracking is performed via the employment of state-of-the-art appearance-based detection coupled with feature-based tracking, using a set of anthropometric constraints.

Experimental results have confirmed the effectiveness of the proposed approach proving that the individual advantages of all involved components are maintained, leading to implementations that combine accuracy, efficiency and robustness. Future work includes tracking hands, faces and facial features of multiple people in the scene.

## REFERENCES

[1] M.H. Yang, D. Kriegman, and D. Ahuja. Detecting faces in images: A survey. *IEEE Trans. PAMI*, 24(1):34–58, 2002.
[2] K. Nickel, E. Seemann, and R. Stiefelhagen. 3d-tracking of head and hands for pointing gesture recognition in a human-robot interaction scenario. In *Proc. IEEE Intl. Conf. on Automatic Face and Gesture Recognition*, pages 565–570, Seoul, Korea, May 2004.
[3] B. Stenger, A. Thayananthan, P. H. S. Torr, and R. Cipolla. Model-based hand tracking using a hierarchical bayesian filter. *IEEE Trans. PAMI*, 28(9):1372–1384,, September 2006.
[4] H. Baltzakis and A. Argyros. Propagation of pixel hypotheses for multiple objects tracking. In *Proc. International Symposium on Visual Computating (ISVC)*, Las Vegas, Nevada, USA, November 2009.
[5] C. Perez, V. Lazcano, P. Estvez, and C. Held. Real-time template based face and iris detection on rotated faces. *International Journal of Optomechatronics*, 3(1):54–67, 2009.
[6] A.L. Yuille, P.W. Hallinan, and D.S. Cohen. Feature extraction from faces using deformable templates. *IJCV*, 8(2):99111, 1992.
[7] R. Herpers and G. Sommer. An attentive processing strategy for the analysis of facial features. *Face recognition: From Theory to Applications*, pages 457–468, 1998.
[8] M. Pardas and M. Losada. Facial parameter extraction system based on active contours. In *ICIP01*, pages I: 1058–1061, 2001.
[9] T. Kawaguchi, M. Rizon, and D. Hidaka. Detection of eyes from human faces by hough transform and separability filter. *Electronics and Communications in Japan*, 88(5):2939, 2005.
[10] X. Zhang, Y. Xu, and L. Du. Locating facial features with color information. volume 2, pages 889 –892 vol.2, 1998.
[11] H.-C. Kim, D. Kim, and S.-Y. Bang. A pca mixture model with an efficient model selection method. In *Proc. Intl. Joint Conf. on Neural Networks (IJCNN '01).*, volume 1, pages 430 –435, 2001.
[12] S. Phimoltares, C. Lursinsap, and K. Chamnongthai. Locating essential facial features using neural visual model. In *Proc. Intl. Conf. on Machine Learning and Cybernetics.*, volume 4, pages 1914 – 1919, nov. 2002.
[13] P. Wilson and J. Fernandez. Facial feature detection using haar classifiers. *J. Comput. Small Coll.*, 21(4):127–133, 2006.
[14] P. Viola and M. Jones. Robust real-time face detection. *IJCV*, 57(2):137–154, 2004.
[15] M. Pateraki, H. Baltzakis, P. Kondaxakis, and P. Trahanias. Tracking of facial features to support human-robot interaction. In *Proc. IEEE International Conference on Robotics and Automation (ICRA '09)*, pages 3755 –3760, may 2009.
[16] Chris Stauffer and W. Eric L. Grimson. Adaptive background mixture models for real-time tracking. In *Proc. IEEE Computer Vision and Pattern Recognition (CVPR)*, pages 2246–2252, June 1999.
[17] Intel. Intel open source computer vision library, v2.0. (Dec 2010).
[18] Y. Tian. Evaluation of face resolution for expression analysis. In *Proc. IEEE Conf. CVPR'04*, pages 82–89. IEEE Computer Society, 2004.
[19] A. S. M. Sohail and P. Bhattacharya. Detection of facial feature points using anthropometric face model. In E. Damiani et al., editor, *Signal Processing for Image Enhancement and Multimedia Processing*, volume 31 of *Multimedia Systems and Applications*, pages 189–200. Springer US, 2008.