



# MUCATAR PROJECT - DELIVERABLE D3 : TOWARDS JOINT TRACKING AND ACTIVITY RECOGNITION

Jean-Marc Odobez <sup>1</sup>      Sileye Ba <sup>1</sup>  
Daniel Gatica-Perez <sup>1</sup>      Kevin Smith <sup>1</sup>  
IDIAP-RR

JANUARY, 2004

Dalle Molle Institute  
for Perceptual Artificial  
Intelligence • P.O.Box 592 •  
Martigny • Valais • Switzerland

phone +41 - 27 - 721 77 11

fax +41 - 27 - 721 77 12

e-mail secre-

tariat@idiap.ch

internet

<http://www.idiap.ch>

---

<sup>1</sup> IDIAP, Martigny, Switzerland

# MUCATAR PROJECT - DELIVERABLE D3 : TOWARDS JOINT TRACKING AND ACTIVITY RECOGNITION

Jean-Marc Odobez

Sileye Ba  
Kevin Smith

Daniel Gatica-Perez

JANUARY, 2004

**Abstract.** This document presents the progress achieved in the MUCATAR (Multiple Camera Tracking and Activity Recognition) IM2 White Paper Project towards human tracking and activity recognition. A generic methodology has been developed to perform joint tracking and recognition using mixed-state models in a particle filter framework. The approach combines in a joint pdf continuous-valued motion parameters (e.g. location, size, aspect ratio of people in the image) with discrete labels (e.g. index of a set of exemplars of appearance).

This reports presents the application of this methodology to the specific problem of joint head tracking and pose estimation. In this framework, the discrete variable characterizes a specific head pose and appearance. The estimated discrete variable can then be exploited to perform simple action recognition such as head turns.

# 1 Introduction

This document presents the progress achieved in the MUCATAR (Multiple Camera Tracking and Activity Recognition) IM2 White Paper Project towards human tracking and activity recognition.

A generic methodology has been developed to perform joint tracking and recognition using mixed-state models in a particle filter framework. The approach combines in a joint pdf continuous-valued motion parameters (e.g. location, size, aspect ratio of people in the image) with discrete labels (e.g. index of a set of exemplars of appearance [9]). In the Bayesian formulation for tracking, the particle filter algorithm should now account for the different nature of the latent variables. The methodology has been applied to the two following problems :

- multi-camera speaker tracking in the meeting room. In this application, the discrete latent variable represents the camera index. Data likelihood terms need to extract speaker information from the camera identified by the discrete variable. The definition of the dynamics needs the specification of camera transition probabilities characterized by likely/unlikely camera switching image regions.
- joint head tracking and head pose estimation. The discrete latent variable describes an appearance exemplar, characterized by a given pose and a given appearance. In this case, the discrete variable indicates the exemplar to which the extracted information will be compared with. More details are given in the next Section.

The estimation of the discrete variable allows for simple action recognition. For instance, in the head tracking case, the sequence of exemplars can easily be used to recognize simple head gestures such as head turns.

The rest of this document is organized as follows. Section 2 describes the methodology applied to the joint head tracking and pose estimation problem. Section 3 describes an example of action recognition (dynamic left and right head turns), while Section 4 concludes the document with the presentation of ongoing work and future work in activity recognition.

## 2 Joint Head Tracking and Pose Estimation with Particle Filter

### 2.1 Introduction

Head detection and tracking are essential components in video applications related to human behaviour understanding. It is commonly used as a first step before applying algorithms for other higher level tasks, such as face and facial expression recognition or gaze direction estimation. At the same time, the estimation of the head pose could be useful for behaviour understanding and to improve the higher level tasks.

Many methods have been proposed to estimate head pose [1], [3], [4], [7], [10],[11]. To our knowledge, the previous work consider tracking and head pose estimation as two

sequential but independent problems. The principle of these methods is to first track the head to extract its location, and then to estimate head orientation by exploiting this location. As a consequence, the head pose estimation process is very dependent on the accuracy of the tracking since, as reported in [1], head pose is very sensitive to the localization of the extracted head box. At the same time, the knowledge of the head pose could improve the head modeling and thus the accuracy of the tracking. This paper addresses these issues by coupling the tracking and head pose estimation processes in a probabilistic setting. For this purpose, a mixed-state particle filter framework is used [9], where a head spatial configuration (e.g. position and scale) and its pose are represented in a joint state-space model. The joint posterior distribution of the state given the sequence of images is estimated at each instant and propagated to the next time instant using the state dynamic. The pose at a given instant is then obtained by marginalizing over the spatial configuration part of the state. As a result, in the approach we propose, the spatial configurations leading to a better pose modeling will have a greater impact on the pose result, leading to a more accurate estimation of the pose than with the tracking *then* pose estimation approach. This is supported by experiments performed on several real sequences.

This report section is organized as follows. Section 2.2 describes our head pose modeling. Section 2.3 shows the embedding of these pose models in a mixed-state particle filter framework. Section 2.4 reports results of pose estimation on still images and tracking results on real sequences. Section 2.5 gives the conclusions of this work.

## 2.2 Head Pose Modeling and Estimation

### 2.2.1 Head Pose Modeling and Learning

The head poses are defined by a pan angle denoted  $\theta$  and ranging from -90 to 90 degrees<sup>1</sup>. Allowed values are discretized with a 22.5 degrees step. Training data patches are extracted from head images by locating a tight bounding box around the head. These patch images are resized to the same  $64 \times 64$  resolution and preprocessed by histogram equalization to reduce the effect of lighting conditions. Four filters, one Gaussian and three rotation invariant Gabor wavelets, are then applied on these patches (Fig. 1). A simple Gabor wavelet is defined by:

$$\psi_{\omega_0, \sigma, \alpha}(x, y) = \exp\left(-\frac{1}{2\sigma^2}(x'^2 + y'^2)\right) \cos(2\pi\omega_0 x')$$

$$x' = x \cos \alpha - y \sin \alpha \text{ and } y' = x \sin \alpha + y \cos \alpha$$

where  $\omega_0$  denotes the angular frequency,  $\sigma$  the scale parameter and  $\alpha$  the orientation of the wavelet. A rotation invariant wavelet is obtained by integrating a simple wavelet over the orientation  $\alpha$ . The rotation invariant Gabor wavelet we used are defined by the scales  $\sigma = 1, 2, 4$  and and angular frequency  $w_0 = \frac{1}{2}, \frac{1}{4}, \frac{1}{8}$ . The resulting images are sampled at 191 points of a grid  $G$  regularly located inside a reference disk  $C$  of center  $(32.5, 32.5)$  and of radius 31.5 (Fig. 1a).

---

<sup>1</sup>An additional tilt angle is also considered but is left aside for brevity in the presentation.



Figure 1: a) Reference grid on a frontal head pose b) The four image features computed from the frontal head pose.

For each filter  $\Psi_i$ , the features computed from an image  $\{f_j^i, j \in G\}$  are normalized to give  $\tilde{f}^i = \{f_j^i = \frac{f_j^i - m_i}{s_i}, j \in G\}$ , where  $m_i$  and  $s_i^2$  represent the mean and variance of the  $i$  features, and are given by :

$$m_i = \frac{1}{|G|} \sum_{j \in G} f_j^i \quad \text{and} \quad s_i^2 = \frac{1}{|G|} \sum_{j \in G} f_j^{i2} - m_i^2 \quad (1)$$

This normalization is made to prevent the features of a filter to dominate the other because their values are higher. These features are then concatenated in a single feature vector  $z = \{\tilde{f}^i, i = 1, 2, 3, 4\}$ .

To learn the model of a head pose we use the CMU PIE database [8], which contain 68 persons at the needed head poses. For each head pose  $\theta$ , the feature vectors are clustered in  $K$  clusters using a Kmeans algorithms. The  $K$  centers of cluster,  $e_k^\theta, k = 1, \dots, K$  are taken to be the models of the head pose. For each head pose the standard deviation of the features  $\sigma_k^\theta$  and the normalized number of element of each cluster  $\pi_k^\theta$  are kept. The Kmeans procedure was preferred to others modeling methods like Gaussian mixture model because our interest is in modeling representative exemplars of head pose and not directly the probability distribution of the features.

### 2.2.2 Head Pose Estimation

The head pose of an input image characterized by its feature  $z$  is estimated using the maximum a posteriori principle :

$$\theta^* = \arg \max_{\theta} p(\theta|z) = \arg \max_{\theta} \frac{p(z|\theta)p(\theta)}{p(z)}. \quad (2)$$

Assuming for static images that  $p(\theta)$  is uniformly distributed, the MAP estimation resume to  $\theta^* = \arg \max_{\theta} p(z|\theta)$ . We assume that for each head pose  $\theta$  the components of the feature vector are independent and can be modeled by a Gaussian mixture having as center the exemplars  $e_k^\theta, k = 1, \dots, K$ , as diagonal covariance matrix  $\Sigma^\theta = \text{diag}(\sigma_k^{\theta2})$  and as probability mixtures  $\pi_k^\theta$ . The probability of data given a head pose is modeled by:

$$p(z|\theta) = \sum_{k=1}^K \pi_k^\theta p(z|k) \quad \text{with:} \quad (3)$$

$$p(z|k) = \prod_i \frac{1}{\sigma_{k,i}^\theta} \exp -\frac{1}{2} \left( \frac{z_i - e_{k,i}^\theta}{\sigma_{k,i}^\theta} \right)^2. \quad (4)$$

As components of a feature vector can be outliers, we will also use the saturated Gaussian likelihood:

$$p_T(z|k) = \prod_i \frac{1}{\sigma_{k,i}^\theta} \max \left\{ \exp -\frac{1}{2} \left( \frac{z_i - e_{k,i}^\theta}{\sigma_{k,i}^\theta} \right)^2, T \right\}. \quad (5)$$

where  $T = \exp^{-3}$  is a lower threshold. This term is useful to avoid local differences between an exemplar and the input image (e.g. in the hair cut) to conduct to a very low likelihood even when the majority of the remaining component features are in good agreement.

## 2.3 Joint Tracking and Head Pose Estimation

Head tracking and pose estimation are performed in a probabilistic framework.

### 2.3.1 Mixed-State Particle Filter.

Particle filtering (PF) implements a recursive Bayesian filter by Monte-Carlo simulations. Let  $X_{0:t} = \{X_l, l = 0, \dots, 1\}$  (resp.  $z_{1:t} = \{z_l, l = 1, \dots, t\}$ ) represents the sequence of states (resp. of observations) up to time  $t$ . Furthermore, let  $\{X_{0:t}^i, w_t^i\}_{i=1}^{N_s}$  denote a set of weighted samples that characterizes the posterior probability density function (pdf)  $p(X_{0:t}|z_{0:t})$ , where  $\{X_{0:t}^i, i = 1, \dots, N_s\}$  is a set of support points with associated weights  $w_t^i$ . The samples and weights can be chosen using the Sequential Importance Sampling (SIS) principle. Assuming that the observations  $\{z_t\}$  are independent given the sequence of states, the state sequence  $X_{0:t}$  follows a first-order Markov chain model, and that the prior distribution  $p(X_{0:t})$  is employed as proposal, we obtain the following recursive update equation [2] for the weight:

$$w_t^i \propto w_{t-1}^i p(z_t|X_t^i) \quad (6)$$

To avoid sampling degeneracy an additional resampling step is necessary [2]. The standard PF is given by :

1. **Initialisation** :  $\forall i \in 1:N_s$ , sample  $X_0^i \sim p(X_0)$ ; set  $t = 1$
2. **IS step**:  $\forall i$  sample  $\tilde{X}_t^i \sim p(X_t^i|X_{t-1}^i)$ ; evaluate  $\tilde{w}_t^i$  using (6).
3. **Selection**: Resample  $N_s$  particles  $\{X_t^i, w_t^i = \frac{1}{N_s}\}$  from the sample set  $\{\tilde{X}_t^i, \tilde{w}_t^i\}$ ; set  $t = t + 1$ ; go to step 2.

In the mixed state particle filter approach of [9], the state  $X = (k, x)$  is the conjunction of a discrete variable  $k$  labeling a discrete set of objects models  $e_k$ , called exemplars and a continuous variable  $x$  specifying the spatial configuration of the object (e.g. position, the size, image rotation). In order to implement the filter, three elements have to be specified: a state model, a dynamical model and an observation model.

### 2.3.2 State space

The state  $X$  is a mixed variable  $X = (k, x)$ . The discrete variable  $k = (\theta, l)$  labels an element of the set of head pose models  $\{e_l^\theta, \theta, l = 1, \dots, K\}$  built in the previous subsection. The continuous variable  $x = (t_x, t_y, s_x, s_y)$  is a vector parameterizing the transform  $\mathcal{T}_x$  defined by:

$$\mathcal{T}_x u = \begin{pmatrix} s_x & 0 \\ 0 & s_y \end{pmatrix} u + \begin{pmatrix} t_x \\ t_y \end{pmatrix}. \quad (7)$$

which characterizes the object configuration, where  $(t_x, t_y)$  specifies the translation of the object in the image plane, and  $(s_x, s_y)$  the scale of the width and the height of the object according to a reference size.

### 2.3.3 Dynamics

The process density on the state sequence  $X_t = (k_t, x_t)$  is modeled as a second order autoregressive process  $P(X_t|X_{t-1}, X_{t-2})$ . We assume that the two components of the states,  $k_t$  and  $x_t$ , are independent. Also it is assumed that a head pose at a given time  $t$ ,  $k_t$ , depends only on the head pose at the previous time  $k_{t-1}$ . Then the equation of the process density is:

$$P(X_t|X_{t-1}, X_{t-2}) = p(k_t|k_{t-1})p(x_t|x_{t-1}, x_{t-2}) \quad (8)$$

The dynamic of the continuous variable  $x$  is modeled as a classical second order auto regressive dynamical mode. The dynamic of the discrete variable  $k$ , defined by the transition process  $p(k_t|k_{t-1}) = p(\theta_t, l_t|\theta_{t-1}, l_{t-1})$ :

$$p(\theta_t, l_t|\theta_{t-1}, l_{t-1}) = p(l_t|\theta_t, l_{t-1}, \theta_{t-1})p(\theta_t|\theta_{t-1}). \quad (9)$$

$p(\theta_t|\theta_{t-1})$  is based on the distance between the two head poses.  $p(l_t|\theta_t, l_{t-1}, \theta_{t-1})$  is a probability table learned using the training set of faces. More precisely, for different head poses, the exemplars are more related when the same persons were used to build them. When  $\theta \neq \theta'$   $p(l|\theta, l', \theta')$  is taken proportional to the number of persons who belong to the class of  $e_l^\theta$  and who are also in the class of  $e_{l'}^{\theta'}$ . When  $\theta = \theta'$ ,  $p(l|\theta, l', \theta')$  is large for  $l = l'$  and small otherwise.

### 2.3.4 Observation model

Finally, let us define the object likelihood. For each state  $X = (k, x)$  the observations are obtained by first extracting an image patch from the image according to  $\mathcal{C}(x) = \{\mathcal{T}_x u, u \in \mathcal{C}\}$ , and then filtering this image patch at the points specified by the grid  $G$  with the four filters defined in the previous subsection, and concatenating the filtered values in a feature vector  $z(x)$ . The likelihood  $p(z|X) = p(z|k, x)$  is finally modeled by  $p(z|k, x) = p_T(z(x)|k)$ ,  $p_T$  referring to Equation 5.

The head pose is then estimated a each time as the mode of the head pose distribution after marginalization over the spatial configuration :

$$\theta_t^* = \arg \max_{\theta} \sum_{i/\theta_i^i=\theta} w_t^i \quad (10)$$

NEP	State of The Art [1]	Gaussian	Sat. Gaussian
1	90%	90%	94%
2	Not Relevant	87.5%	94.8%

Table 1: Recognition rate table for a given number of exemplar per head pose (NEP)

	90	67.5	45	22.5	0	-22.5	-45	-67.5	-90
90	100	0	0	0	0	0	0	0	0
67.5	14.7	85.3	0	0	0	0	0	0	0
45	0	0	100	0	0	0	0	0	0
22.5	0	0	5.9	94.1	0	0	0	0	0
0	0	0	0	0	100	0	0	0	0
-22.5	0	0	0	0	5.9	94.1	0	0	0
-45	0	0	0	0	0	5.9	94.1	0	0
-67.5	0	0	0	0	0	0	3.1	91	5.9
-90	0	0	0	0	0	0	0	0	100

Table 2: Confusion matrix for NEP=2 and saturated Gaussian likelihood

## 2.4 Results

### 2.4.1 Head Pose Estimation Results

To test the efficiency of the pose modeling we used the 68 persons of PIE database and their head pose. For the first experiments we use the same setup than [1]. The 34 first persons were selected and their head poses used to train the head pose models. The half remaining were used to test the models. Table 1 shows the recognition rates when the number of exemplars per head pose are 1 and 2.

This table shows that smoothing the likelihood is indeed very useful, helping in reducing the effect of outlier feature components. Besides, Table 2 displays the confusion matrix of the recognition.

It shows that estimation errors occur in general between close head poses. These errors are still acceptable, there is not a total mismatch of different profile views.

To further study the effect of the number of exemplars, we included in the database 72 persons of the FERET database [6], leading to a total of 140 persons. Then, 70 persons were randomly selected and their head pose used to train the models, and the half remaining used to test the models. We ran this set up 100 times and computed the average and standard deviation of recognition rates for NEP=1,2,3,4. Table 3 gives the best results that were achieved with NEP=3 for the Gaussian likelihood and NEP=4 for the Saturated Gaussian likelihood. These results show that more exemplars improve the recognition and that the saturated Gaussian likelihood is still doing better than the Gaussian likelihood (this indeed true for all NEP).



	Gaussian	Sat. Gaussian
NEP	3	4
Av. of R.R.	67.2%	70.5%
St.dev. of RR	2.4	2.3

Table 3: Best recognition rates for PIE+FERET

## 2.4.2 Tracking Results

The tracking algorithm described previously was tested on several video sequences. None of the tracked persons belonged to the training database. The positions of the camera for the video sequences used to test the algorithm were different than those used in the training database. Also the illumination condition were very different in the training database and in our test video sequences. Despite this mismatch between training and test sets, the tracking was correctly done and the estimation of the head pose was visually very satisfying. Fig. 2 shows tracking results of a typical sequence. Videos of more tracking results are available on the web.

We conducted experiments to compare our method to the traditional sequential head tracking *then* pose estimation approach. We used a color-based state-of-the-art particle filter tracker described in [5], which provides at each time a patch image corresponding to the head. The patch image is processed as described in subsection 2.2 to extract the features, then compared to the exemplars using Equation 5 for pose evaluation. For the sequence of Figure 2 we generated head orientation ground truth by manually extracting a tight bounding box around the head and applying the pose estimation method. At each time  $t$ , the surface of the ground truth box is denoted  $GS(t)$ . We ran the two trackers that output at each time a box containing the head with surface  $TS(t)$  and an estimated head pose. If at each time  $JS(t)$  is the joint surface between the ground truth and the tracker, we choose to measure the tracking error by  $e(t) = \frac{1}{2} \left( \frac{GS(t)-JS(t)}{GS(t)} + \frac{TS(t)-JS(t)}{TS(t)} \right)$ . This error is 0 when tracking is perfect, and 1 when it totally fails. Figure 3 shows tracking errors for the two methods and the estimated head orientations. The results in Fig. 3, left, shows that our method leads to smaller tracking errors in average. The color tracker, on the other side, can be confused by similar color in the background. This results in an over-estimate of the head patch size (cf. the two error peaks at the end of the sequence), which in turn results in pose estimation failures (Fig. 3, right), as illustrated by Fig. 4.

## 2.5 Conclusion on joint head tracking and pose estimation

We described in this section a joint head tracking *and* pose estimation algorithm. The novelty of the approach lie in the coupling of the tracking and head pose estimation processes. This coupling is handled in a probabilistic framework within a mixed state particle filter framework. By implicitly allowing to test multiple head configurations, it reduces the sensitivity of the pose estimation process on the tracking accuracy, a drawback of methods that perform head tracking *then* pose estimation in a sequential manner, and results in more stable and

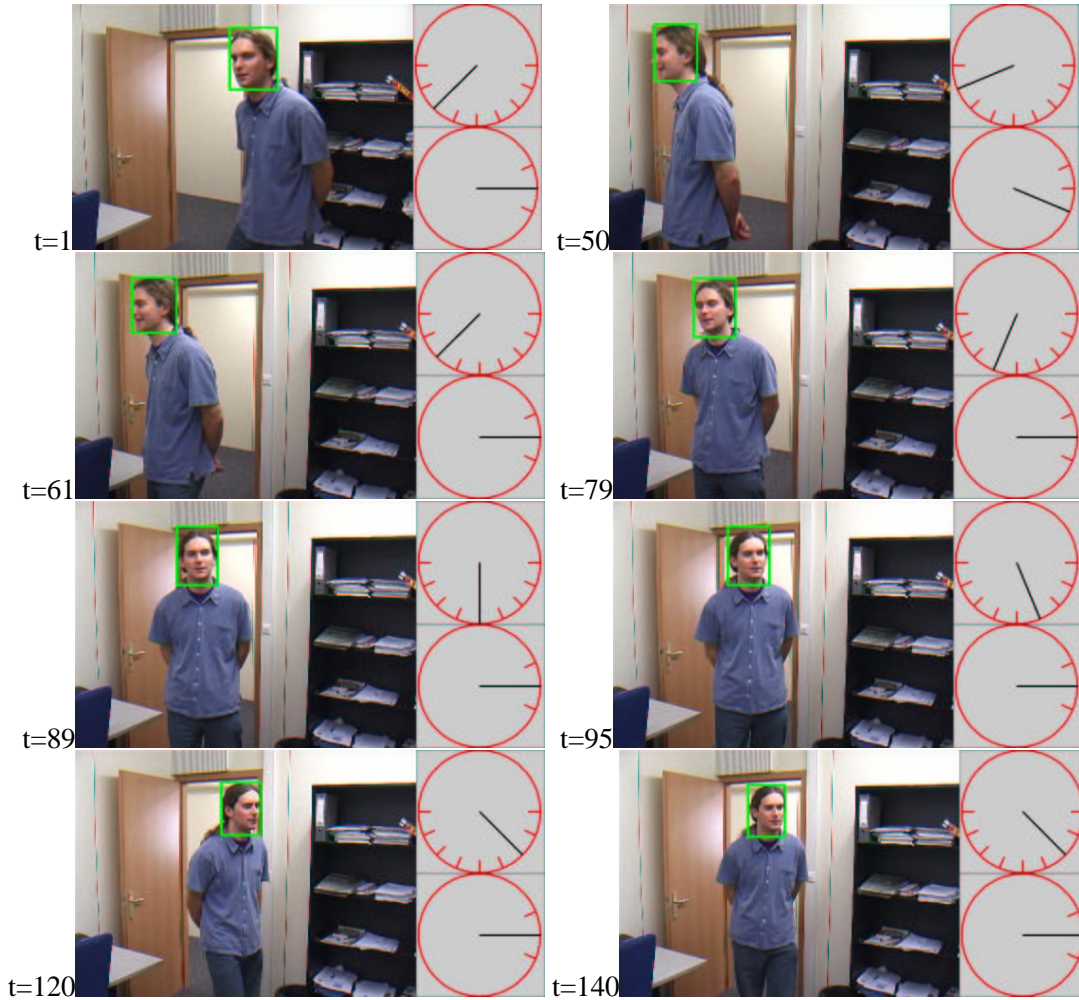


Figure 2: Tracking and head pose estimation results. First clock: pan angle; second clock: tilt angle.

accurate pose estimates.

### 3 Dynamic head turn detection

The previous Section described a methodology to extract the pose information from a head sequence. In this Section, we show on an example how it can be used for dynamic head turn detection. By dynamic head turns, we mean variations of around 15 degrees of head pose within less than a second.

The principle of the approach is the following. The example is taken from a 1 minute 20 sec video of meeting data. Some of the image are shown in Fig. 5, where we are interested in the head turns of the person A on the image right. Starting from the head pose data (top left of Fig. 6) produced by the tracker, a median filter is applied to eliminate the main fluctuations

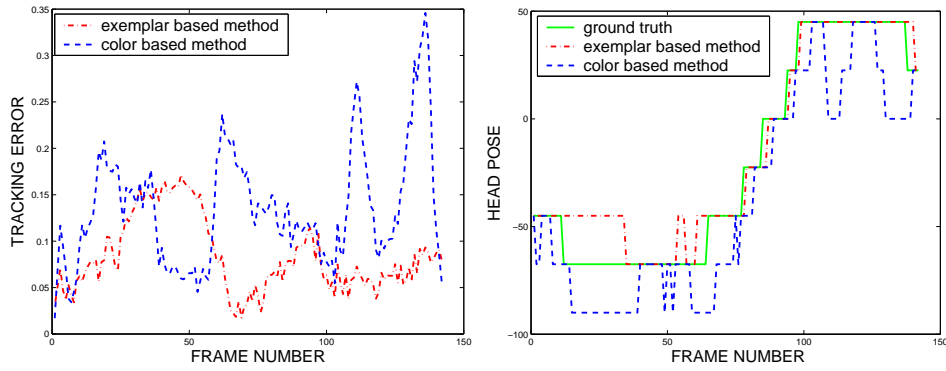


Figure 3: Left: spatial configuration errors. Right: Pan head orientation estimation.

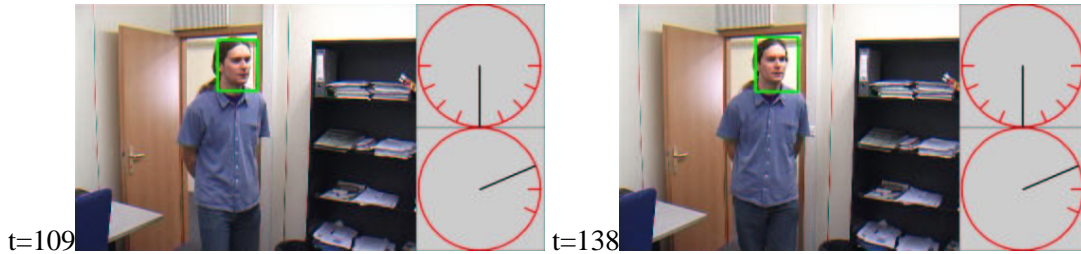


Figure 4: Sample of head pose estimation failure due to bad head location.

of the data curve. These are essentially due to the discrete nature (index of head pose are separated by 22.5 degrees) of the head pose representation in the filter, and the small number of particles that we use and which is not enough to have a very precise pose estimation. A gaussian filter is then applied, leading to a smooth curve (top right of Fig. 6). A derivative operator is then applied and the output is processed by a peak detector, that validate local head turns maxima with a speed above a 20 degrees/s threshold.

The results are shown in Fig. 6, bottom right. Some of the detected head turns of the person A are illustrated in the images of Fig. 5. The first one is a small head turn due to a shift in the focus of attention, from the speaker B (not seen) in front of A, to the neighbor of B. The two last rows illustrates a rapid head pose variation captured by our tracking system (-45 and then +45 degrees in less than 0.5 second).

## 4 Conclusion

In this report, we have shown a methodological framework to perform some joint tracking and action recognition framework. This has been applied to the head tracking and pose estimation problem, for which advantages of this framework can be achieved. The output of this system can be further processed to apply more advanced recognition steps. An example is given for illustration purposes.

In the future, we will investigate the recognition of a larger set of head events. We have identified the following activities :



Figure 5: Detecting head turns. Example of some detected head turns (person on the right). Images at time 240, 250; time 410, 420; time 1600, 1605, 1610, 1620.

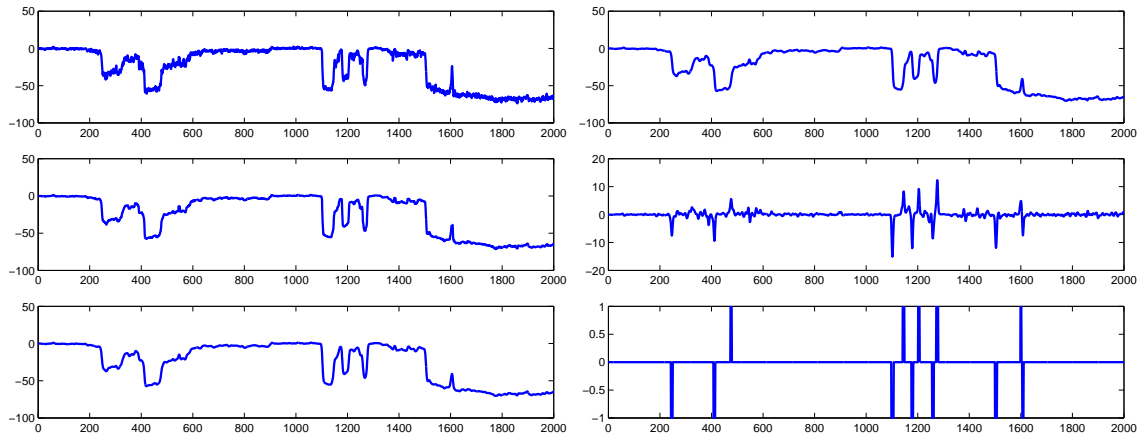


Figure 6: Left: head pose signal. Top : mean value of the tracker head pose distribution (in degree unit). Middle : after filtering with a median filter. Bottom : after filtering with A Gaussian filter. Right : top : head pose after filtering. Middle : derivatives. Bottom : dynamic head turn detection. -1 indicates a right turn, 1 a left head turn. All recognized head turns are correct, and no head turns are missing.

- dynamic head turns (both for pan and tilt angles)
- acknowledgment and head shaking (approval/disapproval), which are head gestures characterized by small oscillatory motion with specific frequencies
- focus of attention identification (specific to a particular set-up) : looking to a particular location (e.g. to the white-board, to the notes, or to other persons/speaker in the room).

Research for accurate and efficient (real-time if possible) algorithms will be undertaken in the upcoming months.

## 5 Acknowledgments

This work was funded by the Swiss NCCR on Interactive Multimodal Information Management (IM)<sup>2</sup>.

## References

- [1] L. Brown and Y. Tian. A study of coarse head pose estimation. *IEEE Workshop on Motion and Video Computing*, pages 125–130, Dec. 2002.
- [2] A. Doucet. On sequential monte carlo method for bayesian filtering. Technical report, University of Cambridge, 1998.
- [3] B. Kruger, S. Bruns, and G. Sommer. Efficient head pose estimation with gabor wavelet. *Proc. of 11th British Machine Vision Conference*, pages 11–14, Sept. 2000.

- [4] S. Niyogi and W. Freeman. Example-based head tracking. *Proc. Int. Conf. on Auto. Face and Gesture Rec.*, Oct. 1996.
- [5] P. Perez, C. Hue, J. Vermaak, and M. Gangnet. Color based probabilistic tracking. *European Conference on Computer Vision*, pages 661–675, 2002.
- [6] P. Phillips, P. R. H. Moon, and S. Rizvi. The feret evaluation methodology for face recognition algorithms. *IEEE Trans. on Pat. Anal. and Machine Intelligence*, 22(10), Oct. 2000.
- [7] R. Rae and H. Ritter. Recognition of human head orientation based on artificial neural networks. *IEEE Trans. on Neural Network*, 9(2):257–265, March 1998.
- [8] T. Sim and S. Baker. The cmu pose, illumination, and expression database. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 50, Oct. 2003.
- [9] K. Toyama and A. Blake. Probabilistic tracking in metric space. *Proc. 7th Int. Conf. on Computer Vision*, Dec. 2001.
- [10] Y. Wu and K. Toyama. Wide range illumination insensitive head orientation estimation. *IEEE Conf. on Automatic Face and Gesture Recognition*, Apr. 2001.
- [11] L. Zhao, G. Pingai, and I. Carlbom. Real-time head orientation estimation using neural networks. *Proc. Int. Conf. on Image Processing*, Sept. 2002.