

# Tracking of Multiple Targets Using Online Learning for Reference Model Adaptation

Franz Pernkopf

**Abstract**—Recently, much work has been done in multiple object tracking on the one hand and on reference model adaptation for a single-object tracker on the other side. In this paper, we do both tracking of multiple objects (faces of people) in a meeting scenario and online learning to incrementally update the models of the tracked objects to account for appearance changes during tracking. Additionally, we automatically initialize and terminate tracking of individual objects based on low-level features, i.e., face color, face size, and object movement. Many methods unlike our approach assume that the target region has been initialized by hand in the first frame. For tracking, a particle filter is incorporated to propagate sample distributions over time. We discuss the close relationship between our implemented tracker based on particle filters and genetic algorithms. Numerous experiments on meeting data demonstrate the capabilities of our tracking approach. Additionally, we provide an empirical verification of the reference model learning during tracking of indoor and outdoor scenes which supports a more robust tracking. Therefore, we report the average of the standard deviation of the trajectories over numerous tracking runs depending on the learning rate.

**Index Terms**—Genetic algorithms (GAs), multiple target tracking, particle filter, reference model learning, visual tracking.

## I. INTRODUCTION

VISUAL tracking of multiple objects is concerned with maintaining the correct identity and location of a variable number of objects over time irrespective of occlusions and visual alterations. Lim *et al.* [1] differentiate between intrinsic and extrinsic appearance variability including pose variation, shape deformation of the object and illumination change, camera movement, occlusions, respectively.

In the past few years, particle filters have become the method of choice for tracking. Isard and Blake [2] introduced particle filtering (condensation algorithm). Many different sampling schemes have been suggested in the meantime. An overview about sampling schemes of particle filters and the relation to Kalman filters is provided in [3].

Recently, the main emphasis is on simultaneously tracking multiple objects and on online learning to adapt the reference models to the appearance changes, e.g., pose variation, illumination change. Lim *et al.* [1] introduce a single-object tracker, where the target representation—a low-dimensional eigenspace representation—is incrementally updated to model the appear-

ance variability. They assume, like most tracking algorithms, that the target region is initialized by hand in the first frame. Jepson *et al.* [4] use a Gaussian mixture model which is adapted using an online expectation maximization (EM) algorithm to account for appearance changes. Their *WSL* tracker uses a wavelet-based object model which is useful for tracking objects where regions of the objects (i.e., faces) are stable while other regions vary, e.g., mouth. McKenna *et al.* [5] employ Gaussian mixtures of the color distributions of the objects as adaptive model. In [6], simple color histograms are used to represent the objects (similar as in [7]). However, they introduce a simple update of the histograms to overcome the appearance changes of the object. All the aforementioned articles are focused on tracking a single object. For tracking multiple objects, most algorithms belong to one of the following three categories: 1) Multiple instances of a single-object tracker are used [8]. 2) All objects of interest are included in the state space [9]. A fixed number of objects is assumed. Varying number of objects result in a dynamic change of the dimension of the state space. 3) Most recently, the framework of particle filters is extended to capture multiple targets using a mixture model [10]. This mixture particle filter—where each component models an individual object—enables interaction between the components by the importance weights. In [11], this approach is extended by the Adaboost algorithm to learn the models of the targets. The information from Adaboost enables detection of objects entering the scene automatically. The mixture particle filter is further extended in [12] to handle mutual occlusions. They introduce a rectification technique to compensate for camera motions, a global nearest neighbor data association method to correctly identify object detections with existing tracks, and a mean-shift algorithm which accounts for more stable trajectories for reliable motion prediction.

In this paper, we do both tracking of multiple persons in a meeting scenario and online adaptation of the models to account for appearance changes during tracking. The tracking is based on low-level features such as skin color, object motion, and object size. Based on these features, automatic initialization and termination of objects are performed. The aim is to use as little prior knowledge as possible. For tracking, a particle filter is incorporated to propagate sample distributions over time. Our implementation is related to the *dual estimation* problem [13], where both the states of multiple objects and the parameters of the reference models are simultaneously estimated given the observations. At every time step, the particle filter estimates the states using the observation likelihood of the current reference models while the online learning of the reference models is based on the current state estimates. Additionally, we discuss

Manuscript received July 24, 2007; revised January 30, 2008. This work was supported by the Austrian Science Fund under Project P19737-N15.

The author is with the Department of Electrical Engineering, Laboratory of Signal Processing and Speech Communication, Graz University of Technology, 8010 Graz, Austria (e-mail: pernkopf@tugraz.at).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TSMCB.2008.927281

94 the similarity between our implemented tracker based on parti-  
 95 cle filters and genetic algorithms (GAs). We want to emphasize  
 96 this close connection since approaches what have indepen-  
 97 dently been developed in one community might turn out to be  
 98 very useful for the other community and vice versa. Numerous  
 99 experiments on meeting data demonstrate the capabilities of our  
 100 tracking approach. Additionally, we empirically show that the  
 101 adaptation of the reference model during tracking of an indoor  
 102 and outdoor scenes results in a more robust tracking. For this,  
 103 we report the average of the standard deviation of the trajecto-  
 104 ries over numerous independent tracking runs depending on the  
 105 learning rate.

106 The proposed approach differs from previous methods in  
 107 several aspects. Recently, much work has been done in multiple  
 108 object tracking on the one hand side and on reference model  
 109 adaptation for a single-object tracker on the other side. In this  
 110 paper, we do both tracking of multiple objects and online learn-  
 111 ing to incrementally update the representation of the tracked ob-  
 112 jects to model appearance changes. We use the Jensen–Shannon  
 113 (JS) divergence [14] to measure the similarity between the  
 114 tracked object and its reference model. Additionally, we discuss  
 115 its advantages compared to the Kullback–Leibler divergence  
 116 [15] and the Bhattacharyya similarity coefficient [16]. We auto-  
 117 matically initialize and terminate tracking of individual objects  
 118 based on low-level features, i.e., face color, face size, and object  
 119 movement. Many methods unlike our approach assume that the  
 120 target region has been initialized in the first frame.

121 This paper is organized as follows. Section II introduces  
 122 the particle filter for multiple object tracking, the state-space  
 123 dynamics, the observation model, automatic initialization and  
 124 termination of objects, and the online learning of the mod-  
 125 els for the tracked objects. Section II-G summarizes the im-  
 126 plemented tracker on the basis of pseudocode. Section III  
 127 sketches the relationship to GA. The tracking results on a  
 128 meeting scenario and for indoor/outdoor scenes are presented in  
 129 Section IV. Additionally, we provide empirical verification of  
 130 the reference model learning in this section. Section V con-  
 131 cludes this paper.

## 132 II. TRACKING USING PARTICLE FILTERS

133 In many applications the states of a dynamic system have  
 134 to be estimated from a time series of noisy observations. The  
 135 Kalman filter [13], [17] is a linear dynamical system [18] that  
 136 provides a linear time-discrete filter that estimates the states  
 137 online over time once observations become available. This  
 138 filter is recursive in a sense that each current state estimate  
 139 is computed from the previous estimate and the current ob-  
 140 served data. In contrast to linear dynamical systems, the hidden  
 141 Markov model [19] assumes a discrete state space. Recently,  
 142 many extensions of the basic linear dynamical system have  
 143 been proposed [13] to overcome the assumption of the linear-  
 144 Gaussian model used for the observations and state transition,  
 145 e.g., the extended Kalman filter, unscented Kalman filter, or  
 146 the switching state-space model [20]. Another approach for  
 147 filtering is to use sequential Monte Carlo methods which are  
 148 also known as particle filters [21]. They are capable to deal with  
 149 any nonlinearity or distribution.

### A. Particle Filter

150

A particle filter is capable to deal with nonlinear non- 151  
 Gaussian processes and has become popular for visual tracking. 152  
 For tracking, the probability distribution that the object is in 153  
 state  $\mathbf{x}_t$  at time  $t$  given the observations  $\mathbf{y}_{0:t}$  up to time  $t$  is of 154  
 interest. Hence,  $p(\mathbf{x}_t|\mathbf{y}_{0:t})$  has to be constructed starting from 155  
 the initial distribution  $p(\mathbf{x}_0|\mathbf{y}_0) = p(\mathbf{x}_0)$ . In Bayesian filtering, 156  
 this can be formulated as iterative recursive process consisting 157  
 of the prediction step 158

$$p(\mathbf{x}_t|\mathbf{y}_{0:t-1}) = \int p(\mathbf{x}_t|\mathbf{x}_{t-1})p(\mathbf{x}_{t-1}|\mathbf{y}_{0:t-1})d\mathbf{x}_{t-1} \quad (1)$$

and of the filtering step 159

$$p(\mathbf{x}_t|\mathbf{y}_{0:t}) = \frac{p(\mathbf{y}_t|\mathbf{x}_t)p(\mathbf{x}_t|\mathbf{y}_{0:t-1})}{\int p(\mathbf{y}_t|\mathbf{x}_t)p(\mathbf{x}_t|\mathbf{y}_{0:t-1})d\mathbf{x}_t} \quad (2)$$

where  $p(\mathbf{x}_t|\mathbf{x}_{t-1})$  is the dynamic model describing the state- 160  
 space evolution which corresponds to the evolution of the 161  
 tracked object (see Section II-B) and  $p(\mathbf{y}_t|\mathbf{x}_t)$  is the likelihood 162  
 of an observation  $\mathbf{y}_t$  given the state  $\mathbf{x}_t$  (see observation model 163  
 in Section II-C). 164

In particle filters  $p(\mathbf{x}_t|\mathbf{y}_{0:t})$  of the filtering step is ap- 165  
 proximated by a finite set of weighted samples, i.e., the 166  
 particles,  $\{\mathbf{x}_t^m, w_t^m\}_{m=1}^M$ , where  $M$  is the number of sam- 167  
 ples. Particles are sampled from a proposal distribution  $\mathbf{x}_t^m \sim$  168  
 $q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{y}_{0:t})$  (importance sampling) [3]. In each iteration, 169  
 the importance weights are updated according to 170

$$w_t^m \propto \frac{p(\mathbf{y}_t|\mathbf{x}_t^m)p(\mathbf{x}_t^m|\mathbf{x}_{t-1}^m)}{q(\mathbf{x}_t^m|\mathbf{x}_{t-1}^m, \mathbf{y}_{0:t})}w_{t-1}^m \sum_{m=1}^M w_t^m = 1. \quad (3)$$

One simple choice for the proposal distribution is to take the 171  
 prior density  $q(\mathbf{x}_t^m|\mathbf{x}_{t-1}^m, \mathbf{y}_{0:t}) = p(\mathbf{x}_t^m|\mathbf{x}_{t-1}^m)$  (bootstrap filter). 172  
 Hence, the weights are proportional to the likelihood model 173  
 $p(\mathbf{y}_t|\mathbf{x}_t^m)$  174

$$w_t^m \propto p(\mathbf{y}_t|\mathbf{x}_t^m)w_{t-1}^m. \quad (4)$$

The posterior filtered density  $p(\mathbf{x}_t|\mathbf{y}_{1:t})$  can be approx- 175  
 imated as 176

$$p(\mathbf{x}_t|\mathbf{y}_{1:t}) \approx \sum_{m=1}^M w_t^m \delta(\mathbf{x}_t - \mathbf{x}_t^m) \quad (5)$$

where  $\delta(\mathbf{x}_t - \mathbf{x}_t^m)$  is the Dirac delta function with mass at  $\mathbf{x}_t^m$ . 177

We use resampling to reduce the *degeneracy problem* [3], 178  
 [21]. We resample the particles  $\{\mathbf{x}_t^m\}_{m=1}^M$  with replacement  $M$  179  
 times according to their weights  $w_t^m$ . The resulting particles 180  
 $\{\mathbf{x}_t^m\}_{m=1}^M$  have uniformly distributed weights  $w_t^m = 1/M$ . 181  
 Similar to the sampling importance resampling filter [3], we 182  
 resample in every time step. This simplifies (4) to  $w_t^m \propto$  183  
 $p(\mathbf{y}_t|\mathbf{x}_t^m)$  since  $w_{t-1}^m = 1/M \quad \forall m$ . 184

In the meeting scenario, we are interested in tracking the 185  
 faces of multiple people. We treat the tracking of multiple 186  
 objects completely independent, i.e., we assign a set of  $M$  187  
 particles to each tracked object  $k$  as  $\{\{\mathbf{x}_t^{m,k}\}_{m=1}^M\}_{k=1}^K$ , where 188  
 $K$  is the total number of tracked objects which dynamically 189

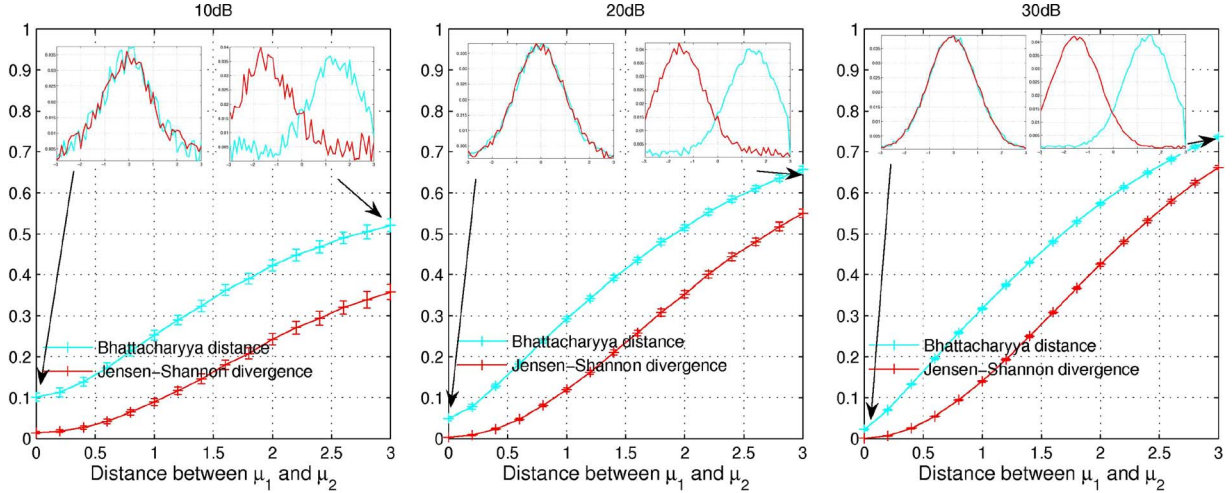


Fig. 1. JS divergence and Bhattacharyya similarity coefficient between two distributions estimated from samples. We added noise at a level of 10, 20, and 30 dB to the distributions.

190 changes over time. Hence, we use multiple instances of a single-  
191 object tracker similar to [8].

### 192 B. State-Space Dynamics

193 The state sequence evolution  $\{\mathbf{x}_t; t \in \mathbb{N}\}$  is assumed to be  
194 a second-order autoregressive process which is used instead  
195 of the first-order formalism ( $p(\mathbf{x}_t|\mathbf{x}_{t-1})$ ) introduced in the  
196 previous section. The second-order dynamics can be written as  
197 first order by extending the state vector at time  $t$  with elements  
198 from the state vector at time  $t-1$ .

199 We define the state vector at time  $t$  as  $\mathbf{x}_t = [x_t \ y_t \ s_t^x \ s_t^y]^T$ .  
200 The location of the target at  $t$  is given as  $x_t, y_t$ , respectively,  
201 and  $s_t^x, s_t^y$  denote the scale of the tracked region in the  $x \times y$   
202 image space. In our tracking approach, the transition model  
203 corresponds to

$$\mathbf{x}_{t+1}^{m,k} = \mathbf{x}_t^{m,k} + C\mathbf{v}_t + \frac{D}{2M} \sum_{m'=1}^M (\mathbf{x}_t^{m',k} - \mathbf{x}_{t-1}^{m',k}) \quad (6)$$

204 where  $\mathbf{v}_t \sim \mathcal{N}(0, \mathbf{I})$  is a simple Gaussian random noise model  
205 and the term  $1/2M \sum_{m'=1}^M (\mathbf{x}_t^{m',k} - \mathbf{x}_{t-1}^{m',k})$  captures the linear  
206 evolution of object  $k$  from the particles of the previous time  
207 step. Factor  $D$  models the influence of the linear evolution,  
208 e.g.,  $D$  is set to 0.5. The parameters of the random noise  
209 model are set to  $C = \text{diag}([10 \ 10 \ 0.03 \ 0.03])$  with the  
210 units of [pixel/frame], [pixel/frame], [1/frame], and [1/frame],  
211 respectively.

### 212 C. Observation Model

213 The shape of the tracked region is determined to be an ellipse  
214 [4] since the tracking is focused on the faces of the individuals.  
215 We assume that the principal axes of the ellipses are aligned  
216 with the coordinate axes of the image. Similarly to [7], we use  
217 the color histograms for modeling the target regions. Therefore,  
218 we transform the image into the hue-saturation-value (HSV)  
219 space [22]. For the sake of readability, we abuse the notation  
220 and write the particle  $\mathbf{x}_t^{m,k}$  as  $\mathbf{x}_t$  in this section. We build  
221 an individual histogram for hue (H)  $h_H^{\mathbf{x}_t}$ , saturation (S)  $h_S^{\mathbf{x}_t}$ ,

and value (V)  $h_V^{\mathbf{x}_t}$  of the elliptic candidate region at  $\mathbf{x}_t$ . The 222  
length of the principal axes of the ellipse are  $A_{\text{ref}}^k s_t^x$  and  $B_{\text{ref}}^k s_t^y$ , 223  
respectively, where  $A_{\text{ref}}^k$  and  $B_{\text{ref}}^k$  are the length of the ellipse 224  
axes of the reference model of object  $k$ . 225

The likelihood of the observation  $k$  model (likelihood model) 226  
 $p(\mathbf{y}_t^{m,k}|\mathbf{x}_t^{m,k})$  must be large for candidate regions with a his- 227  
togram close to the reference histogram. Therefore, we intro- 228  
duce the JS divergence [14] to measure the similarity between 229  
the normalized candidate and reference histograms,  $h_c^{\mathbf{x}_t}$  and 230  
 $h_{c,\text{ref}}^k$ ,  $c \in \{H, S, V\}$ , respectively. Since, the JS divergence 231  
is defined for probability distributions the histograms are nor- 232  
malized, i.e.,  $\sum_N h_c^{\mathbf{x}_t} = 1$ , where  $N$  denotes the number of 233  
histogram bins. In contrast to the Kullback–Leibler divergence 234  
[15], the JS divergence is symmetric and bounded between 0 235  
and 1. The JS divergence between the normalized histograms is 236  
defined as 237

$$\text{JS}_\pi(h_c^{\mathbf{x}_t}, h_{c,\text{ref}}^k) = H(\pi_1 h_c^{\mathbf{x}_t} + \pi_2 h_{c,\text{ref}}^k) - \pi_1 H(h_c^{\mathbf{x}_t}) - \pi_2 H(h_{c,\text{ref}}^k) \quad (7)$$

where  $\pi_1 + \pi_2 = 1, \pi_i \geq 0$  and the function  $H(\cdot)$  is the entropy 238  
[15]. The JS divergence is computed for the histograms of the 239  
H, S, and V space, and the observation likelihood is 240

$$p(\mathbf{y}_t^{m,k}|\mathbf{x}_t^{m,k}) \propto \exp -\lambda \left[ \sum_{c \in \{H, S, V\}} \text{JS}_\pi(h_c^{\mathbf{x}_t}, h_{c,\text{ref}}^k) \right] \quad (8)$$

where parameter  $\lambda$  is chosen to be five and the weight  $\pi_i$  is 241  
uniformly distributed. The number of bins of the histograms is 242  
set to  $N = 50$ . The JS divergence provides a lower and upper 243  
bound to the Bayes error and  $\pi_1$  and  $\pi_2$  can be viewed as 244  
*a priori* probabilities in a classification problem [14]. 245

In contrast to the often used Bhattacharyya similarity coef- 246  
ficient  $\sqrt{1 - \sum_N \sqrt{h_c^{\mathbf{x}_t} h_{c,\text{ref}}^k}}$  [16], the JS divergence is not 247  
so sensitive to local perturbations in the histogram (noise). This 248  
is shown in Fig. 1 where we compute the JS divergence and 249  
Bhattacharyya similarity coefficient on synthetic data. There- 250  
fore, we sample two Gaussian distributions with  $\mu_1 = -\mu_2$ , 251

252 where  $\mu_1$  varies from 0 to 1.5, and unit variance. Noise is added  
 253 to those distributions at a level of 10, 20, and 30 dB. Plots are  
 254 averaged over 100 independent simulations.

#### 255 D. Automatic Initialization of Objects

256 If an object enters the frame, a set of  $M$  particles and a refer-  
 257 ence histogram for this object have to be initialized. Basically,  
 258 the initialization of objects is automatically performed using the  
 259 following simple low-level features.

260 1) Motion: The images are transformed to gray scale  $I_{x_t, y_t}^G$ .  
 261 The motion feature is determined for each pixel located  
 262 at  $x, y$  by the standard deviation over a time window  
 263  $T_w$  as  $\sigma_{x, y}^t = \sigma(I_{x_t - T_w : t, y_t - T_w : t}^G)$ . Applying an adaptive  
 264 threshold  $T_{\text{motion}} = 1/10 \max_{x, y \in IG} \sigma_{x, y}^t$  pixels with  
 265 a value larger  $T_{\text{motion}}$  belong to regions where movement  
 266 happens. However,  $\max_{x, y \in IG} \sigma_{x, y}^t$  has to be sufficiently  
 267 large so that motion exists at all. A binary motion image  
 268  $I_{x_t, y_t}^{B_{\text{motion}}}$  after morphological closing is shown in Fig. 2.

269 2) Skin Color: The skin color of the people is modeled  
 270 by a Gaussian mixture model [23] in the HSV  
 271 color space. A Gaussian mixture model  $p(\mathbf{z}|\Theta)$  is the  
 272 weighted sum of  $L > 1$  Gaussian components,  $p(\mathbf{z}|\Theta) =$   
 273  $\sum_{l=1}^L \alpha_l \mathcal{N}(\mathbf{z}|\mu_l, \Sigma_l)$ , where  $\mathbf{z} = [z_H, z_S, z_V]^T$  is the 3-D  
 274 color vector of one image pixel,  $\alpha_l$  corresponds to the  
 275 weight of each component  $l = 1, \dots, L$ . These weights  
 276 are constrained to be positive  $\alpha_l \geq 0$  and  $\sum_{l=1}^L \alpha_l = 1$ .  
 277 The Gaussian mixture is specified by the set of parameters  
 278  $\Theta = \{\alpha_l, \mu_l, \Sigma_l\}_{l=1}^L$ . These parameters are determined  
 279 by the EM algorithm [24] from a face database.

280 Image pixels  $\mathbf{z} \in I_{x_t, y_t}^{\text{HSV}}$  are classified according to their  
 281 likelihood  $p(\mathbf{z}|\Theta)$  using a threshold  $T_{\text{skin}}$ . The binary  
 282 map  $I_{x_t, y_t}^{B_{\text{skin}}}$  filtered with a morphological closing operator  
 283 is presented in Fig. 2.

284 3) Object Size: We initialize a new object only for skin-  
 285 colored moving regions with a size larger than  $T_{\text{Area}}$ .  
 286 Additionally, we do not allow initialization of a new set of  
 287 particles in regions where currently an object is tracked.  
 288 To this end, a binary map  $I_{x_t, y_t}^{B_{\text{prohibited}}}$  represents the areas  
 289 where initialization is prohibited. The binary combination  
 290 of all images  $I_{x_t, y_t}^B = I_{x_t, y_t}^{B_{\text{motion}}} \cap I_{x_t, y_t}^{B_{\text{skin}}} \cap I_{x_t, y_t}^{B_{\text{prohibited}}}$  is  
 291 used for extracting regions with an area larger  $T_{\text{Area}}$ . Tar-  
 292 get objects are initialized for those regions, i.e., the ellipse  
 293 size  $(A_{\text{ref}}^k, B_{\text{ref}}^k)$  and the histograms  $h_{c, \text{ref}}^k, c \in \{H, S, V\}$   
 294 are determined from the region of the bounding ellipse.

295 Fig. 2 shows an example for the initialization of a new object.  
 296 The original image  $I_{x_t, y_t}^{\text{HSV}}$  is presented in (a). A person entering  
 297 from the right side should be initialized. A second person in  
 298 the middle of the image is already tracked. The binary images  
 299 of the thresholded motion  $I_{x_t, y_t}^{B_{\text{motion}}}$  and the skin-colored areas  
 300  $I_{x_t, y_t}^{B_{\text{skin}}}$  are shown in (b) and (c), respectively. The reflections at  
 301 the table and the movement of the curtain produce noise in the  
 302 motion image. The color of the table and chairs intersects with  
 303 the skin-color model. To guarantee successful initialization the  
 304 lower part of the image—the region of the chairs and desk—has  
 305 to be excluded. This is reasonable since nobody can enter in  
 306 this area. Also, tracking is performed in the area above the

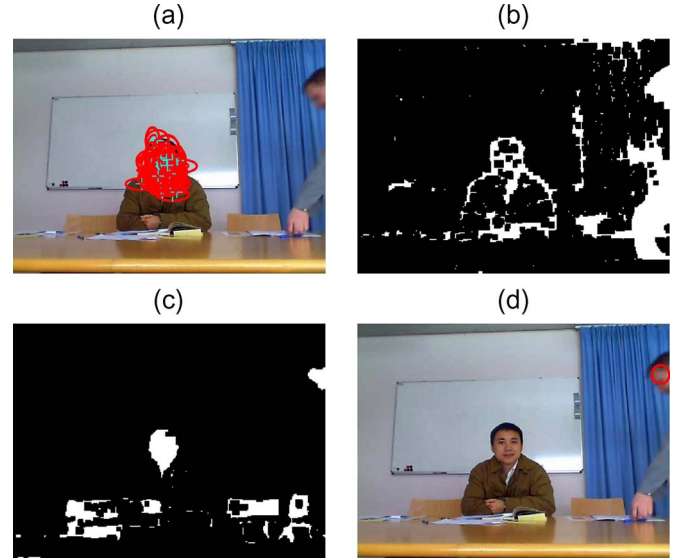


Fig. 2. Initialization of a new object. (a) Original image with one object already tracked. (b) Binary image of the thresholded motion  $I_{x_t, y_t}^{B_{\text{motion}}}$ . (c) Binary image of the skin-colored areas  $I_{x_t, y_t}^{B_{\text{skin}}}$ . (d) Image with region of initialized object.

chairs only. Finally, the region of the new initialized object is  
 307 presented as ellipse in (d). Resizing of the images is performed  
 308 for computing the features to speed up the initialization of  
 309 objects. 310

311 1) Shortcomings: The objects are initialized when they en-  
 312 ter the image. The reference histogram is taken during this  
 313 initialization. There are the following shortcomings during  
 314 initialization. 314

- 315 1) The camera is focused on the people sitting at the table  
 316 and not on people walking behind the chairs. This means  
 317 that walking persons appear blurred. 317
- 318 2) Entering persons are moving relatively fast. This also  
 319 results in a degraded image quality (blurring). 319
- 320 3) During initialization, we normally get the side view of  
 321 the person's head. When the person sits at the table the  
 322 reference histogram is not necessarily a good model for  
 323 the frontal view. 323

324 To deal with these shortcomings, we propose online learning  
 325 to incrementally update the reference models of the tracked  
 326 objects over time (see Section II-F). We perform this only in  
 327 cases where no mutual occlusions between the tracked objects  
 328 are existent. 328

#### 329 E. Automatic Termination of Objects

330 Termination of particles is performed if the observation  
 331 likelihood  $p(\mathbf{y}_t^{m, k} | \mathbf{x}_t^{m, k})$  at state  $\mathbf{x}_t^{m, k}$  drops below a predefined  
 332 threshold  $T_{\text{Kill}}$  (e.g., 0.001), i.e., 332

$$p(\mathbf{y}_t^{m, k} | \mathbf{x}_t^{m, k}) = \begin{cases} 0, & \text{if } p(\mathbf{y}_t^{m, k} | \mathbf{x}_t^{m, k}) < T_{\text{Kill}} \\ p(\mathbf{y}_t^{m, k} | \mathbf{x}_t^{m, k}), & \text{otherwise.} \end{cases} \quad (9)$$

333 Particles with zero probability do not survive during resam-  
 334 pling. If the tracked object leaves the field of view all  $M$  334

335 particles of an object  $k$  are removed, i.e.,  $p(\mathbf{y}_t^{m,k} | \mathbf{x}_t^{m,k}) = 0$   
 336 for all particles of object  $k$ .

### 337 F. Incremental Learning of Object Models

338 To handle the appearance change of the tracked objects over  
 339 time, we use online learning to adapt the reference histograms  
 340  $h_{c,\text{ref}}^k$ ,  $c \in \{H, S, V\}$  (similar to [6]) and ellipse size  $A_{\text{ref}}^k$  and  
 341  $B_{\text{ref}}^k$ . Therefore, a learning rate  $\alpha$  is introduced, and the model  
 342 parameters for target object  $k$  are updated according to

$$h_{c,\text{ref}}^k = \alpha \hat{h}_c^k + (1 - \alpha) h_{c,\text{ref}}^k, \quad c \in \{H, S, V\} \quad (10)$$

$$A_{\text{ref}}^k = \alpha \hat{A}^k + (1 - \alpha) A_{\text{ref}}^k \quad (11)$$

$$B_{\text{ref}}^k = \alpha \hat{B}^k + (1 - \alpha) B_{\text{ref}}^k \quad (12)$$

343 where  $\hat{h}_c^k$  denotes the histogram and  $\hat{A}^k$  and  $\hat{B}^k$  are the prin-  
 344 cipal axes of the bounding ellipse of the nonoccluded (i.e., no  
 345 mutual occlusion between tracked objects) skin-colored region  
 346 of the corresponding tracked object  $k$  located at  $\{\mathbf{x}_t^{m,k}\}_{m=1}^M$ .  
 347 Again, this region has to be larger than  $T_{\text{Area}}$ . No update of  
 348 the reference models is performed in the case where occlusion  
 349 between the tracked objects occurs or the skin-colored region  
 350 is not large enough. The latter condition is a simple way to  
 351 ensure that the model update is only conducted for faces.  
 352 This simplistic assumption can be appropriately extended by  
 353 integrating more advanced face models.

354 The learning rate  $\alpha$  introduces an *exponential forgetting*  
 355 *process*, i.e., the contribution of a specific object exponentially  
 356 decreases as it recedes into the past. Currently, the learning rate  
 357 (value between 0 and 1) is fixed (a good value has been selected  
 358 during experiments). However,  $\alpha$  could be adapted depending  
 359 on the dynamics of the scene.

#### 360 **Algorithm 1** Particle Filter Tracking

361 **Input:**  $I_{x_0:T, y_0:T}^{\text{HSV}}$  (Color image sequence  $0 : T$ ),

362 Skin-color model  $\Theta$

363 **Parameters:**  $M, N, \lambda, C, D, T_w, T_{\text{motion}}, T_{\text{skin}}, T_{\text{Area}},$

364  $T_{\text{Kill}}, \alpha$

365 **Output:**  $\{\{\mathbf{x}_{0:T}^{m,k}\}_{m=1}^M\}_{\forall k}$

366  $t \leftarrow 0$

367  $k \leftarrow 0$

368 **while** InitObjects **do**

369  $k \leftarrow k + 1$

370 Obtain:  $h_{c,\text{ref}}^k : c \in \{H, S, V\}, A_{\text{ref}}^k, B_{\text{ref}}^k, \mathbf{x}_{\text{ref}}^k$

371  $\mathbf{x}_{t+1}^{m,k} \leftarrow \mathbf{x}_{\text{ref}}^k + C\mathbf{v}_t \quad \forall m = 1, \dots, M$  (Generate particles)

372 **end while**

373  $K \leftarrow k$

374 **for**  $t = 1$  **to**  $T$  **do**

375  $w_t^{m,k} \propto p(\mathbf{y}_t^{m,k} | \mathbf{x}_t^{m,k})$

376  $\forall k = 1, \dots, K \quad \forall m = 1, \dots, M$

377 **while** KillObjects **do**

378  $k \leftarrow$  Determine object to terminate

379 Remove  $M$  particles  $x_t^{m,k}$  of object  $k$

380 Remove reference histogram and ellipse size:

381  $h_{c,\text{ref}}^k : c \in \{H, S, V\}, A_{\text{ref}}^k, B_{\text{ref}}^k$

382  $K \leftarrow K - 1$

383 **end while**



Fig. 3. Tracking scene. We track and initialize objects in the red rectangle.

**for**  $k = 1$  **to**  $K$  **do** 384

$w_t^{m,k} \leftarrow w_t^{m,k} / \sum_{m'=1}^M w_t^{m',k} \quad \forall m = 1, \dots, M$  385

$\{\mathbf{x}_t^{m,k}\}_{m=1}^M \leftarrow$  Resampling 386

(with replacement):  $\{\mathbf{x}_t^{m,k}, w_t^{m,k}\}_{m=1}^M$  387

$\mathbf{x}_{t+1}^{m,k} \leftarrow \mathbf{x}_t^{m,k} + C\mathbf{v}_t + (D/2M) \sum_{m'=1}^M (\mathbf{x}_t^{m',k} - \mathbf{x}_{t-1}^{m',k})$  388

$\forall m = 1, \dots, M$  (Apply state-space dynamics) 389

**if** OnlineUpdate **then** 390

Determine:  $\hat{h}_c^k : c \in \{H, S, V\}, \hat{A}^k, \hat{B}^k$  391

$h_{c,\text{ref}}^k \leftarrow \alpha \hat{h}_c^k + (1 - \alpha) h_{c,\text{ref}}^k \quad c \in \{H, S, V\}$  392

$A_{\text{ref}}^k \leftarrow \alpha \hat{A}^k + (1 - \alpha) A_{\text{ref}}^k$  393

$B_{\text{ref}}^k \leftarrow \alpha \hat{B}^k + (1 - \alpha) B_{\text{ref}}^k$  394

**end if** 395

**end for** 396

**while** InitObjects **do** 397

$K \leftarrow K + 1$  398

Obtain:  $h_{c,\text{ref}}^K : c \in \{H, S, V\}, A_{\text{ref}}^K, B_{\text{ref}}^K, \mathbf{x}_{\text{ref}}^K$  399

$\mathbf{x}_{t+1}^{m,K} \leftarrow \mathbf{x}_{\text{ref}}^K + C\mathbf{v}_t \quad \forall m = 1, \dots, M$  (Generate 400

particles) 401

**end while** 402

**end for** 403

### G. Implemented Tracker 404

In the following, we sketch our tracking approach for multi- 405  
 ple objects (see Algorithm 1). The binary variable *InitObject* 406  
 denotes that a new object for tracking has been detected. 407  
*KillObject* is set if an object should be terminated. *OnlineUp-* 408  
*date* indicates that object  $k$  located at  $\{\mathbf{x}_t^{m,k}\}_{m=1}^M$  is nonoc- 409  
 cluded, and the area of the skin-colored region is larger than 410  
 $T_{\text{Area}}$ , i.e., we perform online learning for reference model  $k$ . 411

Our implementation is related to the *dual estimation* problem 412  
 [13], where both the states of multiple objects  $\mathbf{x}_t^{m,k}$  and the 413  
 parameters of the reference models are simultaneously esti- 414  
 mated given the observations. At every time step, the particle 415  
 filter estimates the states using the observation likelihood of 416  
 the current reference models, while the online learning of the 417  
 reference models is based on the current state estimates. 418

## III. RELATIONSHIP TO GAS 419

GAs are optimization algorithms founded upon the principles 420  
 of natural evolution discovered by Darwin. In nature, individ- 421  
 uals have to adapt to their environment in order to survive in 422



Fig. 4. Tracking of people. Frames: 1, 416, 430, 449, 463, 491, 583, 609, 622, 637, 774, 844, 967, 975, 1182, 1400 (the frame number is assigned from left to right and top to bottom).

423 a process of further development. An introduction of GAs can  
 424 be found in [25] and [26]. GA are stochastic procedures which  
 425 have been successfully applied in many optimization tasks.  
 426 GA operate on a population of potential solutions applying the  
 427 principle of *survival of the fittest individual* to produce better  
 428 and better approximations to the solution. At each generation, a  
 429 new set of approximations is created by the process of selecting  
 430 individuals according to their level of fitness in the problem  
 431 domain and assembling them together using operators inspired  
 432 from nature. This leads to the evolution of individuals that are  
 433 better suited to their environment than the parent individuals  
 434 they were created from. GA model the natural processes, such  
 435 as selection, recombination, and mutation. Starting from an  
 436 initial population  $P(0)$ , the sequence  $P(0), P(1), \dots, P(t)$ ,  
 437  $P(t + 1)$  is called population sequence or evolution. The end of  
 438 an artificial evolution process is reached once the termination  
 439 condition is met, and the result of the optimization task is  
 440 available.

441 In this section, we want to point to the close relationship  
 442 between GA and our particle filter for tracking. This analogy  
 443 has been mentioned in [27]. As suggested in Section II, we  
 444 treat the tracking of multiple objects completely independent,  
 445 i.e., we have a set of  $M$  particles for each object  $k$ . In the GA  
 446 framework, we can relate this to  $k$  instantiations of GA, one  
 447 for each tracked object. Hence, each particle  $\mathbf{x}_t^m$  of object  $k$   
 448 represents one individual in the population  $P(t)$  which is value  
 449 encoded. The population size is  $M$ . A new genetic evolution

process is started once a new object is initialized for tracking 450  
 (InitObject). The evolution process of the GA is terminated 451  
 either at the end of the video ( $t = T$ ) or when the set of 452  
 individuals is not supported by the fitness value (KillObject). 453  
 The observation likelihood  $p(\mathbf{y}_t^{m,k} | \mathbf{x}_t^{m,k})$  denotes the fitness 454  
 function to evaluate the individuals. However, the scope of GA 455  
 for tracking is slightly different. GA are generally used to find a 456  
 set of parameters for a given optimization task, i.e., the aim is to 457  
 find the individual with the best fitness after the termination of 458  
 the GA. Whereas, in the tracking case, the focus lies on the evo- 459  
 lution of the individuals, i.e., the trajectory of the tracked object. 460

The selection operator directs the search toward promising 461  
 regions in the search space. *Roulette Wheel Selection* [28] is a 462  
 widely used selection method which is very similar to sampling 463  
 with replacement as used in Section II. To each individual, a re- 464  
 production probability according to  $w_t^m \leftarrow w_t^m / \sum_{m'=1}^M w_t^{m'}$  465  
 is assigned. A roulette wheel is constructed with a slot size cor- 466  
 responding to the individuals reproduction probability. Then, 467  
 $M$  uniformly distributed random numbers on the interval  $[0, 1]$  468  
 are drawn and distributed according to their value around the 469  
 wheel. The slots where they are placed to compose the subse- 470  
 quent population  $P(t)$ . The state-space dynamics of the particle 471  
 filter (see Section II-B) is modeled by the recombination and 472  
 mutation operator. 473

The framework of the GA for tracking one object  $k$  is 474  
 presented in Algorithm 1. The incremental learning of the 475  
 reference model is omitted for the sake of brevity. 476



Fig. 5. Partial occlusions. Frames: 468, 616, 974, 4363 (the frame number is assigned from left to right and top to bottom).

477 **Algorithm 2** GA Tracking  
 478 **Input:**  $I_{x_{t:T}, y_{t:T}}^{\text{HSV}}$  (Color image sequence  $t : T$ ),  
 479 **Parameters:**  $M, N, \lambda, C, D, T_{\text{Kill}}$   
 480 **Output:**  $\{\mathbf{x}_{t:T}^m\}_{m=1}^M$  (set of particle sequences  $t : T$ )  
 481 Initialize population  $P(t)$  :  
 482  $\mathbf{x}_t^m \leftarrow \mathbf{x}_{ref} + C\mathbf{v}_t \quad \forall m = 1, \dots, M$   
 483 **while**  $\text{KillObject} \cap t < T$  (Loop over image sequence) **do**  
 484 Evaluate individuals:  
 485  $w_t^m \leftarrow p(\mathbf{y}_t^m | \mathbf{x}_t^m) \quad \forall m = 1, \dots, M$   
 486 Selection  $P(t)$ :  
 487  $\{\mathbf{x}_t^m\}_{m=1}^M \leftarrow$  (Sampling with replacement)  $\{\mathbf{x}_t^m, w_t^m\}_{m=1}^M$   
 488 Recombination  $P(t+1)$  :  
 489  $\mathbf{x}_{t+1}^m \leftarrow \mathbf{x}_t^m + (D/2M) \sum_{m'=1}^M (\mathbf{x}_t^{m'} - \mathbf{x}_t^{m-1})$   
 490  $\forall m = 1, \dots, M$   
 491 Mutation  $P(t+1)$  :  $\mathbf{x}_{t+1}^m \leftarrow \mathbf{x}_{t+1}^m + C\mathbf{v}_t \quad \forall m = 1, \dots, M$   
 492  $t \leftarrow t + 1$   
 493 **end while**

#### 494 IV. EXPERIMENTS

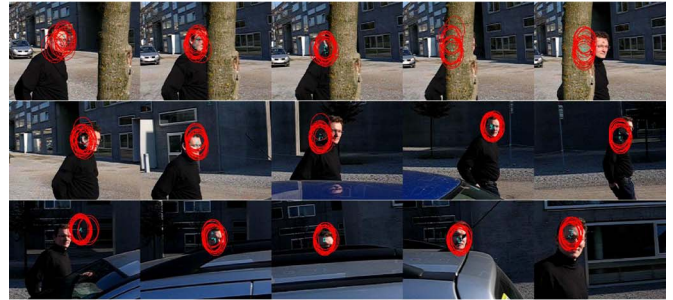
495 We present tracking results on meeting data in Section IV-A  
 496 where we do both tracking of multiple persons and on-  
 497 line adaptation of the reference models during tracking. In  
 498 Section IV-B, we empirically show that the adaptation of the  
 499 reference model during tracking (single object) of an indoor  
 500 and outdoor scene results in a more robust tracking. Finally, in  
 501 Section IV-C, tracking results using reference model adaptation  
 502 for multiple objects of an outdoor scene are presented. For the  
 503 outdoor scenes, we report the average standard deviation of  
 504 the trajectories of independent tracking runs depending on the  
 505 learning rate  $\alpha$ .

##### 506 A. Meeting Scenario

507 The meeting room layout is shown in Fig. 3. The red rec-  
 508 tangle [region of interest (ROI)] in the image marks the frame  
 509 where tracking and initialization of objects are performed. Peo-  
 510 ple may enter and leave on both sides of the image. Currently,  
 511 our tracker initializes a new target even if it enters from the



(a)



(b)



(c)

Fig. 6. Outdoor tracking. Frames: 7, 11, 12, 13, 14, 20, 42, 63, 80, 107, 136, 146, 158, 165, 192 (the frame number is assigned from left to right and top to bottom). (a) Original image sequence. (b) Tracking without reference model adaptation ( $\alpha = 0$ ). (c) Tracking with online reference model learning ( $\alpha = 0.2$ ).

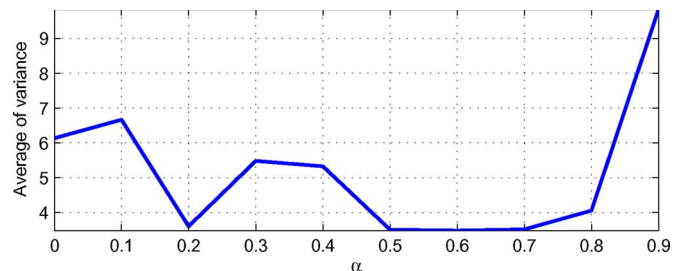


Fig. 7. Averaged standard deviation of the trajectories of 100 tracking runs depending on the reference model learning rate  $\alpha$ .

bottom, e.g., a hand moving from the table into the ROI. The 512  
 strong reflections at the table, chairs, and the white board cause 513  
 noise in the motion image. 514

For testing the performance of our tracking approach, ten 515  
 videos with  $\sim 7000$  frames have been used. The resolution is 516  
 $640 \times 480$  pixels. The meeting room is equipped with a table 517  
 and three chairs. We have different persons in each video. The 518  
 people are coming from both sides into the frame moving 519

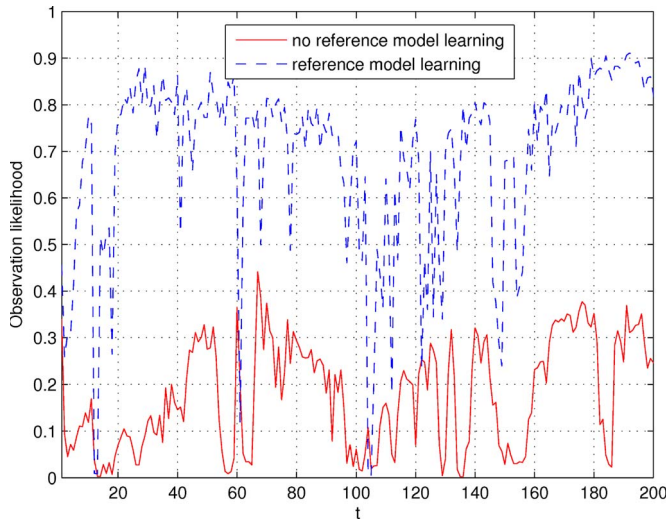


Fig. 8. Observation likelihood of outdoor sequence.

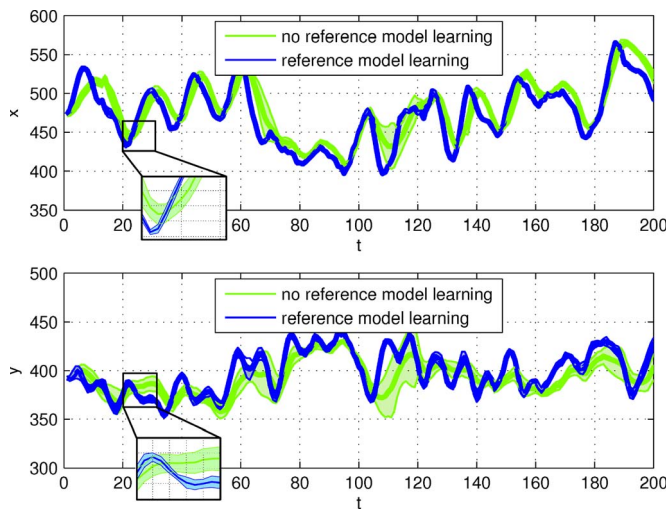


Fig. 9. Averaged trajectory with standard deviation in  $x$  and  $y$  of outdoor sequence (over ten runs).

520 to chairs and sit down. After a short discussion, people are  
 521 sequentially leaving the room, are coming back, sit down at  
 522 different chairs and so on. At the beginning, people may already  
 523 sit at the chairs. In this case, we have to automatically initialize  
 524 multiple objects at the very first frame.

525 Fig. 4 shows the result of the implemented tracker for one  
 526 video. All the initializations and terminations of objects are  
 527 performed automatically. The appearance of an object changes  
 528 over time. When entering the frame, we get the side view of  
 529 the person's head. After sitting down at the table, we have a  
 530 frontal view. We account for this by incrementally updating the  
 531 reference histogram during tracking. We perform this only in  
 532 the case where no mutual occlusions with other tracked objects  
 533 are existent. The participants were successfully tracked over  
 534 long image sequences.

535 First, the person on the left side stands up and leaves the room  
 536 on the right side (frame 416–491). When walking behind the  
 537 two sitting people, partial occlusions occur which do not cause  
 538 problems. Next, the person on the right (frame 583–637) leaves  
 539 the room on the left side. His face is again partially occluded



(a)



(b)



(c)

Fig. 10. Indoor tracking. Frames: 1, 12, 24, 31, 38, 41, 47, 54, 65, 71, 80, 107, 113, 120, 134 (the frame number is assigned from left to right and top to bottom). (a) Original image sequence. (b) Tracking without reference model adaptation ( $\alpha = 0$ ). (c) Tracking with online reference model learning ( $\alpha = 0.2$ ).

by the person in the middle. Then, the person on the center  
 540 chair leaves the room (frame 774). After that, a person on the  
 541 right side enters and sits at the left chair (frame 844). At frame  
 542 967, a small person is entering and moving to the chair in the  
 543 middle. Here, again, a partial occlusion occurs at frame 975, 544  
 544 which is also tackled. Finally, a person enters from the right  
 545 and sits down on the right chair (frame 1182, 1400). The partial  
 546 occlusions are shown in Fig. 5. Also, the blurred face of the  
 547 moving person in the back can be observed in this figure. The  
 548 reference model adaptation enables a more robust tracking. If  
 549 we do not update the models of the tracked objects over time,  
 550 the tracking fails in case of these partial occlusions. In [29],  
 551 occlusions are handled using multiple cameras for tracking  
 552 participants in a meeting. 553

## B. Reference Model Adaptation for Single-Object Tracking 554

In the following, we show the benefit of the reference model  
 555 adaptation during tracking of a short indoor and outdoor se-  
 556 quence. In contrast to the meeting scenario, we restrict the  
 557



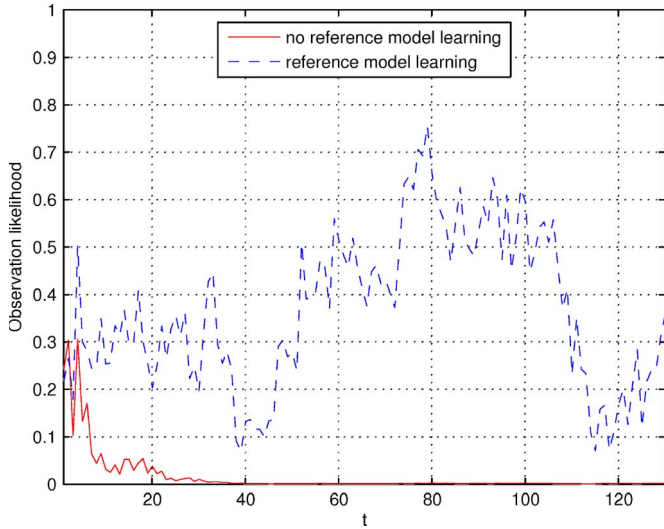


Fig. 11. Observation likelihood of indoor sequence.

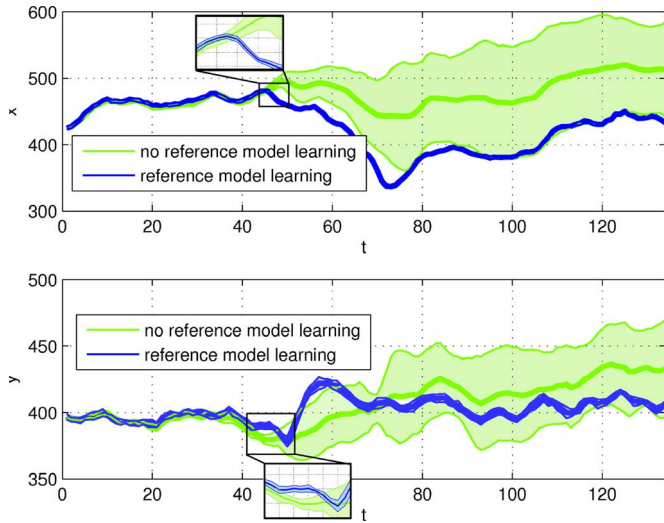


Fig. 12. Averaged trajectory with standard deviation in  $x$  and  $y$  of indoor sequence (over ten runs).

558 tracking to one single object, i.e., face. This means, in particular, that the automatic initialization and termination of the object is disabled. The object is initialized by hand in the first 561 frame.

562 Fig. 6(a) shows a short outdoor sequence where a person is moving behind a tree and two cars with strongly changing lighting conditions. We have a total occlusion of the face in frames 565 12 and 13 and a partial occluded face in frames 146 to 165. We repeated the tracking without and with reference model learning 567 ten times, and a typical result is shown in Fig. 6(b) and (c), 568 respectively. We use  $M = 50$  particles for tracking, whereas only 15 particles with the best observation likelihood are shown 570 in the figures.

571 In Fig. 7, we present the average standard deviation of the trajectories over 100 tracking runs. The reference model learning rate  $\alpha$  has been chosen in the range of  $0, \dots, 0.6$  574 (0 means that there is no learning). The optimal learning rate with respect to a low standard deviation of the trajectories over 576 100 independent runs is  $\alpha = 0.2$  for this outdoor sequence.



(a)



(b)



(c)

Fig. 13. Outdoor tracking of multiple objects. Frames: 1, 12, 30, 47, 49, 51, 53, 57, 59, 79, 105, 107, 109, 111, 149 (the frame number is assigned from left to right and top to bottom). (a) Original image sequence. (b) Tracking without reference model adaptation ( $\alpha = 0$ ). (c) Tracking with online reference model learning ( $\alpha = 0.1$ ).

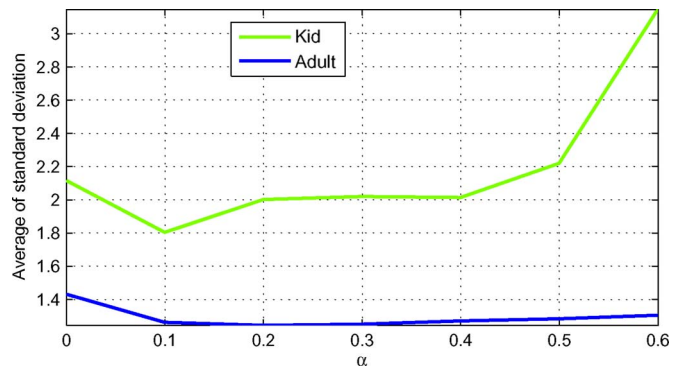


Fig. 14. Averaged standard deviation of the trajectories of ten tracking runs depending on the reference model learning rate  $\alpha$ .

Fig. 8 shows the observation likelihood of the best particle 577 during tracking. At the complete occlusion (frames  $t = 12$  and 578  $t = 13$ ) and the partial occlusion (frames  $t = 145, \dots, 160$ ), the 579 observation likelihood drops, however, with reference model 580 learning a quick recovery is supported. 581

Fig. 9 summarizes the averaged trajectory with the standard deviation over ten different tracking runs performed for the outdoor scene. In the case of reference model learning, we observe in the video sequences that the tracking of the face gives highly similar trajectories. The standard deviation is small and approximately constant over time. However, if no learning of the reference model is performed, the standard deviation is large in certain time segments. This leads to the conclusion that model adaptation results in a more robust tracking.

Fig. 10(a) shows an indoor video where a person is moving on a corridor, and a tree causes partial occlusion of the tracked face. Additionally, the lighting conditions are strongly varying. The face is partially occluded by the tree in frames 37–50 and 110–126. Again, the tracking without and with reference model learning is repeated ten times, and a typical result is shown in Fig. 10(b) and (c), respectively. Only 15 particles with the best observation likelihood are visualized. The parameter setting is the same as in the previous experiments. The tracker without reference model refinement fails during the first occlusion in all ten runs, whereas the tracker with online model update is successful in all cases. The optimal learning rate  $\alpha$  is set to 0.2 (established during experiments).

This can be also observed in the observation likelihood of the best particle over time (see Fig. 11) and in the averaged trajectory over ten tracking results (see Fig. 12).

### C. Reference Model Adaptation for Multiple Object Tracking

We show tracking results for an outdoor scene where a kid is showing an adult dancing steps (see Fig. 13). A typical tracking result without and with reference model learning is shown in Fig. 13(b) and (c), respectively. Again,  $M = 50$  particles are used, whereas only 15 particles with the best observation likelihood are shown in the figures. Similar as in the previous section, we did a repeatability test, i.e., we tracked the objects over ten independent runs. The tracked objects are initialized by hand in the very first frame.

Fig. 14 shows the average standard deviation of the trajectories of ten tracking runs using a learning rate  $\alpha$  in the range of 0, . . . , 0.6. The optimal learning rate for the *Kid* and the *Adult* is  $\alpha = 0.1$  and  $\alpha = 0.2$ , respectively. Currently,  $\alpha$  is fixed for the whole image sequence. Ideally,  $\alpha$  could be adapted depending on the dynamics of the scene.

623

## V. CONCLUSION

We propose a robust visual tracking algorithm for multiple objects (faces of people) in a meeting scenario based on low-level features as skin color, target motion, and target size. Based on these features, automatic initialization and termination of objects is performed. For tracking a sampling importance resampling, particle filter has been used to propagate sample distributions over time. Furthermore, we use online learning of the target models to handle the appearance variability of the objects. We discuss the similarity between our implemented tracker and GAs. Each particle represents an individual in the GA framework. The evaluation function incorporates the observation likelihood model and the individual selection

process maps to the resampling procedure in the particle filter. The state-space dynamics is incorporated in the recombination and mutation operator of the GA. Numerous experiments on meeting data show the capabilities of the tracking approach. The participants were successfully tracked over long image sequences. Partial occlusions are handled by the algorithm. Additionally, we empirically show that the adaptation of the reference model during tracking of indoor and outdoor scenes results in a more robust tracking.

Future work concentrates on extending the tracker to other scenarios and to investigate an adaptive reference model learning rate  $\alpha$  which depends on the dynamics of the scene. Furthermore, we aim to develop approaches for tackling occlusions.

## ACKNOWLEDGMENT

This work was supported by the Austrian Science Fund (Project S106). The author would like to thank M. Grabner and M. Kepesi who collected the data during their involvement in the MISTRAL Project ([www.mistral-project.at](http://www.mistral-project.at)). The MISTRAL Project was funded by the Austrian Research Promotion Agency ([www.ffg.at](http://www.ffg.at)) within the strategic objective FIT-IT under Project 809264/9338. The author would also like to thank C. Kirchstätter for recording the indoor and outdoor video.

## REFERENCES

- [1] J. Lim, D. Ross, R.-S. Lin, and M.-H. Yang, "Incremental learning for visual tracking," in *Advances in Neural Information Processing Systems 17*. Cambridge, MA: MIT Press, 2005.
- [2] M. Isard and A. Blake, "Condensation—Conditional density propagation for visual tracking," *Int. J. Comput. Vis.*, vol. 29, no. 1, pp. 5–28, 1998.
- [3] S. Arulampalam, S. Maskell, N. Gordon, and T. Clapp, "A tutorial on particle filters for on-line non-linear/non-Gaussian Bayesian tracking," *IEEE Trans. Signal Process.*, vol. 50, no. 2, pp. 174–188, Feb. 2002.
- [4] A. D. Jepson, D. J. Fleet, and T. F. El-Maraghi, "Robust online appearance models for visual tracking," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 25, no. 10, pp. 1296–1311, Oct. 2003.
- [5] S. McKenna, Y. Raja, and S. Gong, "Tracking colour objects using adaptive mixture models," *Image Vis. Comput.*, vol. 17, no. 3/4, pp. 225–231, 1999.
- [6] K. Nummiaro, E. Koller-Meier, and L. Van Gool, "An adaptive color-based particle filter," *Image Vis. Comput.*, vol. 21, no. 1, pp. 99–110, Jan. 2003.
- [7] P. Pérez, C. Hue, J. Vermaak, and M. Gangnet, "Color-based probabilistic tracking," in *Proc. ECCV*, 2002, pp. 661–675.
- [8] S. L. Dockstader and A. Tekalp, "Tracking multiple objects in the presence of articulated and occluded motion," in *Proc. Workshop Human Motion*, 2000, pp. 88–98.
- [9] C. Hue, J.-P. Le Cadre, and P. Pérez, "Tracking multiple objects with particle filtering," *IEEE Trans. Aerosp. Electron. Syst.*, vol. 38, no. 3, pp. 791–812, Jul. 2002.
- [10] J. Vermaak, A. Doucet, and P. Pérez, "Maintaining multi-modality through mixture tracking," in *Proc. ICCV*, 2003, pp. 1110–1116.
- [11] K. Okuma, A. Taleghani, N. de Freitas, J. J. Little, and D. G. Lowe, "A boosted particle filter: Multitarget detection and tracking," in *Proc. ECCV*, 2004, pp. 28–39.
- [12] Y. Cai, N. de Freitas, and J. J. Little, "Robust visual tracking for multiple targets," in *Proc. ECCV*, 2006, pp. 107–118.
- [13] S. Haykin, *Kalman Filtering and Neural Networks*. Hoboken, NJ: Wiley, 2001.
- [14] J. Lin, "Divergence measures based on the Shannon entropy," *IEEE Trans. Inf. Theory*, vol. 37, no. 1, pp. 145–151, Jan. 1991.
- [15] T. Cover and J. Thomas, *Elements of Information Theory*. Hoboken, NJ: Wiley, 1991.
- [16] D. Comaniciu, V. Ramesh, and P. Meer, "Real-time tracking of non-rigid objects using mean shift," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2005, pp. 142–149.

- 700 [17] R. E. Kalman, "A new approach to linear filtering and prediction prob-  
701 lems," *Trans. ASME, Ser. D, J. Basic Eng.*, vol. 82, pp. 34–45, 1960.
- 702 [18] C. M. Bishop, *Pattern Recognition and Machine Learning*. Berlin,  
703 Germany: Springer Sci.+Bus. Media, 2006.
- 704 [19] L. R. Rabiner, "A tutorial on hidden Markov models and selected appli-  
705 cations in speech recognition," *Proc. IEEE*, vol. 77, no. 2, pp. 257–286,  
706 Feb. 1989.
- 707 [20] Z. Ghahramani and G. E. Hinton, "Variational learning for switching  
708 state-space models," *Neural Comput.*, vol. 12, no. 4, pp. 963–996, 1998.
- 709 [21] A. Doucet, "On sequential Monte Carlo sampling methods for Bayesian  
710 filtering," Dept. Eng., Cambridge Univ., London, U.K., Tech. Rep.  
711 CUED/F-INFENG/TR. 310, 1998.
- 712 [22] M. Sonka, V. Hlavac, and R. Boyle, *Image Processing, Analysis, and  
713 Machine Vision*. London, U.K.: Int. Thomson, 1999.
- 714 [23] R. O. Duda, P. E. Hart, and D. G. Stork, *Pattern Classification*. Hoboken,  
715 NJ: Wiley, 2000.
- 716 [24] A. Dempster, N. Laird, and D. Rubin, "Maximum likelihood estimation  
717 from incomplete data via the EM algorithm," *J. R. Stat. Soc., Ser. B, Stat.  
718 Methodol.*, vol. 39, pp. 1–38, 1977.
- 719 [25] D. E. Goldberg, *Genetic Algorithms in Search, Optimization and Machine  
720 Learning*. Reading, MA: Addison-Wesley, 1989.
- 721 [26] T. Bäck, *Evolutionary Algorithms in Theory and Practice*. London,  
722 U.K.: Oxford Univ. Press, 1996.
- [27] T. Higuchi, "Monte Carlo filter using the genetic algorithm operators," 723  
*J. Stat. Comput. Simul.*, vol. 59, no. 1, pp. 1–23, Aug. 1997. 724
- [28] L. E. Baker, "Reducing bias and inefficiency in the selection algorithm," 725  
in *Proc. Int. Conf. Genetic Algorithms Appl.*, 1987, pp. 14–21. 726
- [29] H. Nait Charif and S. J. McKenna, "Tracking the activity of participants 727  
in a meeting," *Mach. Vis. Appl.*, vol. 17, no. 2, pp. 83–93, 2006. 728



**Franz Pernkopf** received the M.Sc. (Dipl.Ing.) de- 729  
gree in electrical engineering from the Graz Univer- 730  
sity of Technology, Graz, Austria, in 1999, and the 731  
Ph.D. degree from the University of Leoben, Leoben, 732  
Austria, in 2002. 733

He was a Research Associate with the Department 734  
of Electrical Engineering, University of Washington, 735  
Seattle, from 2004 to 2006. Currently, he is an 736  
Assistant Professor with the Signal Processing and 737  
Speech Communication Laboratory, Graz University 738  
of Technology. His research interests include ma- 739  
chine learning, Bayesian networks, feature selection, finite mixture models, 740  
vision, speech, and statistical pattern recognition. 741  
Dr. Pernkopf was awarded the Erwin Schrödinger Fellowship in 2002. 742

# Tracking of Multiple Targets Using Online Learning for Reference Model Adaptation

Franz Pernkopf

**Abstract**—Recently, much work has been done in multiple object tracking on the one hand and on reference model adaptation for a single-object tracker on the other side. In this paper, we do both tracking of multiple objects (faces of people) in a meeting scenario and online learning to incrementally update the models of the tracked objects to account for appearance changes during tracking. Additionally, we automatically initialize and terminate tracking of individual objects based on low-level features, i.e., face color, face size, and object movement. Many methods unlike our approach assume that the target region has been initialized by hand in the first frame. For tracking, a particle filter is incorporated to propagate sample distributions over time. We discuss the close relationship between our implemented tracker based on particle filters and genetic algorithms. Numerous experiments on meeting data demonstrate the capabilities of our tracking approach. Additionally, we provide an empirical verification of the reference model learning during tracking of indoor and outdoor scenes which supports a more robust tracking. Therefore, we report the average of the standard deviation of the trajectories over numerous tracking runs depending on the learning rate.

**Index Terms**—Genetic algorithms (GAs), multiple target tracking, particle filter, reference model learning, visual tracking.

## I. INTRODUCTION

VISUAL tracking of multiple objects is concerned with maintaining the correct identity and location of a variable number of objects over time irrespective of occlusions and visual alterations. Lim *et al.* [1] differentiate between intrinsic and extrinsic appearance variability including pose variation, shape deformation of the object and illumination change, camera movement, occlusions, respectively.

In the past few years, particle filters have become the method of choice for tracking. Isard and Blake [2] introduced particle filtering (condensation algorithm). Many different sampling schemes have been suggested in the meantime. An overview about sampling schemes of particle filters and the relation to Kalman filters is provided in [3].

Recently, the main emphasis is on simultaneously tracking multiple objects and on online learning to adapt the reference models to the appearance changes, e.g., pose variation, illumination change. Lim *et al.* [1] introduce a single-object tracker, where the target representation—a low-dimensional eigenspace representation—is incrementally updated to model the appear-

ance variability. They assume, like most tracking algorithms, that the target region is initialized by hand in the first frame. Jepson *et al.* [4] use a Gaussian mixture model which is adapted using an online expectation maximization (EM) algorithm to account for appearance changes. Their  $WSL$  tracker uses a wavelet-based object model which is useful for tracking objects where regions of the objects (i.e., faces) are stable while other regions vary, e.g., mouth. McKenna *et al.* [5] employ Gaussian mixtures of the color distributions of the objects as adaptive model. In [6], simple color histograms are used to represent the objects (similar as in [7]). However, they introduce a simple update of the histograms to overcome the appearance changes of the object. All the aforementioned articles are focused on tracking a single object. For tracking multiple objects, most algorithms belong to one of the following three categories: 1) Multiple instances of a single-object tracker are used [8]. 2) All objects of interest are included in the state space [9]. A fixed number of objects is assumed. Varying number of objects result in a dynamic change of the dimension of the state space. 3) Most recently, the framework of particle filters is extended to capture multiple targets using a mixture model [10]. This mixture particle filter—where each component models an individual object—enables interaction between the components by the importance weights. In [11], this approach is extended by the Adaboost algorithm to learn the models of the targets. The information from Adaboost enables detection of objects entering the scene automatically. The mixture particle filter is further extended in [12] to handle mutual occlusions. They introduce a rectification technique to compensate for camera motions, a global nearest neighbor data association method to correctly identify object detections with existing tracks, and a mean-shift algorithm which accounts for more stable trajectories for reliable motion prediction.

In this paper, we do both tracking of multiple persons in a meeting scenario and online adaptation of the models to account for appearance changes during tracking. The tracking is based on low-level features such as skin color, object motion, and object size. Based on these features, automatic initialization and termination of objects are performed. The aim is to use as little prior knowledge as possible. For tracking, a particle filter is incorporated to propagate sample distributions over time. Our implementation is related to the *dual estimation* problem [13], where both the states of multiple objects and the parameters of the reference models are simultaneously estimated given the observations. At every time step, the particle filter estimates the states using the observation likelihood of the current reference models while the online learning of the reference models is based on the current state estimates. Additionally, we discuss

Manuscript received July 24, 2007; revised January 30, 2008. This work was supported by the Austrian Science Fund under Project P19737-N15.

The author is with the Department of Electrical Engineering, Laboratory of Signal Processing and Speech Communication, Graz University of Technology, 8010 Graz, Austria (e-mail: pernkopf@tugraz.at).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TSMCB.2008.927281

94 the similarity between our implemented tracker based on parti-  
 95 cle filters and genetic algorithms (GAs). We want to emphasize  
 96 this close connection since approaches what have indepen-  
 97 dently been developed in one community might turn out to be  
 98 very useful for the other community and vice versa. Numerous  
 99 experiments on meeting data demonstrate the capabilities of our  
 100 tracking approach. Additionally, we empirically show that the  
 101 adaptation of the reference model during tracking of an indoor  
 102 and outdoor scenes results in a more robust tracking. For this,  
 103 we report the average of the standard deviation of the trajecto-  
 104 ries over numerous independent tracking runs depending on the  
 105 learning rate.

106 The proposed approach differs from previous methods in  
 107 several aspects. Recently, much work has been done in multiple  
 108 object tracking on the one hand side and on reference model  
 109 adaptation for a single-object tracker on the other side. In this  
 110 paper, we do both tracking of multiple objects and online learn-  
 111 ing to incrementally update the representation of the tracked ob-  
 112 jects to model appearance changes. We use the Jensen–Shannon  
 113 (JS) divergence [14] to measure the similarity between the  
 114 tracked object and its reference model. Additionally, we discuss  
 115 its advantages compared to the Kullback–Leibler divergence  
 116 [15] and the Bhattacharyya similarity coefficient [16]. We auto-  
 117 matically initialize and terminate tracking of individual objects  
 118 based on low-level features, i.e., face color, face size, and object  
 119 movement. Many methods unlike our approach assume that the  
 120 target region has been initialized in the first frame.

121 This paper is organized as follows. Section II introduces  
 122 the particle filter for multiple object tracking, the state-space  
 123 dynamics, the observation model, automatic initialization and  
 124 termination of objects, and the online learning of the mod-  
 125 els for the tracked objects. Section II-G summarizes the im-  
 126 plemented tracker on the basis of pseudocode. Section III  
 127 sketches the relationship to GA. The tracking results on a  
 128 meeting scenario and for indoor/outdoor scenes are presented in  
 129 Section IV. Additionally, we provide empirical verification of  
 130 the reference model learning in this section. Section V con-  
 131 cludes this paper.

## 132 II. TRACKING USING PARTICLE FILTERS

133 In many applications the states of a dynamic system have  
 134 to be estimated from a time series of noisy observations. The  
 135 Kalman filter [13], [17] is a linear dynamical system [18] that  
 136 provides a linear time-discrete filter that estimates the states  
 137 online over time once observations become available. This  
 138 filter is recursive in a sense that each current state estimate  
 139 is computed from the previous estimate and the current ob-  
 140 served data. In contrast to linear dynamical systems, the hidden  
 141 Markov model [19] assumes a discrete state space. Recently,  
 142 many extensions of the basic linear dynamical system have  
 143 been proposed [13] to overcome the assumption of the linear-  
 144 Gaussian model used for the observations and state transition,  
 145 e.g., the extended Kalman filter, unscented Kalman filter, or  
 146 the switching state-space model [20]. Another approach for  
 147 filtering is to use sequential Monte Carlo methods which are  
 148 also known as particle filters [21]. They are capable to deal with  
 149 any nonlinearity or distribution.

### A. Particle Filter

150

A particle filter is capable to deal with nonlinear non- 151  
 Gaussian processes and has become popular for visual tracking. 152  
 For tracking, the probability distribution that the object is in 153  
 state  $\mathbf{x}_t$  at time  $t$  given the observations  $\mathbf{y}_{0:t}$  up to time  $t$  is of 154  
 interest. Hence,  $p(\mathbf{x}_t|\mathbf{y}_{0:t})$  has to be constructed starting from 155  
 the initial distribution  $p(\mathbf{x}_0|\mathbf{y}_0) = p(\mathbf{x}_0)$ . In Bayesian filtering, 156  
 this can be formulated as iterative recursive process consisting 157  
 of the prediction step 158

$$p(\mathbf{x}_t|\mathbf{y}_{0:t-1}) = \int p(\mathbf{x}_t|\mathbf{x}_{t-1})p(\mathbf{x}_{t-1}|\mathbf{y}_{0:t-1})d\mathbf{x}_{t-1} \quad (1)$$

and of the filtering step 159

$$p(\mathbf{x}_t|\mathbf{y}_{0:t}) = \frac{p(\mathbf{y}_t|\mathbf{x}_t)p(\mathbf{x}_t|\mathbf{y}_{0:t-1})}{\int p(\mathbf{y}_t|\mathbf{x}_t)p(\mathbf{x}_t|\mathbf{y}_{0:t-1})d\mathbf{x}_t} \quad (2)$$

where  $p(\mathbf{x}_t|\mathbf{x}_{t-1})$  is the dynamic model describing the state- 160  
 space evolution which corresponds to the evolution of the 161  
 tracked object (see Section II-B) and  $p(\mathbf{y}_t|\mathbf{x}_t)$  is the likelihood 162  
 of an observation  $\mathbf{y}_t$  given the state  $\mathbf{x}_t$  (see observation model 163  
 in Section II-C). 164

In particle filters  $p(\mathbf{x}_t|\mathbf{y}_{0:t})$  of the filtering step is ap- 165  
 proximated by a finite set of weighted samples, i.e., the 166  
 particles,  $\{\mathbf{x}_t^m, w_t^m\}_{m=1}^M$ , where  $M$  is the number of sam- 167  
 ples. Particles are sampled from a proposal distribution  $\mathbf{x}_t^m \sim$  168  
 $q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{y}_{0:t})$  (importance sampling) [3]. In each iteration, 169  
 the importance weights are updated according to 170

$$w_t^m \propto \frac{p(\mathbf{y}_t|\mathbf{x}_t^m)p(\mathbf{x}_t^m|\mathbf{x}_{t-1}^m)}{q(\mathbf{x}_t^m|\mathbf{x}_{t-1}^m, \mathbf{y}_{0:t})}w_{t-1}^m \sum_{m=1}^M w_t^m = 1. \quad (3)$$

One simple choice for the proposal distribution is to take the 171  
 prior density  $q(\mathbf{x}_t^m|\mathbf{x}_{t-1}^m, \mathbf{y}_{0:t}) = p(\mathbf{x}_t^m|\mathbf{x}_{t-1}^m)$  (bootstrap filter). 172  
 Hence, the weights are proportional to the likelihood model 173  
 $p(\mathbf{y}_t|\mathbf{x}_t^m)$  174

$$w_t^m \propto p(\mathbf{y}_t|\mathbf{x}_t^m)w_{t-1}^m. \quad (4)$$

The posterior filtered density  $p(\mathbf{x}_t|\mathbf{y}_{1:t})$  can be approx- 175  
 imated as 176

$$p(\mathbf{x}_t|\mathbf{y}_{1:t}) \approx \sum_{m=1}^M w_t^m \delta(\mathbf{x}_t - \mathbf{x}_t^m) \quad (5)$$

where  $\delta(\mathbf{x}_t - \mathbf{x}_t^m)$  is the Dirac delta function with mass at  $\mathbf{x}_t^m$ . 177

We use resampling to reduce the *degeneracy problem* [3], 178  
 [21]. We resample the particles  $\{\mathbf{x}_t^m\}_{m=1}^M$  with replacement  $M$  179  
 times according to their weights  $w_t^m$ . The resulting particles 180  
 $\{\mathbf{x}_t^m\}_{m=1}^M$  have uniformly distributed weights  $w_t^m = 1/M$ . 181  
 Similar to the sampling importance resampling filter [3], we 182  
 resample in every time step. This simplifies (4) to  $w_t^m \propto$  183  
 $p(\mathbf{y}_t|\mathbf{x}_t^m)$  since  $w_{t-1}^m = 1/M \quad \forall m$ . 184

In the meeting scenario, we are interested in tracking the 185  
 faces of multiple people. We treat the tracking of multiple 186  
 objects completely independent, i.e., we assign a set of  $M$  187  
 particles to each tracked object  $k$  as  $\{\{\mathbf{x}_t^{m,k}\}_{m=1}^M\}_{k=1}^K$ , where 188  
 $K$  is the total number of tracked objects which dynamically 189

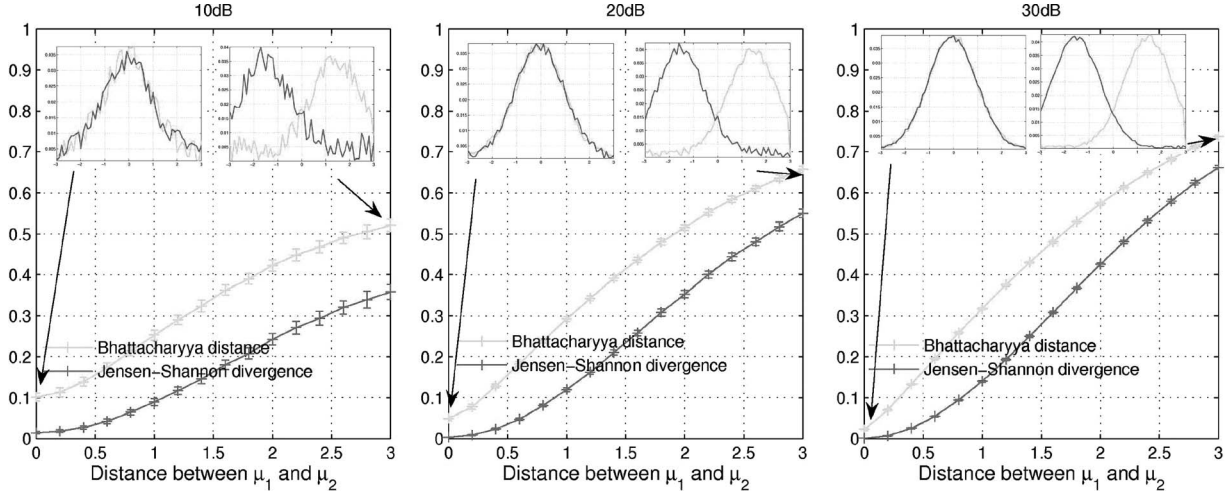


Fig. 1. JS divergence and Bhattacharyya similarity coefficient between two distributions estimated from samples. We added noise at a level of 10, 20, and 30 dB to the distributions.

190 changes over time. Hence, we use multiple instances of a single-  
191 object tracker similar to [8].

### 192 B. State-Space Dynamics

193 The state sequence evolution  $\{\mathbf{x}_t; t \in \mathbb{N}\}$  is assumed to be  
194 a second-order autoregressive process which is used instead  
195 of the first-order formalism ( $p(\mathbf{x}_t|\mathbf{x}_{t-1})$ ) introduced in the  
196 previous section. The second-order dynamics can be written as  
197 first order by extending the state vector at time  $t$  with elements  
198 from the state vector at time  $t-1$ .

199 We define the state vector at time  $t$  as  $\mathbf{x}_t = [x_t \ y_t \ s_t^x \ s_t^y]^T$ .  
200 The location of the target at  $t$  is given as  $x_t, y_t$ , respectively,  
201 and  $s_t^x, s_t^y$  denote the scale of the tracked region in the  $x \times y$   
202 image space. In our tracking approach, the transition model  
203 corresponds to

$$\mathbf{x}_{t+1}^{m,k} = \mathbf{x}_t^{m,k} + C\mathbf{v}_t + \frac{D}{2M} \sum_{m'=1}^M (\mathbf{x}_t^{m',k} - \mathbf{x}_{t-1}^{m',k}) \quad (6)$$

204 where  $\mathbf{v}_t \sim \mathcal{N}(0, \mathbf{I})$  is a simple Gaussian random noise model  
205 and the term  $1/2M \sum_{m'=1}^M (\mathbf{x}_t^{m',k} - \mathbf{x}_{t-1}^{m',k})$  captures the linear  
206 evolution of object  $k$  from the particles of the previous time  
207 step. Factor  $D$  models the influence of the linear evolution,  
208 e.g.,  $D$  is set to 0.5. The parameters of the random noise  
209 model are set to  $C = \text{diag}([10 \ 10 \ 0.03 \ 0.03])$  with the  
210 units of [pixel/frame], [pixel/frame], [1/frame], and [1/frame],  
211 respectively.

### 212 C. Observation Model

213 The shape of the tracked region is determined to be an ellipse  
214 [4] since the tracking is focused on the faces of the individuals.  
215 We assume that the principal axes of the ellipses are aligned  
216 with the coordinate axes of the image. Similarly to [7], we use  
217 the color histograms for modeling the target regions. Therefore,  
218 we transform the image into the hue-saturation-value (HSV)  
219 space [22]. For the sake of readability, we abuse the notation  
220 and write the particle  $\mathbf{x}_t^{m,k}$  as  $\mathbf{x}_t$  in this section. We build  
221 an individual histogram for hue (H)  $h_H^{\mathbf{x}_t}$ , saturation (S)  $h_S^{\mathbf{x}_t}$ ,

and value (V)  $h_V^{\mathbf{x}_t}$  of the elliptic candidate region at  $\mathbf{x}_t$ . The 222  
length of the principal axes of the ellipse are  $A_{\text{ref}}^k s_t^x$  and  $B_{\text{ref}}^k s_t^y$ , 223  
respectively, where  $A_{\text{ref}}^k$  and  $B_{\text{ref}}^k$  are the length of the ellipse 224  
axes of the reference model of object  $k$ . 225

The likelihood of the observation  $k$  model (likelihood model) 226  
 $p(\mathbf{y}_t^{m,k}|\mathbf{x}_t^{m,k})$  must be large for candidate regions with a his- 227  
togram close to the reference histogram. Therefore, we intro- 228  
duce the JS divergence [14] to measure the similarity between 229  
the normalized candidate and reference histograms,  $h_c^{\mathbf{x}_t}$  and 230  
 $h_{c,\text{ref}}^k$ ,  $c \in \{H, S, V\}$ , respectively. Since, the JS divergence 231  
is defined for probability distributions the histograms are nor- 232  
malized, i.e.,  $\sum_N h_c^{\mathbf{x}_t} = 1$ , where  $N$  denotes the number of 233  
histogram bins. In contrast to the Kullback–Leibler divergence 234  
[15], the JS divergence is symmetric and bounded between 0 235  
and 1. The JS divergence between the normalized histograms is 236  
defined as 237

$$\text{JS}_\pi(h_c^{\mathbf{x}_t}, h_{c,\text{ref}}^k) = H(\pi_1 h_c^{\mathbf{x}_t} + \pi_2 h_{c,\text{ref}}^k) - \pi_1 H(h_c^{\mathbf{x}_t}) - \pi_2 H(h_{c,\text{ref}}^k) \quad (7)$$

where  $\pi_1 + \pi_2 = 1, \pi_i \geq 0$  and the function  $H(\cdot)$  is the entropy 238  
[15]. The JS divergence is computed for the histograms of the 239  
H, S, and V space, and the observation likelihood is 240

$$p(\mathbf{y}_t^{m,k}|\mathbf{x}_t^{m,k}) \propto \exp -\lambda \left[ \sum_{c \in \{H, S, V\}} \text{JS}_\pi(h_c^{\mathbf{x}_t}, h_{c,\text{ref}}^k) \right] \quad (8)$$

where parameter  $\lambda$  is chosen to be five and the weight  $\pi_i$  is 241  
uniformly distributed. The number of bins of the histograms is 242  
set to  $N = 50$ . The JS divergence provides a lower and upper 243  
bound to the Bayes error and  $\pi_1$  and  $\pi_2$  can be viewed as 244  
*a priori* probabilities in a classification problem [14]. 245

In contrast to the often used Bhattacharyya similarity coef- 246  
ficient  $\sqrt{1 - \sum_N \sqrt{h_c^{\mathbf{x}_t} h_{c,\text{ref}}^k}}$  [16], the JS divergence is not 247  
so sensitive to local perturbations in the histogram (noise). This 248  
is shown in Fig. 1 where we compute the JS divergence and 249  
Bhattacharyya similarity coefficient on synthetic data. There- 250  
fore, we sample two Gaussian distributions with  $\mu_1 = -\mu_2$ , 251

252 where  $\mu_1$  varies from 0 to 1.5, and unit variance. Noise is added  
 253 to those distributions at a level of 10, 20, and 30 dB. Plots are  
 254 averaged over 100 independent simulations.

#### 255 D. Automatic Initialization of Objects

256 If an object enters the frame, a set of  $M$  particles and a refer-  
 257 ence histogram for this object have to be initialized. Basically,  
 258 the initialization of objects is automatically performed using the  
 259 following simple low-level features.

260 1) Motion: The images are transformed to gray scale  $I_{x_t, y_t}^G$ .  
 261 The motion feature is determined for each pixel located  
 262 at  $x, y$  by the standard deviation over a time window  
 263  $T_w$  as  $\sigma_{x,y}^t = \sigma(I_{x_t-T_w:t, y_t-T_w:t}^G)$ . Applying an adaptive  
 264 threshold  $T_{\text{motion}} = 1/10 \max_{x,y \in IG} \sigma_{x,y}^t$  pixels with  
 265 a value larger  $T_{\text{motion}}$  belong to regions where movement  
 266 happens. However,  $\max_{x,y \in IG} \sigma_{x,y}^t$  has to be sufficiently  
 267 large so that motion exists at all. A binary motion image  
 268  $I_{x_t, y_t}^{B_{\text{motion}}}$  after morphological closing is shown in Fig. 2.

269 2) Skin Color: The skin color of the people is modeled  
 270 by a Gaussian mixture model [23] in the HSV  
 271 color space. A Gaussian mixture model  $p(\mathbf{z}|\Theta)$  is the  
 272 weighted sum of  $L > 1$  Gaussian components,  $p(\mathbf{z}|\Theta) =$   
 273  $\sum_{l=1}^L \alpha_l \mathcal{N}(\mathbf{z}|\mu_l, \Sigma_l)$ , where  $\mathbf{z} = [z_H, z_S, z_V]^T$  is the 3-D  
 274 color vector of one image pixel,  $\alpha_l$  corresponds to the  
 275 weight of each component  $l = 1, \dots, L$ . These weights  
 276 are constrained to be positive  $\alpha_l \geq 0$  and  $\sum_{l=1}^L \alpha_l = 1$ .  
 277 The Gaussian mixture is specified by the set of parameters  
 278  $\Theta = \{\alpha_l, \mu_l, \Sigma_l\}_{l=1}^L$ . These parameters are determined  
 279 by the EM algorithm [24] from a face database.

280 Image pixels  $\mathbf{z} \in I_{x_t, y_t}^{\text{HSV}}$  are classified according to their  
 281 likelihood  $p(\mathbf{z}|\Theta)$  using a threshold  $T_{\text{skin}}$ . The binary  
 282 map  $I_{x_t, y_t}^{B_{\text{skin}}}$  filtered with a morphological closing operator  
 283 is presented in Fig. 2.

284 3) Object Size: We initialize a new object only for skin-  
 285 colored moving regions with a size larger than  $T_{\text{Area}}$ .  
 286 Additionally, we do not allow initialization of a new set of  
 287 particles in regions where currently an object is tracked.  
 288 To this end, a binary map  $I_{x_t, y_t}^{B_{\text{prohibited}}}$  represents the areas  
 289 where initialization is prohibited. The binary combination  
 290 of all images  $I_{x_t, y_t}^B = I_{x_t, y_t}^{B_{\text{motion}}} \cap I_{x_t, y_t}^{B_{\text{skin}}} \cap I_{x_t, y_t}^{B_{\text{prohibited}}}$  is  
 291 used for extracting regions with an area larger  $T_{\text{Area}}$ . Tar-  
 292 get objects are initialized for those regions, i.e., the ellipse  
 293 size  $(A_{\text{ref}}^k, B_{\text{ref}}^k)$  and the histograms  $h_{c, \text{ref}}^k, c \in \{H, S, V\}$   
 294 are determined from the region of the bounding ellipse.

295 Fig. 2 shows an example for the initialization of a new object.  
 296 The original image  $I_{x_t, y_t}^{\text{HSV}}$  is presented in (a). A person entering  
 297 from the right side should be initialized. A second person in  
 298 the middle of the image is already tracked. The binary images  
 299 of the thresholded motion  $I_{x_t, y_t}^{B_{\text{motion}}}$  and the skin-colored areas  
 300  $I_{x_t, y_t}^{B_{\text{skin}}}$  are shown in (b) and (c), respectively. The reflections at  
 301 the table and the movement of the curtain produce noise in the  
 302 motion image. The color of the table and chairs intersects with  
 303 the skin-color model. To guarantee successful initialization the  
 304 lower part of the image—the region of the chairs and desk—has  
 305 to be excluded. This is reasonable since nobody can enter in  
 306 this area. Also, tracking is performed in the area above the

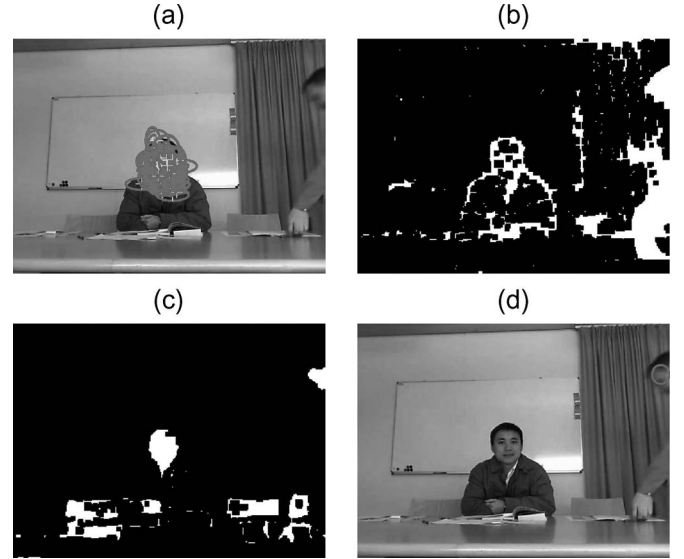


Fig. 2. Initialization of a new object. (a) Original image with one object already tracked. (b) Binary image of the thresholded motion  $I_{x_t, y_t}^{B_{\text{motion}}}$ . (c) Binary image of the skin-colored areas  $I_{x_t, y_t}^{B_{\text{skin}}}$ . (d) Image with region of initialized object.

chairs only. Finally, the region of the new initialized object is  
 307 presented as ellipse in (d). Resizing of the images is performed  
 308 for computing the features to speed up the initialization of  
 309 objects. 310

311 1) *Shortcomings*: The objects are initialized when they en-  
 312 ter the image. The reference histogram is taken during this  
 313 initialization. There are the following shortcomings during  
 314 initialization. 314

- 315 1) The camera is focused on the people sitting at the table  
 316 and not on people walking behind the chairs. This means  
 317 that walking persons appear blurred. 317
- 318 2) Entering persons are moving relatively fast. This also  
 319 results in a degraded image quality (blurring). 319
- 320 3) During initialization, we normally get the side view of  
 321 the person's head. When the person sits at the table the  
 322 reference histogram is not necessarily a good model for  
 323 the frontal view. 323

324 To deal with these shortcomings, we propose online learning  
 325 to incrementally update the reference models of the tracked  
 326 objects over time (see Section II-F). We perform this only in  
 327 cases where no mutual occlusions between the tracked objects  
 328 are existent. 328

#### 329 E. Automatic Termination of Objects

330 Termination of particles is performed if the observation  
 331 likelihood  $p(\mathbf{y}_t^{m,k} | \mathbf{x}_t^{m,k})$  at state  $\mathbf{x}_t^{m,k}$  drops below a predefined  
 332 threshold  $T_{\text{Kill}}$  (e.g., 0.001), i.e., 332

$$p(\mathbf{y}_t^{m,k} | \mathbf{x}_t^{m,k}) = \begin{cases} 0, & \text{if } p(\mathbf{y}_t^{m,k} | \mathbf{x}_t^{m,k}) < T_{\text{Kill}} \\ p(\mathbf{y}_t^{m,k} | \mathbf{x}_t^{m,k}), & \text{otherwise.} \end{cases} \quad (9)$$

333 Particles with zero probability do not survive during resam-  
 334 pling. If the tracked object leaves the field of view all  $M$  334

335 particles of an object  $k$  are removed, i.e.,  $p(\mathbf{y}_t^{m,k} | \mathbf{x}_t^{m,k}) = 0$   
 336 for all particles of object  $k$ .

### 337 F. Incremental Learning of Object Models

338 To handle the appearance change of the tracked objects over  
 339 time, we use online learning to adapt the reference histograms  
 340  $h_{c,\text{ref}}^k$ ,  $c \in \{H, S, V\}$  (similar to [6]) and ellipse size  $A_{\text{ref}}^k$  and  
 341  $B_{\text{ref}}^k$ . Therefore, a learning rate  $\alpha$  is introduced, and the model  
 342 parameters for target object  $k$  are updated according to

$$h_{c,\text{ref}}^k = \alpha \hat{h}_c^k + (1 - \alpha) h_{c,\text{ref}}^k, \quad c \in \{H, S, V\} \quad (10)$$

$$A_{\text{ref}}^k = \alpha \hat{A}^k + (1 - \alpha) A_{\text{ref}}^k \quad (11)$$

$$B_{\text{ref}}^k = \alpha \hat{B}^k + (1 - \alpha) B_{\text{ref}}^k \quad (12)$$

343 where  $\hat{h}_c^k$  denotes the histogram and  $\hat{A}^k$  and  $\hat{B}^k$  are the prin-  
 344 cipal axes of the bounding ellipse of the nonoccluded (i.e., no  
 345 mutual occlusion between tracked objects) skin-colored region  
 346 of the corresponding tracked object  $k$  located at  $\{\mathbf{x}_t^{m,k}\}_{m=1}^M$ .  
 347 Again, this region has to be larger than  $T_{\text{Area}}$ . No update of  
 348 the reference models is performed in the case where occlusion  
 349 between the tracked objects occurs or the skin-colored region  
 350 is not large enough. The latter condition is a simple way to  
 351 ensure that the model update is only conducted for faces.  
 352 This simplistic assumption can be appropriately extended by  
 353 integrating more advanced face models.

354 The learning rate  $\alpha$  introduces an *exponential forgetting*  
 355 *process*, i.e., the contribution of a specific object exponentially  
 356 decreases as it recedes into the past. Currently, the learning rate  
 357 (value between 0 and 1) is fixed (a good value has been selected  
 358 during experiments). However,  $\alpha$  could be adapted depending  
 359 on the dynamics of the scene.

#### 360 **Algorithm 1** Particle Filter Tracking

361 **Input:**  $I_{x_0:T, y_0:T}^{\text{HSV}}$  (Color image sequence  $0 : T$ ),

362 Skin-color model  $\Theta$

363 **Parameters:**  $M, N, \lambda, C, D, T_w, T_{\text{motion}}, T_{\text{skin}}, T_{\text{Area}},$

364  $T_{\text{Kill}}, \alpha$   
 365 **Output:**  $\{\{\mathbf{x}_{0:T}^{m,k}\}_{m=1}^M\}_{\forall k}$

366  $t \leftarrow 0$

367  $k \leftarrow 0$

368 **while** InitObjects **do**

369  $k \leftarrow k + 1$

370 Obtain:  $h_{c,\text{ref}}^k : c \in \{H, S, V\}, A_{\text{ref}}^k, B_{\text{ref}}^k, \mathbf{x}_{\text{ref}}^k$

371  $\mathbf{x}_{t+1}^{m,k} \leftarrow \mathbf{x}_{\text{ref}}^k + C\mathbf{v}_t \quad \forall m = 1, \dots, M$  (Generate particles)

372 **end while**

373  $K \leftarrow k$

374 **for**  $t = 1$  **to**  $T$  **do**

375  $w_t^{m,k} \propto p(\mathbf{y}_t^{m,k} | \mathbf{x}_t^{m,k})$

376  $\forall k = 1, \dots, K \quad \forall m = 1, \dots, M$

377 **while** KillObjects **do**

378  $k \leftarrow$  Determine object to terminate

379 Remove  $M$  particles  $x_t^{m,k}$  of object  $k$

380 Remove reference histogram and ellipse size:

381  $h_{c,\text{ref}}^k : c \in \{H, S, V\}, A_{\text{ref}}^k, B_{\text{ref}}^k$

382  $K \leftarrow K - 1$

383 **end while**

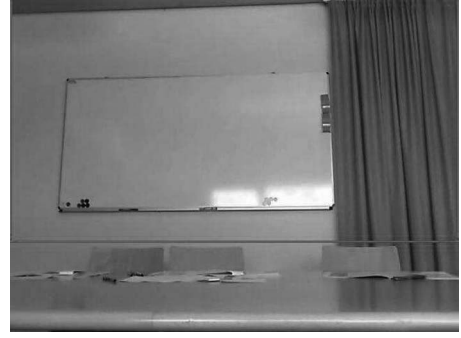


Fig. 3. Tracking scene. We track and initialize objects in the red rectangle.

**for**  $k = 1$  **to**  $K$  **do** 384

$w_t^{m,k} \leftarrow w_t^{m,k} / \sum_{m'=1}^M w_t^{m',k} \quad \forall m = 1, \dots, M$  385

$\{\mathbf{x}_t^{m,k}\}_{m=1}^M \leftarrow$  Resampling 386

(with replacement):  $\{\mathbf{x}_t^{m,k}, w_t^{m,k}\}_{m=1}^M$  387

$\mathbf{x}_{t+1}^{m,k} \leftarrow \mathbf{x}_t^{m,k} + C\mathbf{v}_t + (D/2M) \sum_{m'=1}^M (\mathbf{x}_t^{m',k} - \mathbf{x}_{t-1}^{m',k})$  388

$\forall m = 1, \dots, M$  (Apply state-space dynamics) 389

**if** OnlineUpdate **then** 390

Determine:  $\hat{h}_c^k : c \in \{H, S, V\}, \hat{A}^k, \hat{B}^k$  391

$h_{c,\text{ref}}^k \leftarrow \alpha \hat{h}_c^k + (1 - \alpha) h_{c,\text{ref}}^k \quad c \in \{H, S, V\}$  392

$A_{\text{ref}}^k \leftarrow \alpha \hat{A}^k + (1 - \alpha) A_{\text{ref}}^k$  393

$B_{\text{ref}}^k \leftarrow \alpha \hat{B}^k + (1 - \alpha) B_{\text{ref}}^k$  394

**end if** 395

**end for** 396

**while** InitObjects **do** 397

$K \leftarrow K + 1$  398

Obtain:  $h_{c,\text{ref}}^K : c \in \{H, S, V\}, A_{\text{ref}}^K, B_{\text{ref}}^K, \mathbf{x}_{\text{ref}}^K$  399

$\mathbf{x}_{t+1}^{m,K} \leftarrow \mathbf{x}_{\text{ref}}^K + C\mathbf{v}_t \quad \forall m = 1, \dots, M$  (Generate 400

particles) 401

**end while** 402

**end for** 403

### G. Implemented Tracker 404

In the following, we sketch our tracking approach for multi- 405  
 ple objects (see Algorithm 1). The binary variable *InitObject* 406  
 denotes that a new object for tracking has been detected. 407  
*KillObject* is set if an object should be terminated. *OnlineUp-* 408  
*date* indicates that object  $k$  located at  $\{\mathbf{x}_t^{m,k}\}_{m=1}^M$  is nonoc- 409  
 cluded, and the area of the skin-colored region is larger than 410  
 $T_{\text{Area}}$ , i.e., we perform online learning for reference model  $k$ . 411

Our implementation is related to the *dual estimation* problem 412  
 [13], where both the states of multiple objects  $\mathbf{x}_t^{m,k}$  and the 413  
 parameters of the reference models are simultaneously esti- 414  
 mated given the observations. At every time step, the particle 415  
 filter estimates the states using the observation likelihood of 416  
 the current reference models, while the online learning of the 417  
 reference models is based on the current state estimates. 418

## III. RELATIONSHIP TO GAS 419

GAs are optimization algorithms founded upon the principles 420  
 of natural evolution discovered by Darwin. In nature, individ- 421  
 uals have to adapt to their environment in order to survive in 422





Fig. 4. Tracking of people. Frames: 1, 416, 430, 449, 463, 491, 583, 609, 622, 637, 774, 844, 967, 975, 1182, 1400 (the frame number is assigned from left to right and top to bottom).

423 a process of further development. An introduction of GAs can  
 424 be found in [25] and [26]. GA are stochastic procedures which  
 425 have been successfully applied in many optimization tasks.  
 426 GA operate on a population of potential solutions applying the  
 427 principle of *survival of the fittest individual* to produce better  
 428 and better approximations to the solution. At each generation, a  
 429 new set of approximations is created by the process of selecting  
 430 individuals according to their level of fitness in the problem  
 431 domain and assembling them together using operators inspired  
 432 from nature. This leads to the evolution of individuals that are  
 433 better suited to their environment than the parent individuals  
 434 they were created from. GA model the natural processes, such  
 435 as selection, recombination, and mutation. Starting from an  
 436 initial population  $P(0)$ , the sequence  $P(0), P(1), \dots, P(t)$ ,  
 437  $P(t + 1)$  is called population sequence or evolution. The end of  
 438 an artificial evolution process is reached once the termination  
 439 condition is met, and the result of the optimization task is  
 440 available.

441 In this section, we want to point to the close relationship  
 442 between GA and our particle filter for tracking. This analogy  
 443 has been mentioned in [27]. As suggested in Section II, we  
 444 treat the tracking of multiple objects completely independent,  
 445 i.e., we have a set of  $M$  particles for each object  $k$ . In the GA  
 446 framework, we can relate this to  $k$  instantiations of GA, one  
 447 for each tracked object. Hence, each particle  $\mathbf{x}_t^m$  of object  $k$   
 448 represents one individual in the population  $P(t)$  which is value  
 449 encoded. The population size is  $M$ . A new genetic evolution

process is started once a new object is initialized for tracking 450  
 (InitObject). The evolution process of the GA is terminated 451  
 either at the end of the video ( $t = T$ ) or when the set of 452  
 individuals is not supported by the fitness value (KillObject). 453  
 The observation likelihood  $p(\mathbf{y}_t^{m,k} | \mathbf{x}_t^{m,k})$  denotes the fitness 454  
 function to evaluate the individuals. However, the scope of GA 455  
 for tracking is slightly different. GA are generally used to find a 456  
 set of parameters for a given optimization task, i.e., the aim is to 457  
 find the individual with the best fitness after the termination of 458  
 the GA. Whereas, in the tracking case, the focus lies on the evo- 459  
 lution of the individuals, i.e., the trajectory of the tracked object. 460

The selection operator directs the search toward promising 461  
 regions in the search space. *Roulette Wheel Selection* [28] is a 462  
 widely used selection method which is very similar to sampling 463  
 with replacement as used in Section II. To each individual, a re- 464  
 production probability according to  $w_t^m \leftarrow w_t^m / \sum_{m'=1}^M w_t^{m'}$  465  
 is assigned. A roulette wheel is constructed with a slot size cor- 466  
 responding to the individuals reproduction probability. Then, 467  
 $M$  uniformly distributed random numbers on the interval  $[0, 1]$  468  
 are drawn and distributed according to their value around the 469  
 wheel. The slots where they are placed to compose the subse- 470  
 quent population  $P(t)$ . The state-space dynamics of the particle 471  
 filter (see Section II-B) is modeled by the recombination and 472  
 mutation operator. 473

The framework of the GA for tracking one object  $k$  is 474  
 presented in Algorithm 1. The incremental learning of the 475  
 reference model is omitted for the sake of brevity. 476



Fig. 5. Partial occlusions. Frames: 468, 616, 974, 4363 (the frame number is assigned from left to right and top to bottom).

477 **Algorithm 2** GA Tracking  
 478 **Input:**  $I_{x_t:T, y_t:T}^{\text{HSV}}$  (Color image sequence  $t : T$ ),  
 479 **Parameters:**  $M, N, \lambda, C, D, T_{\text{Kill}}$   
 480 **Output:**  $\{\mathbf{x}_{t:T}^m\}_{m=1}^M$  (set of particle sequences  $t : T$ )  
 481 Initialize population  $P(t)$  :  
 482  $\mathbf{x}_t^m \leftarrow \mathbf{x}_{\text{ref}} + C\mathbf{v}_t \quad \forall m = 1, \dots, M$   
 483 **while**  $\text{KillObject} \cap t < T$  (Loop over image sequence) **do**  
 484 Evaluate individuals:  
 485  $w_t^m \leftarrow p(\mathbf{y}_t^m | \mathbf{x}_t^m) \quad \forall m = 1, \dots, M$   
 486 Selection  $P(t)$ :  
 487  $\{\mathbf{x}_t^m\}_{m=1}^M \leftarrow$  (Sampling with replacement)  $\{\mathbf{x}_t^m, w_t^m\}_{m=1}^M$   
 488 Recombination  $P(t+1)$  :  
 489  $\mathbf{x}_{t+1}^m \leftarrow \mathbf{x}_t^m + (D/2M) \sum_{m'=1}^M (\mathbf{x}_t^{m'} - \mathbf{x}_{t-1}^{m'})$   
 490  $\forall m = 1, \dots, M$   
 491 Mutation  $P(t+1)$  :  $\mathbf{x}_{t+1}^m \leftarrow \mathbf{x}_{t+1}^m + C\mathbf{v}_t \quad \forall m = 1, \dots, M$   
 492  $t \leftarrow t + 1$   
 493 **end while**

#### 494 IV. EXPERIMENTS

495 We present tracking results on meeting data in Section IV-A  
 496 where we do both tracking of multiple persons and on-  
 497 line adaptation of the reference models during tracking. In  
 498 Section IV-B, we empirically show that the adaptation of the  
 499 reference model during tracking (single object) of an indoor  
 500 and outdoor scene results in a more robust tracking. Finally, in  
 501 Section IV-C, tracking results using reference model adaptation  
 502 for multiple objects of an outdoor scene are presented. For the  
 503 outdoor scenes, we report the average standard deviation of  
 504 the trajectories of independent tracking runs depending on the  
 505 learning rate  $\alpha$ .

##### 506 A. Meeting Scenario

507 The meeting room layout is shown in Fig. 3. The red rec-  
 508 tangle [region of interest (ROI)] in the image marks the frame  
 509 where tracking and initialization of objects are performed. Peo-  
 510 ple may enter and leave on both sides of the image. Currently,  
 511 our tracker initializes a new target even if it enters from the



(a)



(b)



(c)

Fig. 6. Outdoor tracking. Frames: 7, 11, 12, 13, 14, 20, 42, 63, 80, 107, 136, 146, 158, 165, 192 (the frame number is assigned from left to right and top to bottom). (a) Original image sequence. (b) Tracking without reference model adaptation ( $\alpha = 0$ ). (c) Tracking with online reference model learning ( $\alpha = 0.2$ ).

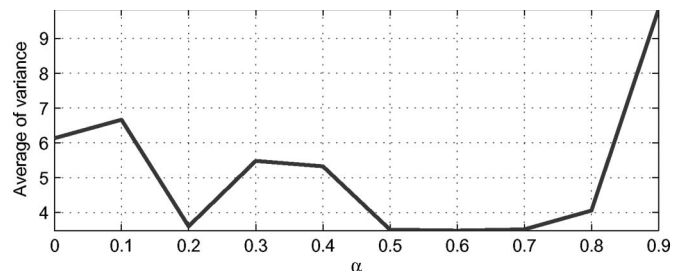


Fig. 7. Averaged standard deviation of the trajectories of 100 tracking runs depending on the reference model learning rate  $\alpha$ .

bottom, e.g., a hand moving from the table into the ROI. The 512  
 strong reflections at the table, chairs, and the white board cause 513  
 noise in the motion image. 514

For testing the performance of our tracking approach, ten 515  
 videos with  $\sim 7000$  frames have been used. The resolution is 516  
 $640 \times 480$  pixels. The meeting room is equipped with a table 517  
 and three chairs. We have different persons in each video. The 518  
 people are coming from both sides into the frame moving 519

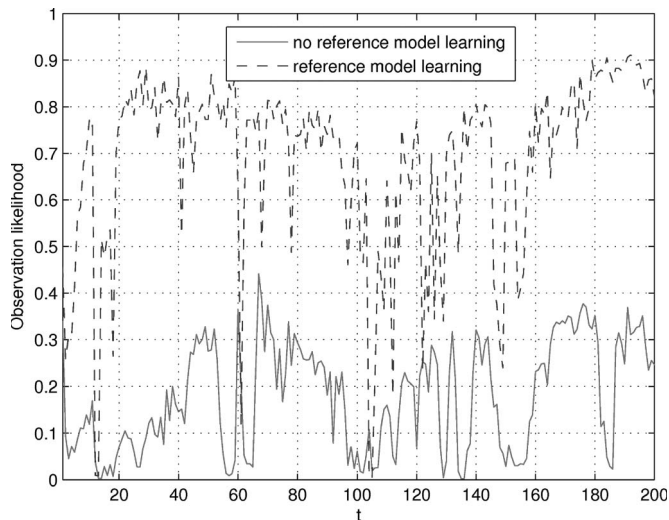


Fig. 8. Observation likelihood of outdoor sequence.

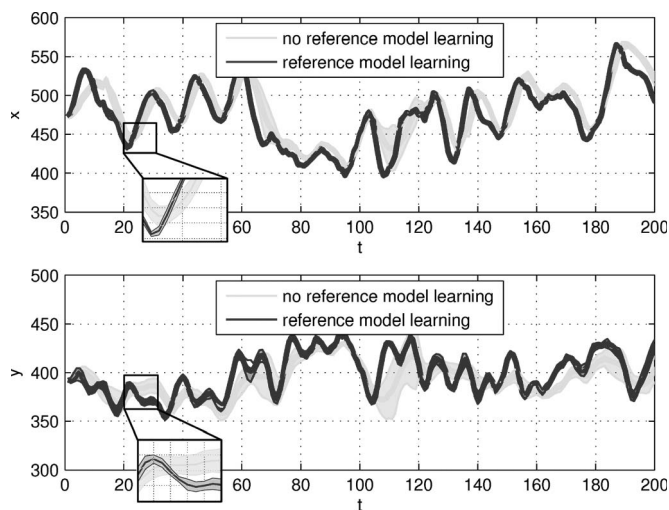


Fig. 9. Averaged trajectory with standard deviation in  $x$  and  $y$  of outdoor sequence (over ten runs).

520 to chairs and sit down. After a short discussion, people are  
521 sequentially leaving the room, are coming back, sit down at  
522 different chairs and so on. At the beginning, people may already  
523 sit at the chairs. In this case, we have to automatically initialize  
524 multiple objects at the very first frame.

525 Fig. 4 shows the result of the implemented tracker for one  
526 video. All the initializations and terminations of objects are  
527 performed automatically. The appearance of an object changes  
528 over time. When entering the frame, we get the side view of  
529 the person's head. After sitting down at the table, we have a  
530 frontal view. We account for this by incrementally updating the  
531 reference histogram during tracking. We perform this only in  
532 the case where no mutual occlusions with other tracked objects  
533 are existent. The participants were successfully tracked over  
534 long image sequences.

535 First, the person on the left side stands up and leaves the room  
536 on the right side (frame 416–491). When walking behind the  
537 two sitting people, partial occlusions occur which do not cause  
538 problems. Next, the person on the right (frame 583–637) leaves  
539 the room on the left side. His face is again partially occluded



(a)



(b)



(c)

Fig. 10. Indoor tracking. Frames: 1, 12, 24, 31, 38, 41, 47, 54, 65, 71, 80, 107, 113, 120, 134 (the frame number is assigned from left to right and top to bottom). (a) Original image sequence. (b) Tracking without reference model adaptation ( $\alpha = 0$ ). (c) Tracking with online reference model learning ( $\alpha = 0.2$ ).

by the person in the middle. Then, the person on the center  
540 chair leaves the room (frame 774). After that, a person on the  
541 right side enters and sits at the left chair (frame 844). At frame  
542 967, a small person is entering and moving to the chair in the  
543 middle. Here, again, a partial occlusion occurs at frame 975, 544  
544 which is also tackled. Finally, a person enters from the right  
545 and sits down on the right chair (frame 1182, 1400). The partial  
546 occlusions are shown in Fig. 5. Also, the blurred face of the  
547 moving person in the back can be observed in this figure. The  
548 reference model adaptation enables a more robust tracking. If  
549 we do not update the models of the tracked objects over time,  
550 the tracking fails in case of these partial occlusions. In [29],  
551 occlusions are handled using multiple cameras for tracking  
552 participants in a meeting. 553

## B. Reference Model Adaptation for Single-Object Tracking 554

In the following, we show the benefit of the reference model  
555 adaptation during tracking of a short indoor and outdoor se-  
556 quence. In contrast to the meeting scenario, we restrict the  
557

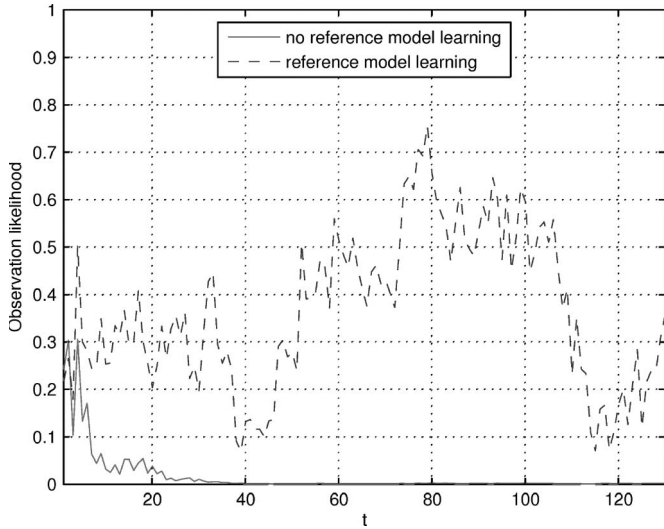


Fig. 11. Observation likelihood of indoor sequence.

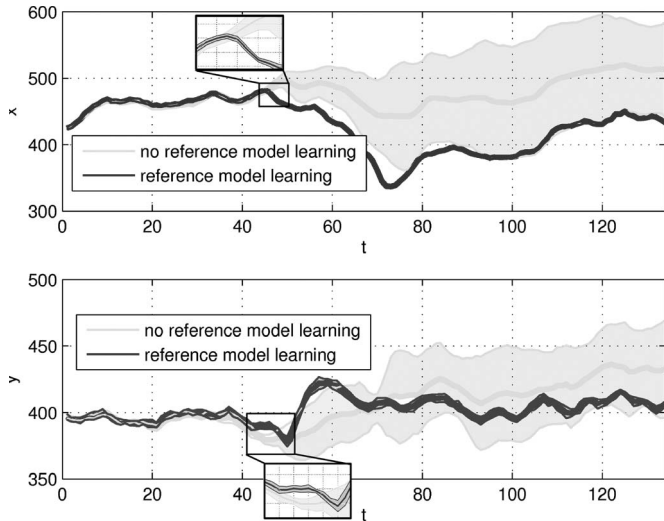
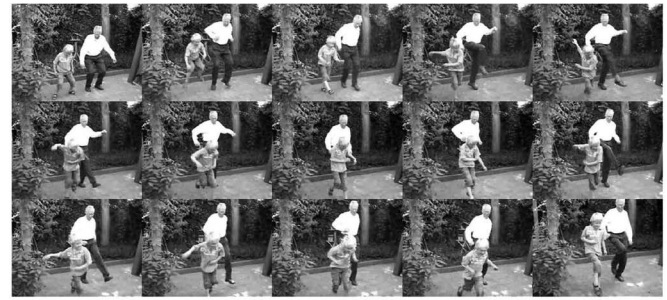


Fig. 12. Averaged trajectory with standard deviation in  $x$  and  $y$  of indoor sequence (over ten runs).

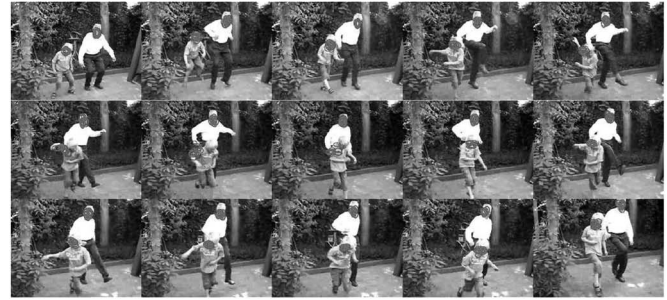
558 tracking to one single object, i.e., face. This means, in par-  
 559 ticular, that the automatic initialization and termination of the  
 560 object is disabled. The object is initialized by hand in the first  
 561 frame.

562 Fig. 6(a) shows a short outdoor sequence where a person is  
 563 moving behind a tree and two cars with strongly changing light-  
 564 ing conditions. We have a total occlusion of the face in frames  
 565 12 and 13 and a partial occluded face in frames 146 to 165. We  
 566 repeated the tracking without and with reference model learning  
 567 ten times, and a typical result is shown in Fig. 6(b) and (c),  
 568 respectively. We use  $M = 50$  particles for tracking, whereas  
 569 only 15 particles with the best observation likelihood are shown  
 570 in the figures.

571 In Fig. 7, we present the average standard deviation of  
 572 the trajectories over 100 tracking runs. The reference model  
 573 learning rate  $\alpha$  has been chosen in the range of  $0, \dots, 0.6$   
 574 (0 means that there is no learning). The optimal learning rate  
 575 with respect to a low standard deviation of the trajectories over  
 576 100 independent runs is  $\alpha = 0.2$  for this outdoor sequence.



(a)



(b)



(c)

Fig. 13. Outdoor tracking of multiple objects. Frames: 1, 12, 30, 47, 49, 51, 53, 57, 59, 79, 105, 107, 109, 111, 149 (the frame number is assigned from left to right and top to bottom). (a) Original image sequence. (b) Tracking without reference model adaptation ( $\alpha = 0$ ). (c) Tracking with online reference model learning ( $\alpha = 0.1$ ).

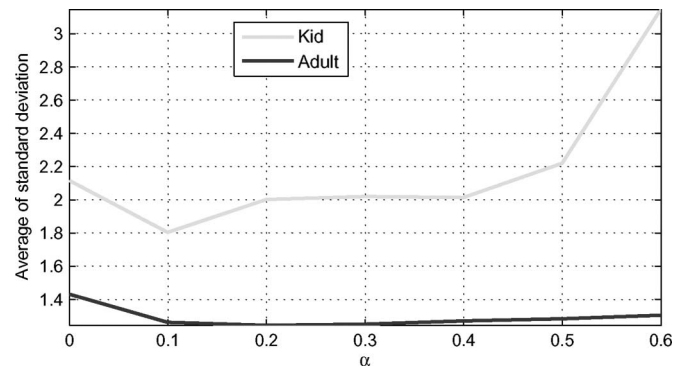


Fig. 14. Averaged standard deviation of the trajectories of ten tracking runs depending on the reference model learning rate  $\alpha$ .

Fig. 8 shows the observation likelihood of the best particle 577 during tracking. At the complete occlusion (frames  $t = 12$  and 578  $t = 13$ ) and the partial occlusion (frames  $t = 145, \dots, 160$ ), the 579 observation likelihood drops, however, with reference model 580 learning a quick recovery is supported. 581

Fig. 9 summarizes the averaged trajectory with the standard deviation over ten different tracking runs performed for the outdoor scene. In the case of reference model learning, we observe in the video sequences that the tracking of the face gives highly similar trajectories. The standard deviation is small and approximately constant over time. However, if no learning of the reference model is performed, the standard deviation is large in certain time segments. This leads to the conclusion that model adaptation results in a more robust tracking.

Fig. 10(a) shows an indoor video where a person is moving on a corridor, and a tree causes partial occlusion of the tracked face. Additionally, the lighting conditions are strongly varying. The face is partially occluded by the tree in frames 37–50 and 110–126. Again, the tracking without and with reference model learning is repeated ten times, and a typical result is shown in Fig. 10(b) and (c), respectively. Only 15 particles with the best observation likelihood are visualized. The parameter setting is the same as in the previous experiments. The tracker without reference model refinement fails during the first occlusion in all ten runs, whereas the tracker with online model update is successful in all cases. The optimal learning rate  $\alpha$  is set to 0.2 (established during experiments).

This can be also observed in the observation likelihood of the best particle over time (see Fig. 11) and in the averaged trajectory over ten tracking results (see Fig. 12).

### C. Reference Model Adaptation for Multiple Object Tracking

We show tracking results for an outdoor scene where a kid is showing an adult dancing steps (see Fig. 13). A typical tracking result without and with reference model learning is shown in Fig. 13(b) and (c), respectively. Again,  $M = 50$  particles are used, whereas only 15 particles with the best observation likelihood are shown in the figures. Similar as in the previous section, we did a repeatability test, i.e., we tracked the objects over ten independent runs. The tracked objects are initialized by hand in the very first frame.

Fig. 14 shows the average standard deviation of the trajectories of ten tracking runs using a learning rate  $\alpha$  in the range of 0, . . . , 0.6. The optimal learning rate for the *Kid* and the *Adult* is  $\alpha = 0.1$  and  $\alpha = 0.2$ , respectively. Currently,  $\alpha$  is fixed for the whole image sequence. Ideally,  $\alpha$  could be adapted depending on the dynamics of the scene.

623

## V. CONCLUSION

We propose a robust visual tracking algorithm for multiple objects (faces of people) in a meeting scenario based on low-level features as skin color, target motion, and target size. Based on these features, automatic initialization and termination of objects is performed. For tracking a sampling importance resampling, particle filter has been used to propagate sample distributions over time. Furthermore, we use online learning of the target models to handle the appearance variability of the objects. We discuss the similarity between our implemented tracker and GAs. Each particle represents an individual in the GA framework. The evaluation function incorporates the observation likelihood model and the individual selection

process maps to the resampling procedure in the particle filter. The state-space dynamics is incorporated in the recombination and mutation operator of the GA. Numerous experiments on meeting data show the capabilities of the tracking approach. The participants were successfully tracked over long image sequences. Partial occlusions are handled by the algorithm. Additionally, we empirically show that the adaptation of the reference model during tracking of indoor and outdoor scenes results in a more robust tracking.

Future work concentrates on extending the tracker to other scenarios and to investigate an adaptive reference model learning rate  $\alpha$  which depends on the dynamics of the scene. Furthermore, we aim to develop approaches for tackling occlusions.

## ACKNOWLEDGMENT

This work was supported by the Austrian Science Fund (Project S106). The author would like to thank M. Grabner and M. Kepesi who collected the data during their involvement in the MISTRAL Project ([www.mistral-project.at](http://www.mistral-project.at)). The MISTRAL Project was funded by the Austrian Research Promotion Agency ([www.ffg.at](http://www.ffg.at)) within the strategic objective FIT-IT under Project 809264/9338. The author would also like to thank C. Kirchstätter for recording the indoor and outdoor video.

## REFERENCES

- [1] J. Lim, D. Ross, R.-S. Lin, and M.-H. Yang, "Incremental learning for visual tracking," in *Advances in Neural Information Processing Systems 17*. Cambridge, MA: MIT Press, 2005.
- [2] M. Isard and A. Blake, "Condensation—Conditional density propagation for visual tracking," *Int. J. Comput. Vis.*, vol. 29, no. 1, pp. 5–28, 1998.
- [3] S. Arulampalam, S. Maskell, N. Gordon, and T. Clapp, "A tutorial on particle filters for on-line non-linear/non-Gaussian Bayesian tracking," *IEEE Trans. Signal Process.*, vol. 50, no. 2, pp. 174–188, Feb. 2002.
- [4] A. D. Jepson, D. J. Fleet, and T. F. El-Maraghi, "Robust online appearance models for visual tracking," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 25, no. 10, pp. 1296–1311, Oct. 2003.
- [5] S. McKenna, Y. Raja, and S. Gong, "Tracking colour objects using adaptive mixture models," *Image Vis. Comput.*, vol. 17, no. 3/4, pp. 225–231, 1999.
- [6] K. Nummiaro, E. Koller-Meier, and L. Van Gool, "An adaptive color-based particle filter," *Image Vis. Comput.*, vol. 21, no. 1, pp. 99–110, Jan. 2003.
- [7] P. Pérez, C. Hue, J. Vermaak, and M. Gangnet, "Color-based probabilistic tracking," in *Proc. ECCV*, 2002, pp. 661–675.
- [8] S. L. Dockstader and A. Tekalp, "Tracking multiple objects in the presence of articulated and occluded motion," in *Proc. Workshop Human Motion*, 2000, pp. 88–98.
- [9] C. Hue, J.-P. Le Cadre, and P. Pérez, "Tracking multiple objects with particle filtering," *IEEE Trans. Aerosp. Electron. Syst.*, vol. 38, no. 3, pp. 791–812, Jul. 2002.
- [10] J. Vermaak, A. Doucet, and P. Pérez, "Maintaining multi-modality through mixture tracking," in *Proc. ICCV*, 2003, pp. 1110–1116.
- [11] K. Okuma, A. Taleghani, N. de Freitas, J. J. Little, and D. G. Lowe, "A boosted particle filter: Multitarget detection and tracking," in *Proc. ECCV*, 2004, pp. 28–39.
- [12] Y. Cai, N. de Freitas, and J. J. Little, "Robust visual tracking for multiple targets," in *Proc. ECCV*, 2006, pp. 107–118.
- [13] S. Haykin, *Kalman Filtering and Neural Networks*. Hoboken, NJ: Wiley, 2001.
- [14] J. Lin, "Divergence measures based on the Shannon entropy," *IEEE Trans. Inf. Theory*, vol. 37, no. 1, pp. 145–151, Jan. 1991.
- [15] T. Cover and J. Thomas, *Elements of Information Theory*. Hoboken, NJ: Wiley, 1991.
- [16] D. Comaniciu, V. Ramesh, and P. Meer, "Real-time tracking of non-rigid objects using mean shift," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2005, pp. 142–149.

- 700 [17] R. E. Kalman, "A new approach to linear filtering and prediction prob-  
701 lems," *Trans. ASME, Ser. D, J. Basic Eng.*, vol. 82, pp. 34–45, 1960.
- 702 [18] C. M. Bishop, *Pattern Recognition and Machine Learning*. Berlin,  
703 Germany: Springer Sci.+Bus. Media, 2006.
- 704 [19] L. R. Rabiner, "A tutorial on hidden Markov models and selected appli-  
705 cations in speech recognition," *Proc. IEEE*, vol. 77, no. 2, pp. 257–286,  
706 Feb. 1989.
- 707 [20] Z. Ghahramani and G. E. Hinton, "Variational learning for switching  
708 state-space models," *Neural Comput.*, vol. 12, no. 4, pp. 963–996, 1998.
- 709 [21] A. Doucet, "On sequential Monte Carlo sampling methods for Bayesian  
710 filtering," Dept. Eng., Cambridge Univ., London, U.K., Tech. Rep.  
711 CUED/F-INFENG/TR. 310, 1998.
- 712 [22] M. Sonka, V. Hlavac, and R. Boyle, *Image Processing, Analysis, and  
713 Machine Vision*. London, U.K.: Int. Thomson, 1999.
- 714 [23] R. O. Duda, P. E. Hart, and D. G. Stork, *Pattern Classification*. Hoboken,  
715 NJ: Wiley, 2000.
- 716 [24] A. Dempster, N. Laird, and D. Rubin, "Maximum likelihood estimation  
717 from incomplete data via the EM algorithm," *J. R. Stat. Soc., Ser. B, Stat.  
718 Methodol.*, vol. 39, pp. 1–38, 1977.
- 719 [25] D. E. Goldberg, *Genetic Algorithms in Search, Optimization and Machine  
720 Learning*. Reading, MA: Addison-Wesley, 1989.
- 721 [26] T. Bäck, *Evolutionary Algorithms in Theory and Practice*. London,  
722 U.K.: Oxford Univ. Press, 1996.
- [27] T. Higuchi, "Monte Carlo filter using the genetic algorithm operators," 723  
*J. Stat. Comput. Simul.*, vol. 59, no. 1, pp. 1–23, Aug. 1997. 724
- [28] L. E. Baker, "Reducing bias and inefficiency in the selection algorithm," 725  
in *Proc. Int. Conf. Genetic Algorithms Appl.*, 1987, pp. 14–21. 726
- [29] H. Nait Charif and S. J. McKenna, "Tracking the activity of participants 727  
in a meeting," *Mach. Vis. Appl.*, vol. 17, no. 2, pp. 83–93, 2006. 728



**Franz Pernkopf** received the M.Sc. (Dipl.Ing.) de- 729  
gree in electrical engineering from the Graz Univer- 730  
sity of Technology, Graz, Austria, in 1999, and the 731  
Ph.D. degree from the University of Leoben, Leoben, 732  
Austria, in 2002. 733

He was a Research Associate with the Department 734  
of Electrical Engineering, University of Washington, 735  
Seattle, from 2004 to 2006. Currently, he is an 736  
Assistant Professor with the Signal Processing and 737  
Speech Communication Laboratory, Graz University 738  
of Technology. His research interests include ma- 739  
chine learning, Bayesian networks, feature selection, finite mixture models, 740  
vision, speech, and statistical pattern recognition. 741

Dr. Pernkopf was awarded the Erwin Schrödinger Fellowship in 2002. 742