

# Face Tracking with Automatic Model Construction

Jesus Nuevo, Luis M. Bergasa\*, David F. Llorca, Manuel Ocaña

*Department of Electronics, Universidad de Alcala. Esc. Politecnica, Crta Madrid-Barcelona, Km 33,600.  
28871 Alcala de Henares, Madrid*

---

## Abstract

~~Driver inattention is one of the major causes of traffic crashes, claiming thousands of lives every year. Face tracking is one of the first stages in safety systems that relay on computer vision to detect inattention.~~ This paper describes an active model with a robust texture model built on-line. The model uses one camera and it is able to operate without active illumination. The texture model is defined by a series of clusters, which are built in a video sequence using previously encountered samples. This model is used to search for the corresponding element in the following frames. An on-line clustering method, named *leaderP* is described and evaluated on an application of face tracking. A 20-point shape model is used. This model is built offline, and a robust fitting function is used to restrict the position of the points. *Our proposal is to serve as one of the stages in a driver monitoring system.* To test it, a new set of sequences of drivers recorded outdoors and in a realistic simulator has been compiled. Experimental results for typical outdoor driving scenarios, with frequent head movement, turns and occlusions are presented. Our approach is tested and compared with the Simultaneous Modeling and Tracking (SMAT) [1], and the recently presented Stacked Trimmed Active Shape Model (STASM) [2], and shows better results than SMAT and similar fitting error levels to STASM, with much faster execution times and improved robustness.

*Keywords:* Face tracking, appearance modeling, incremental clustering, robust fitting, driver monitoring

---

## 1. Introduction

Driver inattention is a major cause of traffic accidents, and it has been found to be involved in some form in 80 percent of the crashes and 65 percent of the near crashes within 3 seconds of the event [3]. Monitoring a driver to detect inattention is a complex problem that involves physiological and behavioural elements. Different works have been presented in recent years, focused mainly in drowsiness, with a broad range of techniques. Physiological

---

\*Corresponding author: Tel. (+34) 91885 6569

*Email addresses:* [jnuevo@depeca.uah.es](mailto:jnuevo@depeca.uah.es) (Jesus Nuevo), [bergasa@depeca.uah.es](mailto:bergasa@depeca.uah.es) (Luis M. Bergasa), [llorca@depeca.uah.es](mailto:llorca@depeca.uah.es) (David F. Llorca), [mocana@depeca.uah.es](mailto:mocana@depeca.uah.es) (Manuel Ocaña)

7 measurements such as electro-encephalography (EEG) [4] or electro-oculography (EOG),  
8 provide the best data for detection [4]. The problem with these techniques is that they are  
9 intrusive to the subject. Moreover, medical equipment is always expensive.

10 Lateral position of the vehicle inside the lane, steering wheel movements and time-to-line  
11 crossing are commonly used, and some commercial systems have been developed [5, 6]. These  
12 techniques are not invasive, and to date they obtain the most reliable results. However, the  
13 measurements they use may not reflect behaviors such as the so-called micro-sleeps [7]. They  
14 also require a training period for each person, and thus are not applicable to the occasional  
15 driver.

16 Drivers in fatigue exhibit changes in the way their eyes perform some actions, like moving  
17 or blinking. These actions are known as *visual behaviors*, and are readily observable in drowsy  
18 and distracted drivers. Face pose [8] and gaze direction also contain information and have  
19 been used as another element of inattention detection systems [9]. Computer vision has  
20 been the tool of choice for many researchers to be used to monitor visual behaviours, as it  
21 is non-intrusive. Most systems use one or two cameras to track the head and eyes of the  
22 subject [10, 11, 12, 13, 14]. A few companies commercialize systems [15, 16] as accessories  
23 for installation in vehicles. These systems require user-specific calibration, and some of them  
24 use near-IR lighting, which is known to produce eye fatigue. Reliability of these systems is  
25 still not high enough for car companies to take on the responsibility of its production and  
26 possible liability in case of malfunctioning.

27 Face location and tracking are the first processing stages of most computer vision systems  
28 for driver monitoring. Some of the most successful systems to date use near-IR active  
29 illumination [17, 18, 19], to simplify the detection of the eyes thanks to the *bright pupil*  
30 effect. Near-IR illumination is not as useful during the day because sunlight also has a  
31 near-IR component. As mentioned above, near-IR can produce eye fatigue and thus limits  
32 the amount of time these systems can be used on a person.

33 *Given the complexity of the problem, it has been divided in parts and in this work only*  
34 *the problem of face tracking is addressed.*

35 *This paper presents a new active model with the texture model built incrementally. We*  
36 *use it to characterize and track the face in video sequences.* The tracker can operate without  
37 active illumination. The texture model of the face is created online, and thus specific for each  
38 person without requiring a training phase. A new online clustering algorithm is described,  
39 and its performance compared with the method proposed in [1]. Two shape models, trained  
40 online and off-line, are compared. This paper also presents a new video sequence database,  
41 recorded in a car moving outdoors and in a simulator. The database is used to assess  
42 the performance of the proposed face tracking method in the challenging environment a  
43 driver monitoring application would meet. *No evaluations of face pose estimation and driver*  
44 *inattention detection are performed.*

45 The rest of the paper is structured as follows. Section 2 presents a few remarkable  
46 works in face tracking in the literature that are related to our proposal. Section 3 describes  
47 our approach. Section 4 describes the video dataset used for performance evaluation, and  
48 experimental results. This paper closes with conclusions and future work.

## 2. Background

Human face tracking is a broad field in computing research [20], and a myriad of techniques have been developed in the last decades. It is of the greatest interest, as vast amounts of information are contained in face features, movements and gestures, which are constantly used for human communication. Systems that work on such data often use face tracking [21, 22].

Non-rigid object tracking has been a major focus of research in latter years, and general purpose template-based trackers have been used to track faces in the literature with success. Several efficient approaches have been presented [23, 24, 25, 26].

Statistical models have been used for face modeling and tracking. Active Shape Models [27] (ASM) are similar to the active contours (*snakes*), but include constraints from a Point Distribution Model (PDM) [28] computed in advance from a training set. Advances in late years have increased their robustness and precision to remarkable levels (STASM,[2]). Extensions of ASM that include modeling of texture have been presented, of which Active Appearance Models (AAMs) [29] are arguably the best known. Active Appearance Models are global models in the sense that the minimization is performed over all pixels that fall inside the mesh defined by the mean of the PDM. All these models have an offline training phase, which require comprehensive training sets so they can generalize properly to unseen instances of the object. This is time consuming process, and there is still the risk that perfectly valid instances of the object would not be modeled correctly.

Several methods that work without *a priori* models have been presented in the literature. Most of them focus on patch tracking on a video sequence. The classic approach is to use the image patch extracted on the first frame of the sequence to search for similar patches on the following frames. Lukas-Kanade method [30] was one of the first proposed solutions and it is still widely used. Jepson *et al.* [31] presented a system with appearance model based on three components: a stable component that is learned over a long period based on wavelets, a 2-frame tracker and an outlier rejection process. Yin and Collins [32] build an adaptive view-dependent appearance model on-line. The model is made of patches selected around Harris corners. Model and target patches are matched using correlation, and the change in position, rotation and scale is obtained with the Procrustes algorithm.

Another successful line of work in object tracking without *a priori* training is based on classification instead of modeling. Collins and Liu [33] presented a system based on background/foreground discrimination. Avidan [34] presents one of the many systems that use machine learning to classify patches [35, 36]. Avidan uses weak classifiers trained every frame and AdaBoost to combine them. Pilet *et al.* [37] train keypoint classifiers using Random Trees that are able to recognize hundreds of keypoints in real-time.

Simultaneous Modeling and Tracking (SMAT) [1] is in line with methods like Lucas-Kanade, relying on matching to track patches. Lukas-Kanade extracts a template at the beginning of the sequence and uses it for tracking, and will fail if the appearance of the patch changes considerably. Matthews *et al.* [38] proposed an *strategic update* of the template, which keeps the template from the first frame to correct errors that appear in the localization. When the error is too high, the update is blocked. In [39], a solution is proposed with fixed

91 template that adaptively detected and selected the window around the features. SMAT  
 92 builds a more complex model based on incremental clustering.

93 In this paper we combine concepts from active models with the incremental clustering  
 94 proposed in SMAT. The texture model is created online, making the model adaptative,  
 95 while the shape model is learnt offline. The clustering used by SMAT has some limitations,  
 96 and we propose some modifications to obtain a more robust model and better tracking. We  
 97 name the approach *Robust SMAT* for this reason.

98 Evaluation of face tracking methods is performed in most works with images captured  
 99 indoors. Some authors use freely available image sets, but most of them test on internal  
 100 datasets created by them, which limits the validity of a comparison with other systems. Only  
 101 a few authors [40][41] have used images recorded in a vehicle, but the number of samples is  
 102 limited. To the best of our knowledge, there is no publicly available video dataset of people  
 103 driving, either in a simulator or in a real road. We propose a new dataset that covers such  
 104 scenarios.

### 105 3. Robust Simultaneous Modeling and Tracking

106 This section describes the Simultaneous Modeling and Tracking (SMAT) of Dowson and  
 107 Bowden [1], and some modifications we propose to improve its performance. SMAT tries to  
 108 build a model of appearance of features and how their positions are related (the structure  
 109 model, or *shape*), from samples of texture and shape obtained in previous frames.

110 The models of appearance and shape are independent. Fitting is performed in the same  
 111 fashion of ASM: the features are first found separately using correlation, and then their  
 112 final positions are constrained by the shape model. If the final positions are found to be  
 113 reliable and not caused by fitting errors, the appearance model is updated, otherwise it is  
 114 left unchanged. Figure 1 shows a flow chart of the algorithm.

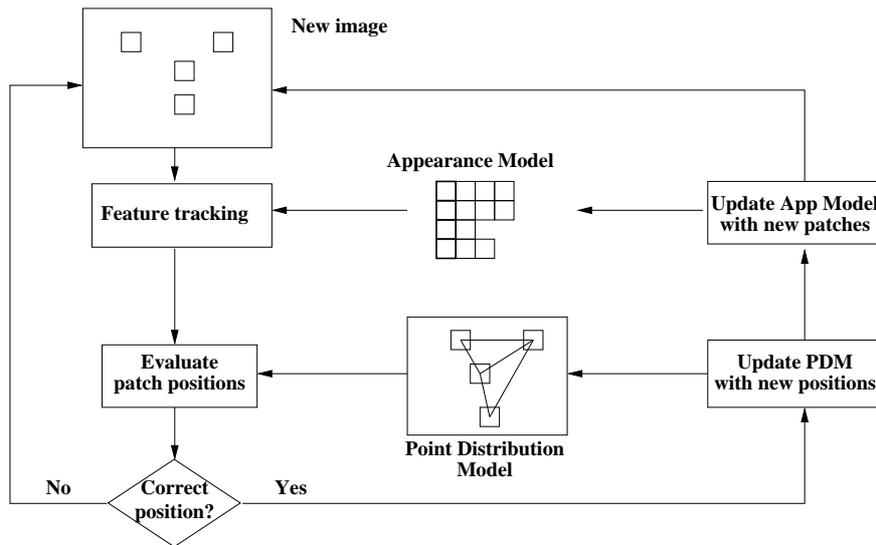


Figure 1: SMAT block diagram

115 *3.1. Appearance modeling*

116 Each one of the possible appearances of an object, or a feature of it, can be considered  
 117 as a point in a feature space. Similar appearances will be close in this space, away from  
 118 other points representing dissimilar appearances of the object. These groups of points, or  
 119 clusters, form a mixture model that can be used to define the appearance of the object.

120 SMAT builds a library of exemplars obtained from previous frames, image patches in this  
 121 case. Dowson and Bowden defined a series of clusters by their median patch, also known  
 122 as *representative*, and their variance. A new incoming patch is made part of the cluster if  
 123 the distance between it and the median of the cluster is below a threshold that is a function  
 124 of the variance. The median and variance of a cluster are recalculated everytime a patch is  
 125 added to it. Up to  $M$  exemplars per cluster are kept. If the size limit is reached, the most  
 126 distant element from the representative is removed.

127 Everytime a cluster is updated, the weight of the clusters is recalculated as in equation  
 128 1:

$$w_k^{(t+1)} = \begin{cases} (w_k^{(t)} + \alpha) \frac{1}{1+\alpha} & \text{if } k = k_u \\ w_k^{(t)} \frac{1}{1+\alpha} & \text{otherwise} \end{cases} \quad (1)$$

129 where  $\alpha \in [0, 1)$  is the learning rate, and  $k_u$  is the index of the updated cluster. The  
 130 number of clusters is also limited to  $K$ . If  $K$  is reached, the cluster with the lowest weight  
 131 is discarded.

132 In a later work, Dowson *et al.* [42], introduced a different condition for membership,  
 133 that compares the probability of the exemplar belonging to foreground (a cluster) or to the  
 134 background

$$\frac{p(fg | d(x, \mu_n), \sigma_{fg_n})}{p(bg | d(x, \mu_n), \sigma_{bg_n})} \quad (2)$$

135 where  $\sigma_{fg_n}$  is obtained from the distances between the representative and the other exemplars  
 136 in the cluster, and  $\sigma_{bg_n}$  is obtained from the distances between the representative and the  
 137 exemplars in the cluster offset by 1 pixel.

138 We have found that this clustering method can be improved in several ways. The adapt-  
 139 ing nature of the clusters could theoretically lead two or more clusters to overlap. However,  
 140 in our tests we have observed that the opposite is much more frequent: the representative  
 141 of the cluster rarely changes after the cluster has reached a certain number of elements.

142 Outliers can be introduced in the model in the event of an occussion of the face by a  
 143 hand or other elements like a scarf. In most cases, these exemplars would be far away from  
 144 the representative in the cluster. To remove them and reduce memory footprint, SMAT  
 145 keeps up to  $M$  exemplars per cluster. If the size limit is reached, the most distant element  
 146 from the representative is removed. When very similar patches are constantly introduced,  
 147 one of them will be finally chosen as the median, and the variance will decrease, overfitting  
 148 the cluster and discarding valuable exemplars. At a frame rate of 30 fps, with  $M$  set to 50,  
 149 the cluster will overfit in less than 2 seconds. This would happen even if the exemplar to be  
 150 removed is chosen randomly. This procedure will discard valuable information and future,  
 151 subtle changes to the feature will lead to the creation of another cluster.

152 We propose an alternative clustering method, named *leaderP*, to partially solve these and  
 153 other problems. The method is a modification of the *leader* algorithm [43, 44], arguably the  
 154 simplest and most frequently used incremental clustering method. In *leader*, each cluster  $\mathcal{C}_i$   
 155 is defined by only one exemplar, and a fixed membership threshold  $T$ . It starts by making the  
 156 first exemplar the *representative* of a cluster. If an incoming exemplar fulfills being within  
 157 the threshold  $T$  it is marked as member of that cluster, otherwise it becomes a cluster on  
 158 its own. The pseudocode is shown in algorithm 1.

---

**Algorithm 1** Leader clustering

---

```

1: Let  $C = \{\mathcal{C}_1, \dots, \mathcal{C}_n\}$  be a set of  $n$  clusters, with weights  $\{w_1^t, \dots, w_n^t\}$ 
2: procedure LEADER( $E, C$ ) ▷ cluster patch  $E$ 
3:   for all  $\mathcal{C}_i \in C$  do
4:     if  $d(\mathcal{C}_k, E) < T$  then ▷ Check if patch  $E \in \mathcal{C}_k$ 
5:       UPDATEWEIGHTS( $w_1^t, \dots, w_n^t$ ) ▷ As in equation 1
6:       return
7:     end if
8:   end for
9:   Create new cluster  $\mathcal{C}_{n+1}$ , with  $E$  as representative.
10:  Set  $w_{n+1}^{t+1} \leftarrow 0$  ▷ Weight of new cluster  $\mathcal{C}_{n+1}$ 
11:   $C \leftarrow C \cup \mathcal{C}_{n+1}$  ▷ Add new cluster to the model
12:  if  $n + 1 > K$  then ▷ Remove the cluster with lowest weight
13:    Find  $\mathcal{C}_k \mid w_k \leq w_i \quad i = 1, \dots, n$ 
14:     $C \leftarrow C \setminus \mathcal{C}_k$ 
15:  end if
16: end procedure

```

---

159 On the other hand, *leaderP* keeps the first few exemplars added to the cluster are kept,  
 160 up to  $P$ . The median of the cluster is chosen as the representative, as in the original  
 161 clustering of Dowson and Bowden. When the number of exemplars in the cluster reaches  $P$ ,  
 162 all exemplars but the representative are discarded, and it starts to work under the leader  
 163 algorithm.  $P$  is chosen as a small number (we use  $P = 10$ ). The membership threshold is  
 164 however flexible: the distances between the representative and each of the exemplars that  
 165 are found to be members of the cluster is saved, and the variance of those distances is used  
 166 to calculate the threshold. Because the representative is fixed and distance is a scalar, many  
 167 values can be kept in memory without having a impact on the overall performance. Keeping  
 168 more values reduces the risk of overfitting.

169 The original proposal of SMAT used Mutual Information (MI) as a distance measure to  
 170 compare the image patches, and found it to perform better than Sum of Squared Differences  
 171 (SSD), and slightly better than correlation in some tests. Any definition of distance could be  
 172 used. We have also tested Zero-mean Normalized Cross-Correlation (ZNCC). Several types  
 173 of warping were tested in [42]: translation, euclidean, similarity and affine. The results  
 174 showed an increasing failure rate as the degrees of freedom of the warps increased. Based

175 on this, we have chosen to use the simplest, and the patches are only translated depending  
 176 on the point distribution model.

### 177 3.2. Shape model

178 In the original SMAT of Dowson and Bowden, the shape was also learned on-line. The  
 179 same clustering algorithm was used, but the membership of a new shape to a cluster was  
 180 calculated using Mahalanobis distance.

181 Our method relies on the pre-learned shape model. The restrictions on using a pre-  
 182 learned model for shape are less than those for an appearance model, as it is of lower  
 183 dimensionality and the deformations are easier to model. It has been shown [45] that  
 184 location and tracking errors are mainly due to appearance, and that a generic shape model  
 185 for faces is easier to construct. We use the method of classic ASM [27], which applies PCA  
 186 to a set of samples created by hand and extracts the mean  $\mathbf{s}_0$  and an orthogonal vector  
 187 basis  $(\mathbf{s}_1, \dots, \mathbf{s}_N)$ . The shapes are first normalized and aligned using Generalized Procrustes  
 188 Analysis [46].

189 Let  $\mathbf{s} = (x_0, y_0, \dots, x_{n-1}, y_{n-1})$  be a shape. A shape can be generated from this base as

$$\mathbf{s} = \mathbf{s}_0 + \sum_{i=1}^m p_i \cdot \mathbf{s}_i \quad (3)$$

190 Using  $L_2$  norm, the coefficients  $\mathbf{p} = (p_1, \dots, p_N)$  can be obtained for a given shape  $\mathbf{s}$  as  
 191 a projection of  $\mathbf{s}$  on the vector basis

$$\mathbf{p} = \mathbf{S}^T(\mathbf{s} - \mathbf{s}_0), \quad p_i = (\mathbf{s} - \mathbf{s}_0) \cdot \mathbf{s}_i \quad (4)$$

192 where  $\mathbf{S}$  is a matrix with the eigenvectors  $\mathbf{s}_i$  as rows. The estimation of  $\mathbf{p}$  with equation  
 193 4 is very sensitive to the presence of outlier points: a high error value from one point will  
 194 severely influence the values of  $\mathbf{p}$ . We use M-estimators [47] to solve this problem. This  
 195 technique has been applied to ASM and AAM in previous works [48, 49], so it is only briefly  
 196 presented here.

197 Let  $\mathbf{s}$  be a shape, obtained by fitting each feature independently. The function to mini-  
 198 mize is

$$\arg \min_{\mathbf{p}} \sum_{i=1}^{2n} \rho(r^i, \theta) \quad (5)$$

199 where  $\rho : \mathbb{R} \times \mathbb{R}^+ \rightarrow \mathbb{R}^+$  is an M-estimator, and  $\theta$  is obtained from the standard deviation  
 200 of the residues [50].  $r^i$  is the residue for coordinate  $i$  of the shape

$$r^i = \mathbf{x}^i - (\mathbf{s}_0^i + \sum_{j=1}^m p_j \mathbf{s}_j^i) \quad (6)$$

201 where  $\mathbf{x}^i$  are the points of the shape  $\mathbf{s}$ , and  $\mathbf{s}_j^i$  is the  $i$ th element of the vector  $\mathbf{s}_j$ .

202 Minimizing function 5 is a case of re-weighted least squared. The weight decreases more  
 203 rapidly than the square of the residue, and thus a point with error tending to infinite will  
 204 have zero weight in the estimation.

205 Several robust estimators have been tested: *Huber*, *Cauchy*, *Gaussian* and *Tukey* func-  
206 tions [50]. A study was made in [19] that resulted in similar performance for all of them  
207 in a similar scenario to that of this paper, and Huber function was chosen. Huber function  
208 performs correctly up to a number of outliers of 50% of the points.

209 We use the 20-point distribution of the BioID database [51]. Data from this database  
210 was used to train the model. This distribution places the points in some of the most salient  
211 locations of the face, and has been used in several other works [40].

## 212 4. Tests and results

213 This section presents the video sequences used to test different tracking algorithms in  
214 a driving scenario. The dataset contains most actions that appear in everyday driving  
215 situations. A comparison between our approach and SMAT is presented. Additionally, we  
216 compare R-SMAT results with the recently introduced Stacked Trimmed ASM (STASM).

### 217 4.1. Test set

218 Driving scenarios present a series of challenges for a face tracking algorithm. Drivers  
219 move constantly, rotate their head (self-occluding part of the face) or occlude their face with  
220 their hands (or other elements such as glasses). If other people are in the car, talking and  
221 gesturing are common. There are also constant background changes and, more importantly,  
222 frequent illumination changes, produced by shadows of trees or buildings, streets lights,  
223 other vehicles, etc. A considerable amount of test data is needed to properly evaluate the  
224 performance of a system under all these situations.

225 A new video dataset has been created, with sequences of subjects driving outdoor, and  
226 in a simulator. The RobeSafe Driver Monitoring Video (RS-DMV) dataset contains 10  
227 sequences, 7 recorded outdoors (*Type A*) and 3 in a simulator (*Type B*).

228 Outdoor sequences were recorded on RobeSafe’s vehicle moving at the campus of the  
229 University of Alcalá. Drivers were fully awake, talked frequently with other passengers in  
230 the vehicle and were asked to look regularly to the rear-view mirrors and operate the car  
231 sound system. The cameras are placed over the dashboard, to avoid occlusions caused by  
232 the wheel. All subjects drove the same streets, shown in figure 2.



(a) Trajectory of the vehicle during recordings

Figure 2: Trajectory of the vehicle (map from *maps.google.com*)

233 The length of the track is around 1.1 km. The weather conditions during the recordings  
234 were mostly sunny, which made noticeable shadows appear on the face. Figure 3 shows a  
235 few samples from these video sequences.



Figure 3: Samples of outdoor videos

236 *Type B* sequences were recorded in a realistic truck simulator. Drivers were fully awake,  
237 and were presented with a demanding driving environment where many other vehicles were  
238 present and potentially dangerous situations took place. These situations increase the prob-  
239 ability of small periods of distraction leading to crashes or near-crashes. The sequences try  
240 to capture both distracted behaviour and the reaction to dangerous driving situations.

241 A few images from *Type B* sequences can be seen in figure 4. The recording took place  
242 in a low-light scenario that approached nighttime conditions. This forced the camera to  
243 increase exposure time to a maximum, which led to motion blur being present during head  
244 movements. Low power near-IR illumination was used in some of the sequences to increase  
245 the available light.



Figure 4: Samples of sequences in simulator

246 The outdoor sequences are around 2 minutes long, and sequences in the simulator are  
247 close to 10 minutes in length. The algorithms in this paper were tested on images of ap-  
248 proximately  $320 \times 240$  pixels, but high resolution images were acquired so they can be used  
249 in other research projects. The images are  $960 \times 480$  pixels for the outdoor sequences and  
250  $1392 \times 480$  for the simulator sequences, and are stored without compression. Frame rate is  
251 30 frames per second in both cases. The camera has a  $2/3''$  sensor, and used 9mm standard

252 lenses. Images are grayscale. The recording software controlled camera gain using values of  
 253 the pixels that fell directly on the face of the driver.

254 The RS-DMV is publicly available, free of charge, for research purposes. Samples and  
 255 information on how to obtain the database are available at the authors' webpage<sup>1</sup>.

#### 256 4.2. Performance evaluation

257 Performance of the algorithms is evaluated as the error between the estimated position  
 258 of the features and their actual position, as given by a human operator. Hand-marking is a  
 259 time consuming task, and thus not all frames in all videos have been marked. Approximately  
 260 1 in 30 frames (1 per second) has been marked in the sequences in RS-DMV. We call this  
 261 frames *keyframes*.

262 We used the metric  $m_e$ , introduced by Cristinacce and Cootes [40]. Let  $\mathbf{x}^i$  be the points  
 263 of the ground-truth shape  $\mathbf{s}$ , and let  $\hat{\mathbf{x}}^i$  be the points of the estimated shape  $\hat{\mathbf{s}}$ . Then,

$$m_e = \frac{1}{ns} \sum_{i=1}^n d^i, \quad d^i = \sqrt{(\mathbf{x}^i - \hat{\mathbf{x}}^i)^T (\mathbf{x}^i - \hat{\mathbf{x}}^i)} \quad (7)$$

264 where  $n$  is the number of points and  $s$  is the inter-ocular distance. We also discard the point  
 265 on the chin and the exterior of the eyes, because their location changes much from person  
 266 to person. Moreover, the variance of their position when marked by human operators is  
 267 greater than for the other points. Because only 17 points are used, we note the metric as  
 268  $m_{e17}$ . In the event of a tracking loss, or if the face can not be found, the value of  $m_{e17}$  for  
 269 that frame is set to  $\infty$ .

270 During head turns, the inter-eye distance reduces with the cosine of the angle. In these  
 271 frames,  $s$  is not valid and is calculated from its value on previous frames.

272 Handmarked points and software used to ease the marking process are distributed with  
 273 the RS-DMV dataset.

#### 274 4.3. Results

275 We tested the performance of R-SMAT approach on the RS-DMV dataset, as well as  
 276 that of SMAT. We compared these results with those obtained by STASM, using the imple-  
 277 mentation in [2].

278 One of the most remarkable problems of (R-)SMAT is that it needs to be properly  
 279 initialized, and the first frames of the sequence are key to building a good model. We  
 280 propose STASM to initialize (R-)SMAT in the first frame. STASM has been shown to be  
 281 very accurate when the face is frontal. Nonetheless, a slightly incorrect initialization will  
 282 make (R-)SMAT track the (slightly) erroneous points. To decouple this error from the  
 283 evaluation of accuracy of (R-)SMAT in the tests, the shape was initialized in the first frame  
 284 with positions from the ground-truth data. At the end of this section, the performance of  
 285 R-SMAT with automatic initialization is evaluated.

---

<sup>1</sup>[www.robosafe.com/personal/jnuevo](http://www.robosafe.com/personal/jnuevo)

286 First, a comparison of the shape models is presented. With the best shape model, the  
 287 original clustering algorithm and the proposed alternative are evaluated. Results are pre-  
 288 sented for outdoor and simulator sequences separately, as each has specific characteristics  
 289 on their own.

290 The incremental shape model of SMAT was found to produce much higher error than  
 291 the pre-learned model. Figure 5 shows the cumulative distribution error of the incremental  
 292 shape model (*on-line*) with the robust pre-learned model (using Huber function) (*robust*).  
 293 For comparison purposes, the figure also shows the performance for the pre-learned shape  
 294 model fitted using a  $L_2$  norm (*non-robust*). All models use *leaderP* clustering, and patches  
 295 of  $15 \times 15$  pixels.

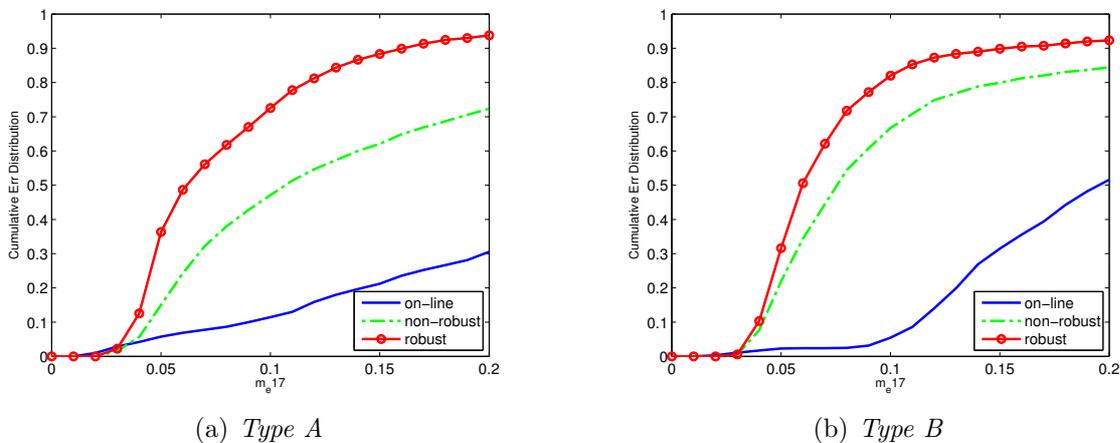


Figure 5: Performance of different shape models, with *leaderP* clustering

296 Clear improvements in performance are made by the change to a pre-learned model with  
 297 robust fitting. The robust, pre-learned shape model is very important in the first frames,  
 298 because it allows the model to have bigger certainties that the patches that are being included  
 299 correspond to correct positions. Robust shape model is used in the rest of the experiments  
 300 in this paper. Figure 6 shows the plot of the  $m_e17$  distance of both models in a sequence. A  
 301 clear example of the benefits of the robust model is depicted in figure 7. The online model  
 302 diverges as soon as a few points are occluded by the hand, while the robust model keeps  
 303 track of the face. The method is also able to keep track of the face during head rotations,  
 304 although with increased fitting error. This is quite remarkable for a model that has only  
 305 been trained with fully frontal faces.

306 Figure 8 shows the performance of the original SMAT clustering compared with the  
 307 proposed *leaderP* clustering algorithm, as implemented in R-SMAT.

308 R-SMAT presents much better performance than the original SMAT clustering. This is  
 309 specially clear in 8(b). We stated in 3.1 that the original clustering method could lead to  
 310 overfitting, and *type B* sequences are specially prone to this: patches are usually dark and  
 311 do not change much from frame to frame, and the subject does not move frequently. When a  
 312 movement takes place, it leads to high error values, because the model has problems finding  
 313 the features.

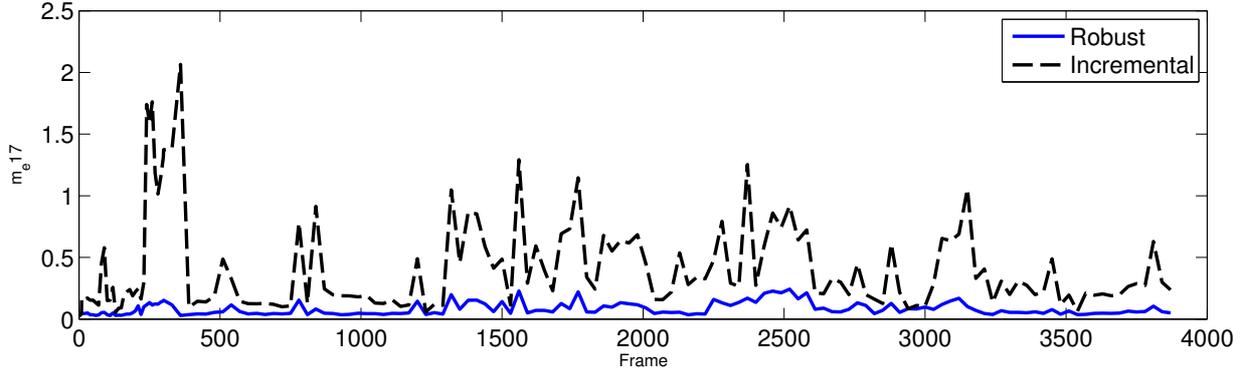


Figure 6:  $m_e17$  error for a sequence

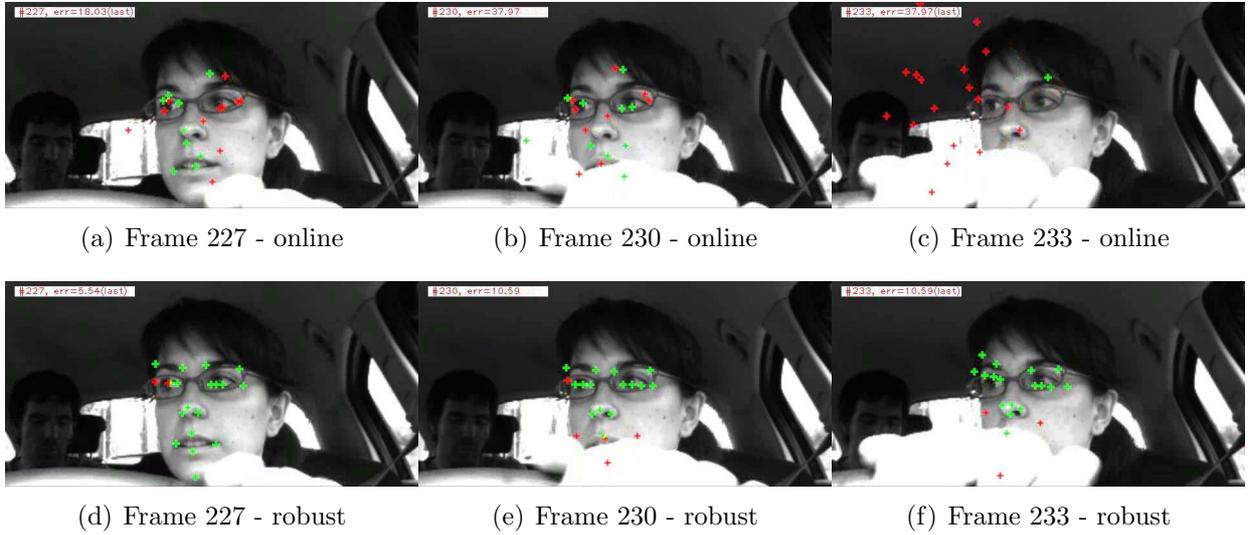


Figure 7: Samples of *type A* sequence #1. Outlier points are drawn in red.

314 Table 1 shows the tracking losses of SMAT and R-SMAT, as a percentage of the *keyframes*  
 315 in the sequences. Tracking losses were monitored by counting the points inside the face area,  
 316 detected with Viola&Jones algorithm [52]. Tracking was considered lost when more than  
 317 33% of the points were out of the box, or when the rotation of the model exceeded a pre-set  
 318 value. The model was then repositioned, simply by centering it on the Viola&Jones box.

		Mean	Maximum	Minimum
R-SMAT	<i>Type A</i>	0.99%	2.08%(seq. #7)	0%(seq. #1,#2,#6)
	<i>Type B</i>	0.71%	1.96%(seq. #9)	0%(seq. #10)
SMAT	<i>Type A</i>	1.77%	5.03%(seq. #4)	0%(seq. #1,#2,#5)
	<i>Type B</i>	1.03%	2.45%(seq. #9)	0%(seq. #10)

Table 1: Track losses for different clustering methods

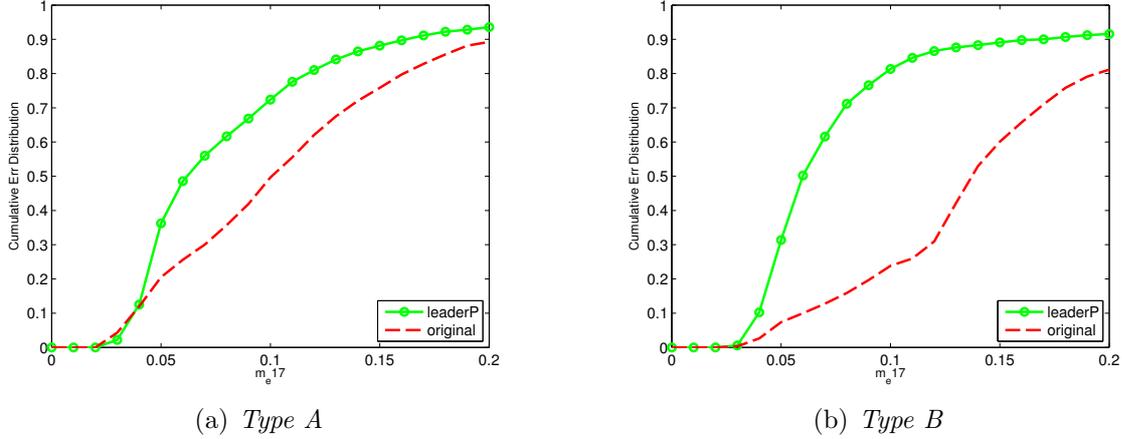


Figure 8: Comparison of the performance of clustering algorithms

319 *4.3.1. R-SMAT with automatic initialization*

320 Results presented above have been obtained initializing SMAT and R-SMAT with land-  
 321 marks from the handmarked ground-truth data. In a real scenario, an automatic algorithm  
 322 would be used to initialize SMAT and R-SMAT.

323 We have used STASM for this task. STASM has demonstrated high accuracy, but only  
 324 works properly when the face is frontal to the camera, and does not find the face otherwise.  
 325 Another problem, critical to our application, is that it does not work in real time. However,  
 326 a one-time delay can be considered acceptable. STASM was run on the first frame of each  
 327 sequence, and its estimation seeded the position of R-SMAT in that video. Figure 9 plots  
 328 the error distributions of R-SMAT when initialized with STASM and from ground-truth  
 329 (manual). The error of STASM is also shown.

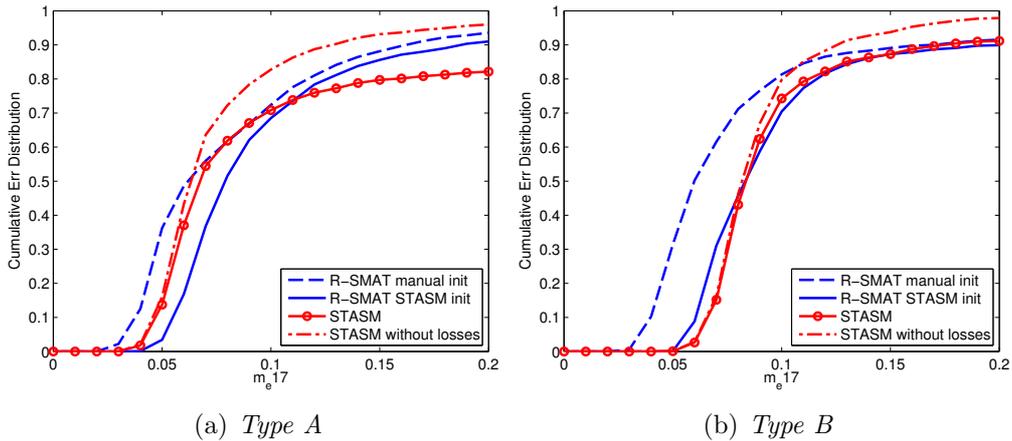


Figure 9: Comparison of the performance of STASM and SMAT

330 As expected, the figure shows the results worsen for both types of sequences. But the  
 331 lost accuracy is relatively small, with a 5% loss at  $m_{e17} = 0.1$  for *type A* sequences, and

332 10% loss for *type B*. The mean of the  $m_{e17}$  error of STASM in the first frame is 0.0571 for  
 333 *type A* sequences and 0.0805 for *type B*.

334 STASM is plotted in figure 9 with and without considering frames where the face was  
 335 not found (losses). For all types of sequences, R-SMAT initialized manually outperforms  
 336 STASM when losses are considered. Expectedly, STASM shows better accuracy than R-  
 337 SMAT when tracking is not lost (i.e., when the face is frontal). R-SMAT initialized with  
 338 STASM performs almost identically as STASM in figure 9(b), and slightly worse for *type A*  
 339 sequences.



Figure 10: Examples from sequences with R-SMAT fitted

340 Figure 10 shows a few frames with R-SMAT fitted to the face of the drivers in moments  
 341 of the sequences that reflect some of the challenges the system has to face. Figures 10(a)-  
 342 10(d) contain drivers talking and gesturing. Drivers in sequences #5 and #6 talk frequently  
 343 and no tracking loss results from these actions. The system is able to work with drivers that  
 344 also wear a beard, as in figure 10(d). Examples of head turns appear in figures 10(e) and  
 345 10(f). Most occlusions in the sequences are caused by the presence of a hand in front of the  
 346 camera. If the occlusion is partial as in figure 10(g), the tracker is able to correctly position  
 347 the model.

348 Samples of R-SMAT fitted to one of the simulator sequences appear in figures 10(h) and  
 349 10(i). Despite the low-light conditions, a low fitting error is obtained. The driver in the  
 350 sequences talks and gestures frequently.

351 Figures 11 and 12 plot the error for R-SMAT and STASM for sequences #6 and #7.  
 352 Dots on the STASM curve mark keyframes where the face was not found. A case of quick  
 353 illumination change, due to shadows of trees by the road is found in sequence #6 around  
 354 frame 1400. A R-SMAT track loss can be observed around frame 2400 of the latter sequence,  
 355 due to a total occlusion of the face. Figure 12 also shows that fitting error is higher during  
 356 head turns, but in most cases track is not lost.

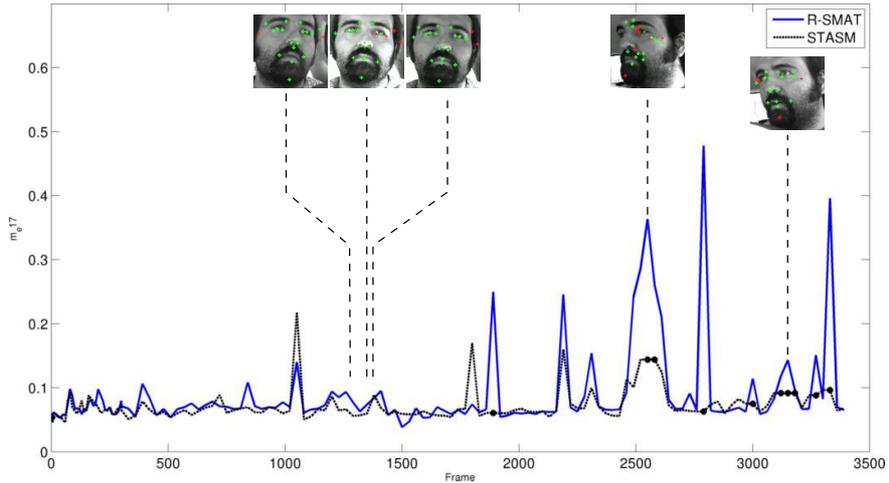


Figure 11: Error plots for STASM and R-SMAT in sequence #6

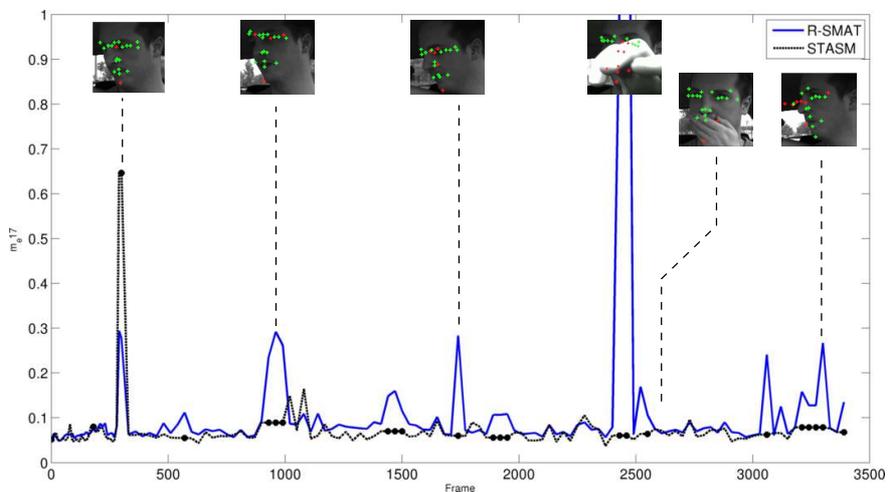


Figure 12: Error plots for STASM and R-SMAT in sequence # 7

#### 357 4.4. Timings

358 One of the most important requirements for R-SMAT is for it to run in real-time. Table  
 359 2 summarizes the average execution speed for R-SMAT in frames per second, for some  
 360 representative configurations. The worst frame processing times are close to the limit, but

361 these are extreme cases that occur infrequently. Processing times for STASM are also  
 362 included for comparison. STASM executes the whole initialization process for each frame,  
 363 and does not use the position found in the previous frame. This would result in a shorter  
 364 search time.

Configuration	Mean (fps)	Sdv (fps)	Worst frame (fps)
R-SMAT	112.56	32.80	36.86
STASM	2.17	0.13	1.96

Table 2: Execution time of R-SMAT and STASM in frames per second

365 The tests were run on a Xeon 2.2 GHz, running GNU/Linux, with GCC 4.2 as compiler.  
 366 Multi-threading was not used and compiler optimizations were disabled (-O0). Times on  
 367 the table refer to the actual tracking and the tracking loss detection, and do not consider  
 368 time employed in the display of results, loading of data and saving results to the hard drive.

## 369 5. Conclusions and future work

370 This paper has presented a face tracking method based on automatic appearance mod-  
 371 eling, to be used as part of a driver monitoring application.

372 Monitoring a driver with computer vision is a complex task, and proper performance  
 373 evaluation of a method requires a comprehensive set of test data. A new video dataset  
 374 (RS-DMV) has been created, comprised of sequences recorded in a real scenario, and in a  
 375 truck simulator in low light conditions. In the first set of sequences, subjects were asked to  
 376 drive a vehicle at the University campus. Drivers in the simulator were fully awake, and were  
 377 presented with dangerous situations that would highlight distractions. RS-DMV dataset has  
 378 been used to test the methods in this paper, and is freely available for research purposes.

379 The proposed face tracking method is an active model. The shape model is built offline  
 380 from handmarked data, and Huber function is used for fitting. The texture model is created  
 381 online using incremental clustering, in a similar fashion to SMAT. We have presented and  
 382 alternative incremental clustering algorithm, which addresses some of the weaknesses of  
 383 the SMAT proposal. The improvements of R-SMAT over the original SMAT have been  
 384 evaluated, and the performance of R-SMAT and STASM has been compared on the sequences  
 385 in RS-DMV. R-SMAT is able to process more than 100 frames per second, and obtains similar  
 386 accuracy to STASM. The source code of R-SMAT is available from the authors.

387 Future work will explore ways to make R-SMAT fully autonomous by improving the  
 388 incremental shape model. Texture clustering has shown to be reliable, but better techniques  
 389 to remove outliers from the model are needed. Including a multi-scale approach to appear-  
 390 ance modeling and model fitting would be of help in other scenarios where the size of the  
 391 face changes noticeably. The R-SMAT has been used to track and model faces in this paper,  
 392 but can be extended to other deformable objects. The RS-DMV dataset will be extended  
 393 with more sequences, more drivers and more diverse scenarios. Finally, the R-SMAT is to  
 394 be made part of a driver monitoring system.

## 395 Acknowledgments

396 The authors would like to thank Pedro Jiménez, Ivan García and Noelia Hernández of  
397 RobeSafe for their work in recording the sequences, Sebastian Bronte for his help in marking  
398 the images, as well as the drivers that participated. The outdoor recordings were made  
399 under project MOVI<sup>2</sup>CON (TRA2005-08529-C02-02) and the simulator recordings under  
400 CABINTEC project (PSE-370100-2007-2). This work was supported in part by the Spanish  
401 Ministry of Science and Innovation under DRIVER-ALERT Project (TRA2008-03600), and  
402 Comunidad de Madrid under project RoboCity2030 (S-0505/CPI/000176). J. Nuevo was  
403 working under a researcher training grant from the Education Department of the Comunidad  
404 de Madrid and the European Social Fund.

- 405 [1] R. Dowson, N.D.H.; Bowden, Simultaneous modeling and tracking (SMAT) of feature sets, in: IEEE  
406 Conference on Computer Vision and Pattern Recognition 2005, Vol. 2, 2005, pp. 99–105.
- 407 [2] S. Milborrow, F. Nicolls, Locating facial features with an extended active shape model, in: ECCV,  
408 2008, pp. 504–513, <http://www.milbo.users.sonic.net/stasm>.
- 409 [3] T. A. Dingus, S. Klauer, V. L. Neale, A. Petersen, S. E. Lee, J. Sudweeks, M. A. Perez, J. Hankey,  
410 D. Ramsey, S. Gupta, C. Bucher, Z. R. Doerzaph, J. Jermeland, R. Knipling, The 100-car naturalistic  
411 driving study, Tech. rep., Virginia Tech Transportation Institute, NHTSA (Apr. 2006).
- 412 [4] A. Kircher, M. Uddman, J. Sandin, Vehicle control and drowsiness, Tech. Rep. VTI-922A, Swedish  
413 National Road and Transport Research Institute (2002).
- 414 [5] Volvo Car Corp., Driver alert control, <http://www.volvocars.com> (2008).
- 415 [6] DaimlerAG, Attention assist, <http://www.daimler.com> (Jun. 2009).
- 416 [7] H. Ueno, M. Kaneda, M. Tsukino, Development of drowsiness detection system, in: Proceedings of  
417 Vehicle Navigation and Information Systems Conference, 1994, pp. 15–20.
- 418 [8] E. Murphy-Chutorian, M. Trivedi, Head pose estimation in computer vision: A survey, *IEEE Trans.*  
419 *Pattern Anal. Mach. Intell.* 31 (4) (2009) 607–626. doi:10.1109/TPAMI.2008.106.
- 420 [9] P. Smith, M. Shah, N. da Vitoria Lobo, Determining driver visual attention with one camera, *IEEE*  
421 *Trans. Intell. Transp. Syst.* 4 (4) (2003) 205–218.
- 422 [10] J. Wu, T. Chen, Development of a drowsiness warning system based on the fuzzy logic images analysis,  
423 *Expert Systems with Applications* 34 (2) (2008) 1556–1561.
- 424 [11] Y. Matsumoto, A. Zelinsky, An algorithm for real-time stereo vision implementation of head pose and  
425 gaze direction measurements, in: *Procs. IEEE 4th Int. Conf. Face and Gesture Recognition*, 2000, pp.  
426 499–505.
- 427 [12] T. Victor, O. Blomberg, A. Zelinsky, Automating the measurement of driver visual behaviours using  
428 passive stereo vision, in: *Procs. Int. Conf. Series Vision in Vehicles VIV9*, Brisbane, Australia, 2001.
- 429 [13] M. Kuttila, *Methods for Machine Vision Based Driver Monitoring Applications*, Ph.D. thesis, VTT  
430 Technical Research Centre of Finland (2006).
- 431 [14] T. D’Orazio, M. Leo, C. Guaragnella, A. Distante, A visual approach for driver inattention detection,  
432 *Pattern Recognition* 40 (8) (2007) 2341–2355.
- 433 [15] Seeing Machines, Driver state sensor, <http://www.seeingmachines.com/dss.html> (Aug. 2007).
- 434 [16] SmartEyeAG, AntiSleep, [www.smarteye.se](http://www.smarteye.se) (2009).
- 435 [17] Smart Eye AB, Image capturing device with reflex reduction, Patent, EP 1349487B1 (Dec. 2001).
- 436 [18] Q. Ji, Z. Zhu, P. Lan, Real-time nonintrusive monitoring and prediction of driver fatigue, *IEEE Trans.*  
437 *Veh. Technol.* 53 (4).
- 438 [19] L. M. Bergasa, J. Nuevo, M. A. Sotelo, R. Barea, E. López, Real-time system for monitoring driver  
439 vigilance, *IEEE Trans. Intell. Transp. Syst.* 7 (1) (2006) 1524–1538.
- 440 [20] M. Yang, D. Kriegman, N. Ahuja, Detecting faces in images: a survey, *IEEE Trans. Pattern Anal.*  
441 *Mach. Intell.* 24 (1) (2002) 34–58. doi:<http://doi.ieeecomputersociety.org/10.1109/34.982883>.
- 442 [21] P. Belhumeur, J. Hespanha, D. Kriegman, Eigenfaces vs. fisherfaces: Recognition using class specific

- 443 linear projection, *IEEE Trans. Pattern Anal. Mach. Intell.* 19 (7) (1997) 711–720, special Issue on Face  
444 Recognition.
- 445 [22] J. M. Buenaposada, E. Muñoz, L. Baumela, Recognising facial expressions in video sequences, *Pattern*  
446 *Analysis and Applications* 11 (1) (2008) 101–116.
- 447 [23] G. D. Hager, P. N. Belhumeur, Efficient region tracking with parametric models of geom-  
448 etry and illumination, *IEEE Trans. Pattern Anal. Mach. Intell.* 20 (10) (1998) 1025–1039.  
449 doi:<http://dx.doi.org/10.1109/34.722606>.
- 450 [24] J. M. Buenaposada, E. Muñoz, L. Baumela, Efficiently estimating facial expression and illumination in  
451 appearance-based tracking, in: *Proc. British Machine Vision Conference*, Vol. 1, 2006, pp. 57–66.
- 452 [25] F. Jurie, M. Dhome, Hyperplane approximation for template matching, *IEEE Trans. Pattern Anal.*  
453 *Mach. Intell.* (2002) 996–1000.
- 454 [26] S. Baker, I. Matthews, Equivalence and efficiency of image alignment algorithms, in: *IEEE Computer*  
455 *Society Conference on Computer Vision and Pattern Recognition*, Vol. 1, 2001, pp. 1090–1097.
- 456 [27] T. Cootes, C. Taylor, D. Cooper, J. Graham, Active Shape Models-Their Training and Application,  
457 *Computer Vision and Image Understanding* 61 (1) (1995) 38–59.
- 458 [28] I. Dryden, K. Mardia, *Statistical Shape Analysis*, John Wiley & Sons, 1998.
- 459 [29] T. F. Cootes, G. J. Edwards, C. J. Taylor, Active appearance models, *IEEE Trans. Pattern Anal.*  
460 *Mach. Intell.* 23 (2001) 681–685.
- 461 [30] B. Lucas, T. Kanade, An iterative image registration technique with an application to stereo vision,  
462 in: *International Joint Conference on Artificial Intelligence*, Vol. 3, 1981, pp. 674–679.
- 463 [31] A. Jepson, D. Fleet, T. El-Maraghi, Robust Online Appearance Models for Visual Tracking, *IEEE*  
464 *Trans. Pattern Anal. Mach. Intell.* (2003) 1296–1311.
- 465 [32] Z. Yin, R. Collins, On-the-fly Object Modeling while Tracking, in: *Computer Vision and Pattern*  
466 *Recognition, 2007. CVPR'07. IEEE Conference on*, 2007, pp. 1–8.
- 467 [33] R. Collins, Y. Liu, M. Leordeanu, Online Selection of Discriminative Tracking Features, *IEEE Trans.*  
468 *Pattern Anal. Mach. Intell.* (2005) 1631–1643.
- 469 [34] S. Avidan, Ensemble Tracking, *IEEE Trans. Pattern Anal. Mach. Intell.* (2007) 261–271.
- 470 [35] H. Grabner, H. Bischof, On-line boosting and vision, in: *Proc. CVPR*, Vol. 1, 2006, pp. 260–267.
- 471 [36] M. Pham, T. Cham, Online Learning Asymmetric Boosted Classifiers for Object Detection, in: *Com-*  
472 *puter Vision and Pattern Recognition CVPR, 2007*, pp. 1–8.
- 473 [37] J. Pilet, V. Lepetit, P. Fua, Real-time non-rigid surface detection, in: *IEEE Conference on Computer*  
474 *Vision and Pattern Recognition 2007*, San Diego, CA, 2005, pp. 822–828.
- 475 [38] I. Matthews, T. Ishikawa, S. Baker, The Template Update Problem, *IEEE Trans. Pattern Anal. Mach.*  
476 *Intell.* (2004) 810–815.
- 477 [39] S. Segvic, A. Remazeilles, F. Chaumette, Enhancing the point feature tracker by adaptive modelling of  
478 the feature support, in: *European Conf. on Computer Vision, ECCV'2006*, Vol. 3952 of *Lecture Notes*  
479 *in Computer Science*, Graz, Austria, 2006, pp. 112–124.
- 480 [40] D. Cristinacce, T. Cootes, Feature Detection and Tracking with Constrained Local Models, in: *17th*  
481 *British Machine Vision Conference*, 2006, pp. 929–938.
- 482 [41] L. Bergasa, J. Buenaposada, J. Nuevo, P. Jimenez, L. Baumela, Analysing Driver's Attention Level  
483 using Computer Vision, in: *Intelligent Transportation Systems. ITSC 2008. 11th International IEEE*  
484 *Conference on*, 2008, pp. 1149–1154.
- 485 [42] N. Dowson, R. Bowden, N-tier simultaneous modelling and tracking for arbitrary warps, in: *Proc. of*  
486 *the 17th British Machine Vision Conference. British Machine Vision Association*, Vol. 1, 2006, p. 6.
- 487 [43] J. Hartigan, *Clustering Algorithms*, John Wiley & Sons, Inc. New York, NY, USA, 1975.
- 488 [44] H. Spath, V. Bull, *Cluster analysis algorithms for data reduction and classification of objects*, Ellis  
489 Horwood, 1980.
- 490 [45] R. Gross, I. Matthews, S. Baker, Generic vs. person specific active appearance models, *Image and*  
491 *Vision Computing* 23 (11) (2005) 1080–1093.
- 492 [46] J. Gower, Generalized procrustes analysis, *Psychometrika* 40 (1) (1975) 33–51.
- 493 [47] P. J. Huber, *Robust Statistics*, Wiley Series in Probability and Mathematical Statistics, Wiley-

- 494 Interscience, 1981.
- 495 [48] M. Rogers, J. Graham, Robust active shape model search, Lecture Notes in Computer Science (2002)  
496 517–530.
- 497 [49] R. Gross, I. Matthews, S. Baker, Constructing and fitting active appearance models with occlusion, in:  
498 Proceedings of the IEEE Workshop on Face Processing in Video, 2004, p. 72.
- 499 [50] Z. Zhang, Parameter estimation techniques: A tutorial with application to conic fitting, Image and  
500 vision Computing 15 (1) (1997) 59–76.
- 501 [51] O. Jesorsky, K. Kirchberg, R. Frischholz, et al., Robust face detection using the haus-  
502 dorff distance, Proceedings of Audio and Video based Person Authentication (2001) 90–  
503 95<http://www.bioid.com/downloads/facedb/>.
- 504 [52] P. Viola, M. Jones, Robust real-time face detection, International Journal of Computer Vision 57 (2)  
505 (2004) 137–154.