

Exploiting structural constraints for visual object tracking

Wassim Bouachir^{a,*}, Guillaume-Alexandre Bilodeau^a

^a*LITIV lab., Department of Computer and Software Engineering,
École Polytechnique de Montréal,
P.O. Box 6079, Station Centre-ville, Montréal
(Québec), Canada, H3C 3A7*

Abstract

This paper presents a novel structure-aware method for visual tracking. The proposed tracker relies on keypoint regions as salient and stable elements that encode the object structure efficiently. In addition to the object structural properties, the appearance model also includes global color features that we first use in a probabilistic approach to reduce the search space. The second step of our tracking procedure is based on keypoint matching to provide a preliminary prediction of the target state. Final prediction is then achieved by exploiting object structural constraints, where target keypoints vote for the corrected object location. Once the object location is obtained, we update the appearance model and structural properties, allowing to track targets with changing appearance and non-rigid structures. Extensive experiments demonstrate that the proposed Structure-Aware Tracker (SAT) outperforms recent state-of-the-art trackers in challenging scenarios, especially when the

*Corresponding author

Email addresses: `wassim.bouachir@polymtl.ca` (Wassim Bouachir),
`gabilodeau@polymtl.ca` (Guillaume-Alexandre Bilodeau)

target is partly occluded and in moderately crowded scenes.

Keywords: Object tracking, Structure-aware tracker, keypoint, SIFT, keypoint layout.

1. Introduction

Model-free visual tracking is one of the most active research areas in computer vision [1, 2, 3]. With a *model-free* tracker, the only available input is the target state annotated in the first video frame. Tracking an object is thus a challenging task due to (1) the lack of sufficient information on object appearance, (2) the inaccuracy in distinguishing the target from the background (which is generally done using a geometric shape), and (3) the object appearance change caused by various perturbation factors (*e.g.* noise, occlusion, motion, illumination, etc.).

This work aims to develop a novel visual tracking method to handle real life difficulties, particularly when tracking an object in a moderately crowded scene in the presence of distracting objects similar to the target, and in the case of severe partial occlusion. The robustness of a tracking algorithm in handling these situations is determined by two major aspects: the target representation and the search strategy. The target representation refers to the appearance model that represents the object characteristics while the search strategy deals with how the search of the target is performed on every processed frame. The main contributions and differences of our work from previous works are on both aspects. In the proposed tracker, the target representation includes color features for coarse localization of the target, and keypoints for encoding the object structure while adding distinctiveness and

robustness to occlusions. In our search strategy, probabilistic tracking and deterministic keypoint matching are used sequentially to provide a preliminary estimate of the target state. Object internal structural constraints are then applied in a correction step to find an accurate prediction. Our approach for representing the object structure is related to previous works on *context tracking* [4, 5, 6, 7, 8]. The main idea of *context tracking* is to consider the spatial context of the target including neighboring elements whose motion is correlated with the target. While the proposed approach is inspired by the idea of *context tracking*, in our work we exploit the spatial layout of keypoints to encode the internal structure of the target. More specifically, our contributions are:

1. A novel target representation model where local features are stored in a reservoir encoding recent and old structural properties of the target;
2. A new threefold search strategy that reduces the search space, tracks keypoints, and corrects prediction sequentially;
3. A discriminative approach that evaluates tracking quality online to determine if potential new target properties should be learned.

Extensive experiments on challenging video sequences show the validity of the proposed Structure-Aware Tracker (SAT) and its competitiveness with state-of-the-art trackers. A previous version of this work was presented at a conference [9]. This paper extends this previous work with a more complete review of related works, more details and depth in the explanation of the method, and additional experiments analyzing the tracker behavior in several situations.

46 This paper is organized as follows. In the next section, we review recent
47 works on keypoint tracking and context tracking which are related to our
48 algorithm. The proposed SAT algorithm is presented in section 3. Experi-
49 mental results are given and discussed in section 4. Section 5 concludes the
50 paper.

51 2. Related works

52 2.1. Keypoint tracking: from object context to object structure

53 Many tracking algorithms achieved good performances at a low complex-
54 ity by using a geometric shape to contain the target, and global features for
55 modeling [10, 11, 12]. Nevertheless, this approach is not designed to han-
56 dle occlusions, unless representing the target by multiple fragments to be
57 matched. Keypoint methods can handle the occlusion problem by establish-
58 ing partial correspondences that allow locating the occluded target. Unlike
59 fragment-based methods (where the target image region is divided randomly
60 or according to a regular grid), keypoint locations correspond to salient and
61 stable patches that can be invariantly detected under various perturbation
62 factors. Moreover, their spatial layout naturally encodes structural proper-
63 ties that can enhance the target model.

64 Due to these characteristics, keypoint-based methods have attracted much
65 attention during the last decade. In this approach, objects are modeled as a
66 set of keypoints detected by an external mechanism (i.e. a keypoint detec-
67 tor) [13, 14, 15]. After computing their descriptors, the object localization
68 can be achieved according to two possible approaches: matching in the case
69 of a generative approach, and classification in the case of a discriminative

70 approach. Generative trackers use a database where keypoint descriptors are
 71 stored. The descriptors are designed to be stable and invariant, and can be
 72 matched in a nearest-neighbor fashion. Discriminative approaches consider
 73 matching as a binary classification problem. Every feature is thus classified
 74 as belonging to the background, or to the tracked object. The classifier is
 75 built either via online learning, or offline, considering the background and
 76 the target observed under various transformations.

77 Some recent works on object tracking rely on target context to predict its
 78 state, which is often referred as *context tracking* [4, 5, 16, 7, 17]. According to
 79 this approach, it is necessary to consider target context to ensure the tracker
 80 robustness in most real life video surveillance applications. Following this
 81 principle, the authors in [4] use a *compagnion* to improve object tracking.
 82 This corresponds to image regions around the tracked object with the same
 83 movements as those of the target. In [5] the spatial context that can help
 84 the tracker includes multiple *auxiliary objects*. These objects have consis-
 85 tent motion correlation with the tracked target and thus help to avoid the
 86 drifting problem. In [16], Gu and Tomasi consider the spatial relationship
 87 between the target and similar objects and track all of them simultaneously
 88 to eliminate target confusion. In a more general approach, Grabner et al.
 89 introduced the notion of *supporters* defined as "*useful features for predicting*
 90 *the target object position*" [7]. These features do not belong to the target, but
 91 they move in a way that is statistically related to the motion of the target.
 92 They developed a method for discovering these local image features around
 93 the target, and demonstrated that motion coupling of *supporters* may allow
 94 locating the target even if it is completely occluded. In a later work, Dinh et

95 al. [17] used *supporters* for context tracking, and added the concept of *dis-*
96 *tracters* which are regions co-occurring with the target while having a similar
97 appearance. Their tracker explicitly handles situations where several objects
98 similar to the target are present.

99 *Context tracking* methods expanded the target model by exploiting the
100 motion correlation information in the scene. However, finding motion cor-
101 relation between objects is a costly task that often requires detecting and
102 analyzing features on the whole image, as in [18] where the authors detect
103 and analyze all local features in the scene, to keep only features which move
104 along with the target object. Furthermore, most of the proposed track-
105 ers were tested only on specific scenarios and in constrained environments,
106 where almost all the experiments were limited to proofs of concept. Our idea
107 of using structural constraints in the target appearance model is inspired
108 by *context tracking* methods. However, our motivations differ in an impor-
109 tant aspect since our model incorporates the internal structural information
110 of the target, and not the structural layout of different scene elements. In
111 our work, we show that the structural information of the target, encoded by
112 the keypoint spatial layout, allows achieving accurate tracking and handling
113 partial occlusion by inferring the position of the target using the unoccluded
114 features.

115 2.2. Tracking objects by structure

116 The idea of exploiting object structure for tracking was present, more
117 or less explicitly, in recent works. This is the so called *part-based tracking*
118 that relies on local components for target representation. The most common
119 way to encode object structure is the sparse representation such as in [19]

120 and [20]. In [19], the authors propose to use a histogram-based model that
 121 encodes the spatial information of the object patches. In a similar manner,
 122 Jia et al. sample a set of overlapped patches on the tracked object [20]. Their
 123 strategy includes an occlusion handling module allowing target localization
 124 by using only visible image patches.

125 Another approach for encoding structure consists in using keypoints, since
 126 they are more significant than random overlapped patches. In this direction,
 127 the authors in [21] model the target by a set of keypoint *manifolds* organized
 128 as a graph to explicitly represent the target structure. Each feature *manifolds*
 129 includes, in addition to the keypoint descriptor, a set of synthetic descriptors
 130 simulating possible variations of the original feature (under viewpoint and
 131 scale change). The target location is found by detecting keypoints on the
 132 current frame, matching them with those of the target model, and computing
 133 a homography for the correspondences. In [22], the authors include both
 134 random patches and keypoints in the target model. The random patches are
 135 described by their RGB color histograms and LBP (Local Binary Patterns)
 136 descriptors to form an appearance model. Keypoints are characterized by
 137 their spatial histograms to be considered as a structural model. Tracking
 138 then implies matching detected keypoints in the current frame with those of
 139 the object in the previous frame. Matched keypoints are utilized to construct
 140 a spatial histogram, which is used jointly with LBP and RGB histograms
 141 to locate the target. This approach exploits multiple object characteristics
 142 (LBP, color, Keypoints), but the object structural model captures only recent
 143 structural properties, as the spatial histogram considers only the keypoints
 144 that are matched with those of the target in the last frame.

145 In our work, we argue and demonstrate through our experiments that
 146 keypoint regions are more efficient than random patches in encoding the
 147 structure, as they correspond to salient and stable patches invariably de-
 148 tectable under several perturbation factors. Unlike in [22] where random
 149 regions are analyzed to extract local features, and [21] where keypoints are
 150 extracted from a region with a fixed size (with the assumption of small dis-
 151 placements), we use a probabilistic method to reduce the search space to the
 152 most likely image regions, based on the target’s global color features. Con-
 153 cerning the target structure, our structural model is not limited, like in [22]
 154 to recent properties, which would make it strongly related to the last predic-
 155 tion (and thus may be completely contaminated if the tracker drifts from the
 156 target). Instead, our representation includes both recent and old structural
 157 constraints in a reservoir of features. The local features and their structural
 158 constraints are learned online during tracking. The deletion of a given fea-
 159 ture is related to its persistence (not to its moment of occurrence), while the
 160 impact of its constraint depends on the persistence as well as the consistence
 161 of the feature. Every local feature expresses its structural constraint individ-
 162 ually by voting to possible target locations. Thus, our voting-based method
 163 preserves the object structure without requiring building and updating com-
 164 plex keypoint graphs, neither calculating homographies such as in [21]. Our
 165 method takes into consideration the temporal information of all the target’s
 166 model components. The target model is thus updated to reflect the object
 167 appearance changes including structure changes, which allows tracking ob-
 168 jects with non-rigid structures.

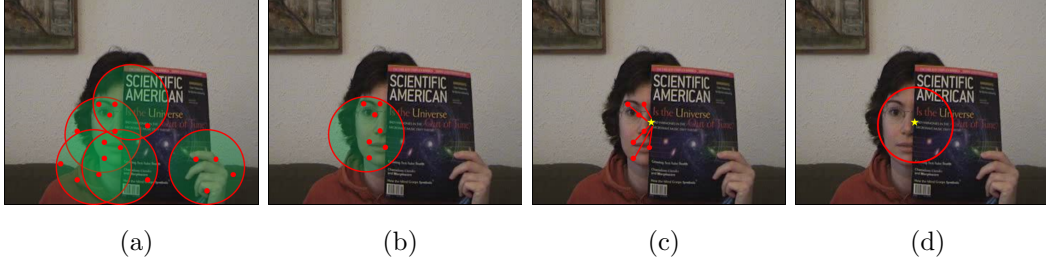


Figure 1: Illustration of the SAT algorithm steps when tracking a partly occluded face. **(a)**: Reducing the search space with a probabilistic method, based on color. Local features (red dots) are computed only on the obtained areas. **(b)**: Predicting a preliminary target state based on feature matching. **(c)**: Visible features vote for a new position (yellow star) by applying their structural constraints. **(d)**: The target state is corrected based on the new location

3. Proposed algorithm

3.1. Motivation and overview

The proposed method is illustrated in figure 1 where we aim to track a partly occluded face. First, we apply a color-based particle filtering. This allows to reduce the search space and provides a coarse estimation by considering only the best particles. Keypoints are then detected by analyzing the reduced search space as shown in figure 1a. The detected keypoints are matched with those of the target model, which leads to a preliminary estimate of the target location (see figure 1b).

Note that the preliminary prediction considers only the matching scores of the particles and thus does not guarantee an accurate localization. This is illustrated in figure 1b, where the circular shape representing the best particle includes pixels from the background and from the occluding object. Knowing

182 the internal structure of the target, our idea is to perform a correction step
 183 by applying internal structural constraints to improve target prediction. In
 184 practice, this is carried out by a voting mechanism where available features
 185 (unoccluded) determine the exact position of the target (figure 1c and 1d).
 186 Once the target is predicted, the appearance model including keypoints and
 187 their structural constraints is updated according to an evaluation criterion
 188 (that we define in section 3.5). The newly detected keypoints are added
 189 to the model while existing keypoints are re-evaluated based on two proper-
 190 ties. First, we consider the individual keypoint persistence represented by its
 191 weight value. The second property is the spatial consistency of the keypoint
 192 that depends on the motion correlation with the target center. If a keypoint
 193 of the background is erroneously included in the target model, these two
 194 voting parameters will reduce the effect of its vote until its removal from the
 195 model when its persistence decreases significantly. Our algorithm steps are
 196 explained in details in the following.

197 *3.2. Appearance Model*

198 Our appearance model describes the image region delimited by the circle
 199 that circumscribes the target. This is a multi-features model including (1)
 200 the color probability distribution represented by a weighted histogram, (2)
 201 a set of local descriptors computed for the detected keypoints within the
 202 target region, and (3) the target structural properties encoded by the voting
 203 parameters of keypoints. By constructing a m-bin histogram $\hat{\mathbf{q}} = \{\hat{q}\}_{u=1\dots m}$,
 204 with $\sum_{u=1}^m \hat{q}_u = 1$, some parts of the background may lie inside the circular
 205 kernel. As discussed in [23], these pixels will affect the color distribution and
 206 may cause tracking drift. To reduce the effect of these pixels, we use a kernel

207 function $k(x)$ that assigns smaller weights to pixels farther from the center.
 208 The color histogram is thus computed for the h pixels inside the target region
 209 according to the equation:

$$\hat{q}_u = \frac{1}{\sum_{i=1}^h k(d_i)} \sum_{i=1}^h k(d_i) \delta[c_i - u] \quad (1)$$

210 where $d_i \in [0, 1]$ is the normalized distance from the pixel x_i to the kernel
 211 center, c_i is the bin index for x_i in the quantized space, δ is the Kronecker
 212 delta function, and $k(d_i)$ is the tricube kernel profile defined by:

$$k(d_i) = \frac{70}{81} (1 - d_i^3)^3. \quad (2)$$

213 Note that the tricube function was selected among various kernel func-
 214 tions, as it allows the best experimental result. We also note that any other
 215 color space could be used instead of RGB.

216 The proposed system should be able to handle many difficult scenarios,
 217 such as occlusions and the presence of distracting objects. For, example, it
 218 has been shown that even for individuals of different races, the skin color
 219 distributions are very similar [24]. To ensure a more robust and distinctive
 220 feature set, the target reference model also includes SIFT keypoints [25] de-
 221 tected in the target region and stored in a *Reservoir of Features (RF)*. SIFT
 222 features increase the distinctiveness of the tracking algorithm to distinguish
 223 the target from other similar objects that may enter the field of view. In fact,
 224 SIFT was successfully used for distinguishing between multiple instances of
 225 the same object such as in the face recognition problem [26, 27, 28]. In this
 226 way, we implicitly handle situations where objects of the same category as
 227 the target co-occur (*e.g.* tracking a face in the presence of several faces), and

Algorithm 1 Reducing the search space at frame t

Input: frame t , particle states after processing frame $t - 1$

Output: reduced search space, new particle states

Assumption: processing frame t with $t > 2$

```

1: for  $i = 1$  to  $N$  do
2:   - generate a random number  $r_i \in [0, 1]$ 
3:   - find the particle  $s_{t-1}^{(j)}$  with the smallest  $j$  verifying  $c_{t-1}^{(j)} \geq r_i$ 
4:   - generate a new particle  $s_t^{(i)}$  for the selected particle  $s_{t-1}^{(j)}$ , with  $s_t^{(i)} =$ 
       $f(s_{t-1}^{(j)})$ 
5:   - evaluate similarity between  $\hat{p}_t^{(i)}$  and  $\hat{\mathbf{q}}$  {Eq. 3 and 4}
6:   - compute the weight  $\pi_t^{(i)}$  for  $s_t^{(i)}$ 
7: end for
8: - select the  $N^*$  best particles
9: - normalize weights  $\pi_t^{(n)}$  to get  $\sum_{n=1}^{N^*} \pi_t^{(n)} = 1$ 
10: - compute cumulative probabilities  $c_t^{(n)}$ 

```

228 thus we avoid using an additional mechanism to track and distinguish *dis-*
229 *tracters* as in [17]. Other than the keypoint descriptors, we also exploit the
230 spatial layout of keypoints to encode structural properties of objects. The
231 target structural constraints and the voting method that we use for predic-
232 tion correction are explained later. We note that our method is not specific
233 to SIFT. Even faster keypoint detector/descriptor combination may be used,
234 although SIFT remains one of the most reliable methods under various image
235 transformations [29].

236 *3.3. Reducing the search space*

237 The target search is firstly guided by particle filtering [30]. Each particle
 238 is a circular region characterized by its color distribution as explained above.
 239 The possible target states at frame t are represented by N weighted particles
 240 $\{s_t^{(i)} : i = 1, \dots, N\}$ where the weight $\pi_t^{(i)}$ reflects the importance of the
 241 particle. The weight of a generated particle $s_t^{(i)}$ depends on the similarity
 242 between its color distribution $\hat{p}_t^{(i)}$ and the reference color model $\hat{\mathbf{q}}$. We define
 243 the distance between the two distributions as:

$$d(\hat{\mathbf{q}}, \hat{p}_t^{(i)}) = \sqrt{1 - \rho[\hat{q}, \hat{p}_t^{(i)}]} \quad (3)$$

244 where

$$\rho[\hat{q}, \hat{p}_t^{(i)}] = \sum_{u=1}^m \sqrt{\hat{q}_u \cdot \hat{p}_{u,t}^{(i)}} \quad (4)$$

245 is the Bhattacharyya coefficient between $\hat{\mathbf{q}}$ and $\hat{p}_t^{(i)}$.

246 After generating N particles on the current frame, the area covered by the
 247 N^* best particles (i.e. the particles having the highest weights) is considered
 248 as a coarse estimation of the target state, and thus constitutes a reduced
 249 search space where keypoints will be detected and matched. Moreover, we
 250 use the N^* states selected at frame t for generating N particles at frame
 251 $t + 1$. Note that to simplify computations, we assign a cumulative weight
 252 $c^{(n)}$ to each pair $(s^{(n)}, \pi^{(n)})$ where $c^{(N^*)} = 1$. The cumulative weight $c^{(n)}$ for
 253 the n^{th} particle is calculated as $c^{(n)} = c^{(n-1)} + \pi^{(n)}$, where $c^{(1)} = \pi^{(1)}$. In this
 254 manner, for each particle $s^{(n)}$ we assign the interval $[c^{(n-1)}, c^{(n)}] \subset [0, 1]$ to
 255 allow a random particle selection (see steps 2 and 3 in Alg. 1). Our space
 256 reduction algorithm is summarized in Alg. 1.

257 3.4. Tracking keypoints

258 Keypoint detection and matching will consider only the reduced search
 259 space defined by the N^* best particles. By reducing the search region to the
 260 most important candidate particles, we avoid detecting features, computing
 261 local descriptors and matching them on the entire image.

262 The detected descriptors are then matched with those of the target model
 263 (features from the reservoir RF) based on the Euclidian distance. Similarly
 264 to the criterion used in [25], we determine if a match is correct by evaluating
 265 the ratio of distance from the closest neighbor to the distance of the second
 266 closest. For our algorithm, we keep only the matches for which the distance
 267 ratio is less than $\theta_m = 0.7$. Given the final set of matched pairs, we con-
 268 sider the particle having the highest matching score as a preliminary state
 269 of the target (see figure 1b). A more formal description of the preliminary
 270 prediction is provided in Alg. 2. Since the preliminary prediction considers
 271 only matching scores, without guaranteeing an accurate localization of the
 272 selected particle, the structural properties of the predicted region will be an-
 273 alyzed in a correction step to provide an accurate estimation of the target
 274 location.

275 3.5. Applying structural constraints

276 In this step, we aim to correct the preliminary prediction by applying a
 277 learned structural model of the target. The model is learned from reliable
 278 measurements (*i.e.* when a good tracking is achieved), and the internal
 279 structural properties are considered as a part of the object appearance model.

280 **Internal structural model.** The target keypoints extracted on the
 281 target region at different times of its lifecycle are stored in the reservoir

Algorithm 2 Preliminary prediction at frame t

```
1: - detect features on the reduced search space
2: for all detected_features  $f^{(i)}$  do
3:   - compute Euclidian distance with features from  $RF$ 
4:   - compute  $dist\_ratio = \frac{dist(f^{(i)}, closest\_neighbor)}{dist(f^{(i)}, 2^{nd\_closest\_neighbor})}$ 
5:   if  $dist\_ratio \leq \theta_m$  then
6:     - match  $f^{(i)}$  with  $closest\_neighbor$ 
7:     - update matching scores for the particles containing  $f^{(i)}$ 
8:   end if
9: end for
10: -  $preliminary\_prediction_t$  = the particle having the highest score
```

282 of features RF . Instead of automatically eliminating old keypoints, we only
283 remove those that become "non-persistent". RF is thus formed by recent and
284 old keypoints, representing both old and recent object properties. Other than
285 its descriptor summarizing the local gradient information, every keypoint is
286 characterized by a *voting profile* (μ, w, Σ) where:

- 287 • $\mu = [\Delta_x, \Delta_y]$ is the average offset vector that describes the keypoint's
288 location with respect to the target region center;
- 289 • w is the keypoint's weight considered as a persistence indicator to reflect
290 the feature co-occurrence with the target, and to allow eliminating "bad"
291 keypoints;
- 292 • Σ is the covariance matrix used as a spatial consistency indicator, de-
293 pending on the motion correlation with the target center.

294 **Voting.** Every matched keypoint f that is located on the preliminary
 295 target region votes for the potential object position \mathbf{x} by $P(\mathbf{x}|f)$. Note that
 296 we accumulate the votes for all the pixel positions inside the reduced search
 297 space. Given the *voting profile* of the feature f , we estimate the voting of f
 298 with the Gaussian probability density function:

$$P(\mathbf{x}|f) \propto \frac{1}{\sqrt{2\pi|\Sigma|}} \exp(-0.5(\mathbf{x}_f - \mu)^\top \Sigma^{-1}(\mathbf{x}_f - \mu)), \quad (5)$$

299 where \mathbf{x}_f is the relative location of \mathbf{x} with respect to the keypoint coordi-
 300 nates. The probability of a given pixel in the voting space is estimated by
 301 accumulating the votes of keypoints weighted by their persistence indicators
 302 w . The probability for a given pixel position \mathbf{x} in the voting space at time t
 303 is estimated by:

$$P_t(\mathbf{x}) \propto \sum_{i=1}^{|RF|} w_t^{(i)} P_t(\mathbf{x}|f^{(i)}) \mathbb{1}_{\{f^{(i)} \in F_t\}}, \quad (6)$$

304 where $\mathbb{1}_{\{f^{(i)} \in F_t\}}$ is the indicator function defined on the set RF (reservoir of
 305 features), indicating if the considered feature $f^{(i)}$ is among the matched target
 306 features set F_t at frame t . The target position is then found by analyzing
 307 the voting space and selecting its peak to obtain the corrected target state
 308 as shown in figure 1c.

309 **Update.** It has been previously shown that an adaptive target model,
 310 evolving during the tracking, is the key to good performance [31]. In our
 311 algorithm, the target model (including color, keypoints, and structural con-
 312 straints) is updated every time we achieve a good tracking using a discrimina-
 313 tive approach. Our definition of a good tracking is inspired by the Bayesian
 314 evaluation method used in [32], referred as *histogram filtering*. Using the
 315 target histogram $\hat{\mathbf{q}}$ (calculated for the target region annotated in the first

316 frame), and the background histogram $\hat{\mathbf{q}}_{bg}$ (calculated for the area outside
 317 the reduced search space), we compute a filtered histogram $\hat{\mathbf{q}}_{filt} = \hat{\mathbf{q}}/\hat{\mathbf{q}}_{bg}$ in
 318 every iteration. The latter represents the likelihood ratios of pixels belonging
 319 to the target. The likelihood ratios are used to calculate a backprojection
 320 map on the target region. Quality evaluation is done by analyzing the back-
 321 projection map and thresholding it to determine the percentage of pixels
 322 belonging to the target. Every time the evaluation procedure shows suffi-
 323 cient tracking quality, the target model is updated at frame t with a learning
 324 factor α as follows:

$$\hat{q}_t = (1 - \alpha)\hat{q}_{t-1} + \alpha\hat{q}_{new} \quad (7)$$

$$\hat{q}_{bg,t} = (1 - \alpha)\hat{q}_{bg,t-1} + \alpha\hat{q}_{bg,new} \quad (8)$$

$$w_t^{(i)} = (1 - \alpha)w_{t-1}^{(i)} + \alpha\mathbf{1}_{\{f^{(i)} \in F_t\}} \quad (9)$$

$$\Delta_{x,t}^{(i)} = (1 - \alpha)\Delta_{x,t-1}^{(i)} + \alpha\Delta_{x,new}^{(i)} \quad (10)$$

$$\Delta_{y,t}^{(i)} = (1 - \alpha)\Delta_{y,t-1}^{(i)} + \alpha\Delta_{y,new}^{(i)} \quad (11)$$

325 where $\mu_{new}^{(i)} = [\Delta_{x,new}^{(i)}, \Delta_{y,new}^{(i)}]$ is the current estimate of the voting vector
 326 for the feature $f^{(i)}$. After updating the feature weights, we remove from
 327 RF all the features whose the persistence indicators become less than the
 328 persistence threshold θ_p (*i.e.* $w_t^{(i)} \leq \theta_p$) regardless if they are recent or old,
 329 and we add the newly detected features with initial weight w_0 . Further,

we update the covariance matrix to determine the spatial consistency of the feature by applying:

$$\Sigma_t^{(i)} = (1 - \alpha)\Sigma_{t-1}^{(i)} + \alpha\Sigma_{new}^{(i)}, \quad (12)$$

where the new correlation estimate is:

$$\Sigma_{new}^{(i)} = (\mu_{new}^{(i)} - \mu_t^{(i)})(\mu_{new}^{(i)} - \mu_t^{(i)})^\top, \quad (13)$$

with $\mu_t^{(i)} = [\Delta_{x,t}^{(i)}, \Delta_{y,t}^{(i)}]$. Note that for the newly detected features, the preliminary persistence indicator is initialized to the covariance matrix $\Sigma = \sigma_0^2 I_2$, where I_2 is a 2 x 2 identity matrix. For consistent features, Σ decreases during the tracking, and thus their votes become more concentrated in the voting space. The overall algorithm is presented in Alg. 3.

4. Experiments

4.1. Experimental setup

We evaluated our SAT tracker by comparing it with four recent state-of-the-art methods on 11 challenging video sequences. Seven sequences of the dataset are publicly available and commonly used in the literature, while four are our own sequences¹. The *Tiger 1*, *Tiger2* and *Cliff bar* are provided in [1] and the *David indoor* and *Sylvester* are from [33]. The *Girl* and *occluded face 1* video sequences are respectively from [34] and [35]. The sequences *jp1*, *jp2*, *wdesk*, and *wbook* (with 608, 229, 709, and 581 frames respectively) were captured in our laboratory using a Sony SNC-RZ50N camera. The video

¹Our sequences are available at <http://www.polymtl.ca/litiv/en/vid/>.

frames are 320x240 pixels captured at a frame rate of 15 fps. For quantitative evaluation, we manually labeled the ground truth of our four sequences. Some of the sequences are available only in grayscale format (*Tiger 1*, *Tiger2*, *Sylvester*, and *Cliff bar*). For these videos, we slightly adapted our algorithm (especially the color model) to use grayscale information instead of RGB color information.

The four methods that we used for our comparison are the SuperPixel Tracker (SPT) [36], the Sparsity-Based Collaborative Tracker (SBCT) [19], the Adaptive Structural Tracker (AST) [20], and the Online Multiple Support Instance Tracker (OMSIT) [37]. The source codes of these trackers are available on the authors' respective websites. The authors also provide various parameter combinations. For fairness, we tuned the parameters of their methods so that for every video sequence, we always use the best combination among the ones that they proposed. Most of the parameters of SAT were set to default values for all the sequences, and only three parameters were tuned to optimize the performance of the tracker:

- N^* : the number of particles defining the reduced search space.
- θ_u : the threshold on the percentage of pixels belonging to the target that is required to update the appearance model.
- θ_p the persistence threshold used to determine if the keypoint should be removed from the reservoir.

table 1 shows the optimized parameter values for 5 video sequences from our dataset.

parameters	<i>girl</i>	<i>tiger 1</i>	<i>David indoor</i>	<i>occluded face 1</i>	<i>Wdesk</i>
N^*	30	100	100	40	80
θ_u	0.6	0.75	0.55	0.7	0.65
θ_p	0.3	0.4	0.2	0.2	0.3

Table 1: The optimized parameter values used in SAT with each video from the subset including *girl*, *tiger 1*, *David indoor*, *occluded face 1*, and *Wdesk*.

371 We quantitatively evaluated the performance of the trackers using the
 372 success rate and the average location error. To measure the success rate, we
 373 calculate for each frame the Overlap Ratio $OR = \frac{area(P_r \cap G_r)}{area(P_r \cup G_r)}$, where P_r is the
 374 predicted target region and G_r is the ground truth target region. Tracking is
 375 considered as a success for a given frame, if OR is larger than 0.5. The eval-
 376 uation of the Center Location Error (CLE) is based on the relative position
 377 errors between the center of the tracking result and that of the ground truth.
 378 Table 2 presents the success rates and the average center location errors for
 379 the compared methods. In order to analyze in depth the compared meth-
 380 ods on several video sequences, we also prepared two plots for every video
 381 sequence: 1) the center location error versus the frame number presented in
 382 figure 6, and 2) the overlap ratio versus the frame number presented in figure
 383 7. These plots are useful for understanding more in details the behavior of
 384 the trackers since the success rate and the average location error just sum-
 385 marize the performance of the tracker on a given sequence. Note that we
 386 averaged the results over five runs in all our experiments.

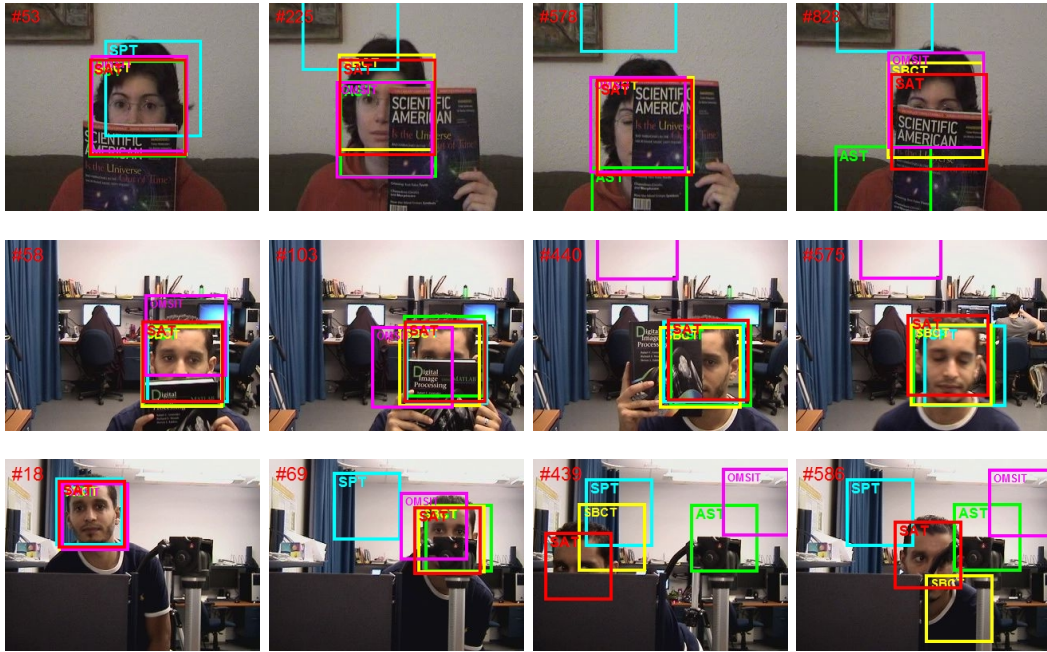


Figure 2: Tracking results for video sequences with long-term occlusions: *Occluded face 1*, *Wbook*, *Wdesk*. Green, magenta, yellow, cyan, and red rectangles correspond to results from AST, OMSIT, SBCT, SPT, SAT.

Sequence	SPT		SBCT		AST		OMSIT		SAT	
	S	E	S	E	S	E	S	E	S	E
<i>David indoor</i>	62	36	60	34	38	69	<i>63</i>	<i>27</i>	100	10
<i>girl</i>	<i>84</i>	9	2	201	18	53	1	66	85	<i>10</i>
<i>occluded face 1</i>	6	117	100	5	26	85	81	23	100	<i>14</i>
<i>tiger 1</i>	61	<i>17</i>	25	108	31	38	3	75	<i>51</i>	15
<i>tiger 2</i>	<i>46</i>	<i>23</i>	16	189	31	29	6	45	70	16
<i>Sylvester</i>	39	32	49	34	<i>73</i>	10	3	99	79	<i>14</i>
<i>Cliff bar</i>	52	22	24	77	70	35	8	74	<i>60</i>	<i>25</i>
<i>Jp1</i>	18	35	78	18	<i>84</i>	<i>17</i>	4	97	89	7
<i>Jp2</i>	39	<i>31</i>	<i>55</i>	69	<i>55</i>	45	17	39	94	7
<i>Wdesk</i>	14	80	<i>57</i>	<i>34</i>	32	81	10	123	90	11
<i>Wbook</i>	99	11	100	5	100	<i>9</i>	9	132	100	12
<i>average</i>	47	<i>38</i>	<i>52</i>	70	51	43	19	73	84	13

Table 2: Success rate (S) and average location error (E) results for SAT and the four other trackers: **Bold red** font indicates best results, *blue italics* indicates second best.

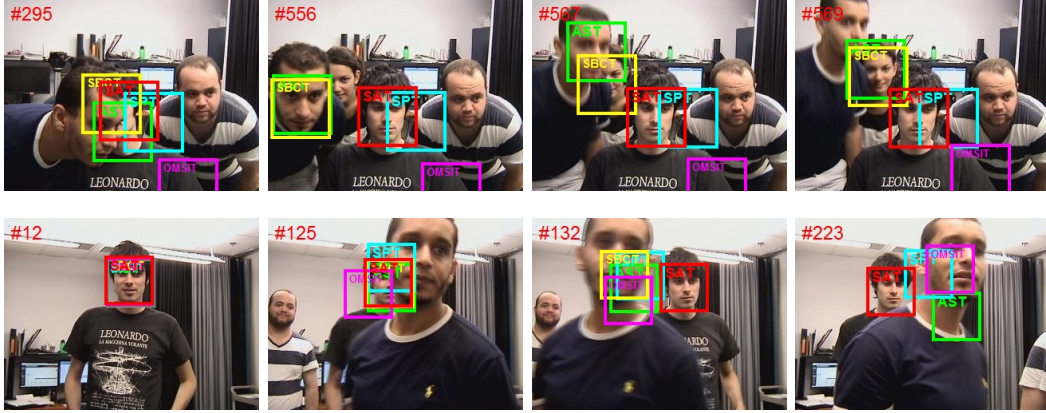


Figure 3: Screenshots of face tracking in moderately crowded scenes under short-term occlusions. In the *Jp1* sequence (first row), the tracked face is the one that is in the center of the scene. The same person is tracked while he is walking in the *Jp2* sequence. Green, magenta, yellow, cyan, and red rectangles correspond to results from AST, OMSIT, SBCT, SPT, SAT.

4.2. Experimental results

Long-time occlusion: Figure 2 demonstrates the performance of the compared trackers when tracking faces under long-time partial occlusions. In the *Occluded face 1* and the *wbook* sequences, the target faces remain partially occluded for several seconds while they barely move. The corresponding plots in figures 6 and 7 show that some trackers drift away from the target face, to track the occluding object (*e.g.* between frames 200 and 400 in *Occluded face 1*). Because it is specifically designed to handle partial occlusions via its structure-based model, our tracker was able to track the faces successfully in practically all the frames. SBCT has also achieved a good performance with a slightly lower average location error. In fact, SBCT is also designed to handle occlusions using a scheme that considers only the

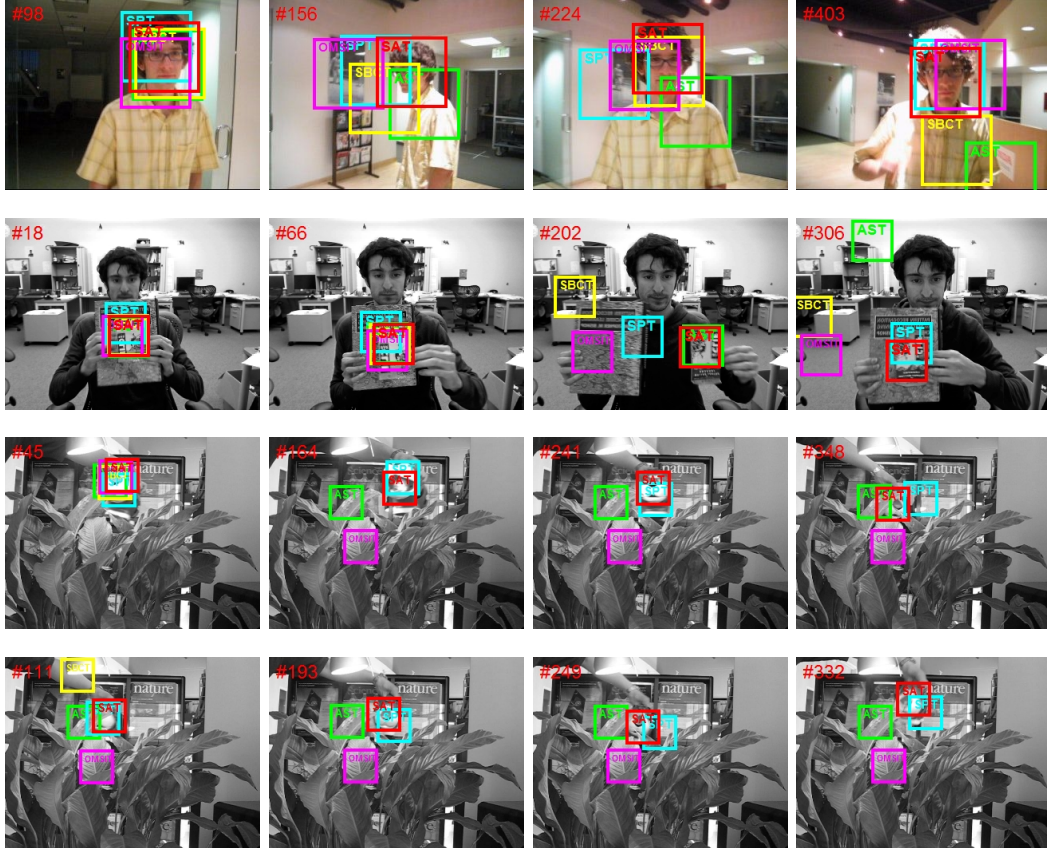


Figure 4: Screenshots of tracking results for some of the sequences with illumination change (*david indoor*) and background clutter (*Cliff bar*, *Tiger1*, *Tiger2*). Green, magenta, yellow, cyan, and red rectangles correspond to results from AST, OMSIT, SBCT, SPT, SAT.

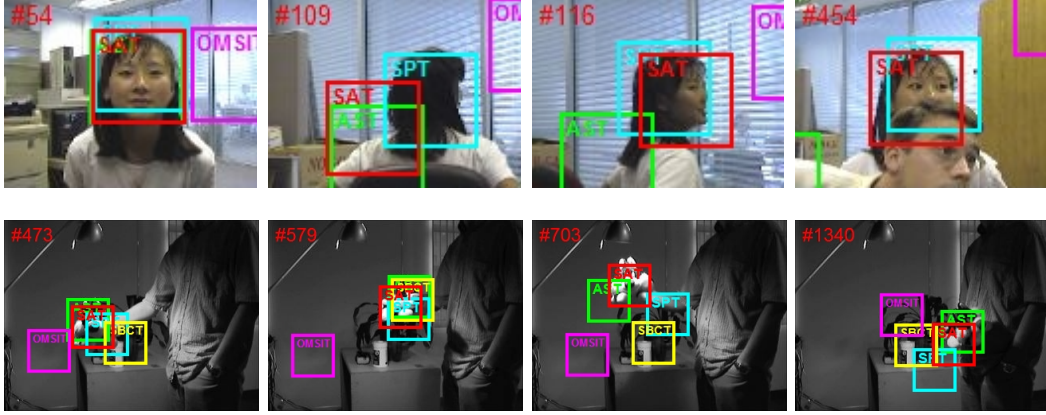


Figure 5: Tracking results for video sequences with abrupt motion and/or out of plane rotation: *Girl* and *Sylvester* sequences. Green, magenta, yellow, cyan, and red rectangles correspond to results from AST, OMSIT, SBCT, SPT, SAT.

399 patches that are not occluded. The target face in *Wdesk* undergoes severe
 400 partial occlusions many times while moving behind structures of the back-
 401 ground. SAT and SBCT track the target correctly until frame 400. At this
 402 point the person performs large displacements, and SBCT drifts away from
 403 the face. Nevertheless, our tracker continues the tracking successfully while
 404 the tracked person is trying to hide behind structures of the background,
 405 achieving a success rate of 90%. The superiority of the proposed method
 406 in this experiment highlights the importance of using structural constraints
 407 defined by keypoint regions that are more invariant than the patches used in
 408 SBCT when such a situation occurs.

409 **Moderately crowded scenes:** Figure 3 presents the results of face
 410 tracking in a moderately crowded scene (four persons). In the *Jp1* video,
 411 we aim to track a target face in presence of other faces that may partially
 412 occlude the target. Although the success rates of 84% and 78% respectively

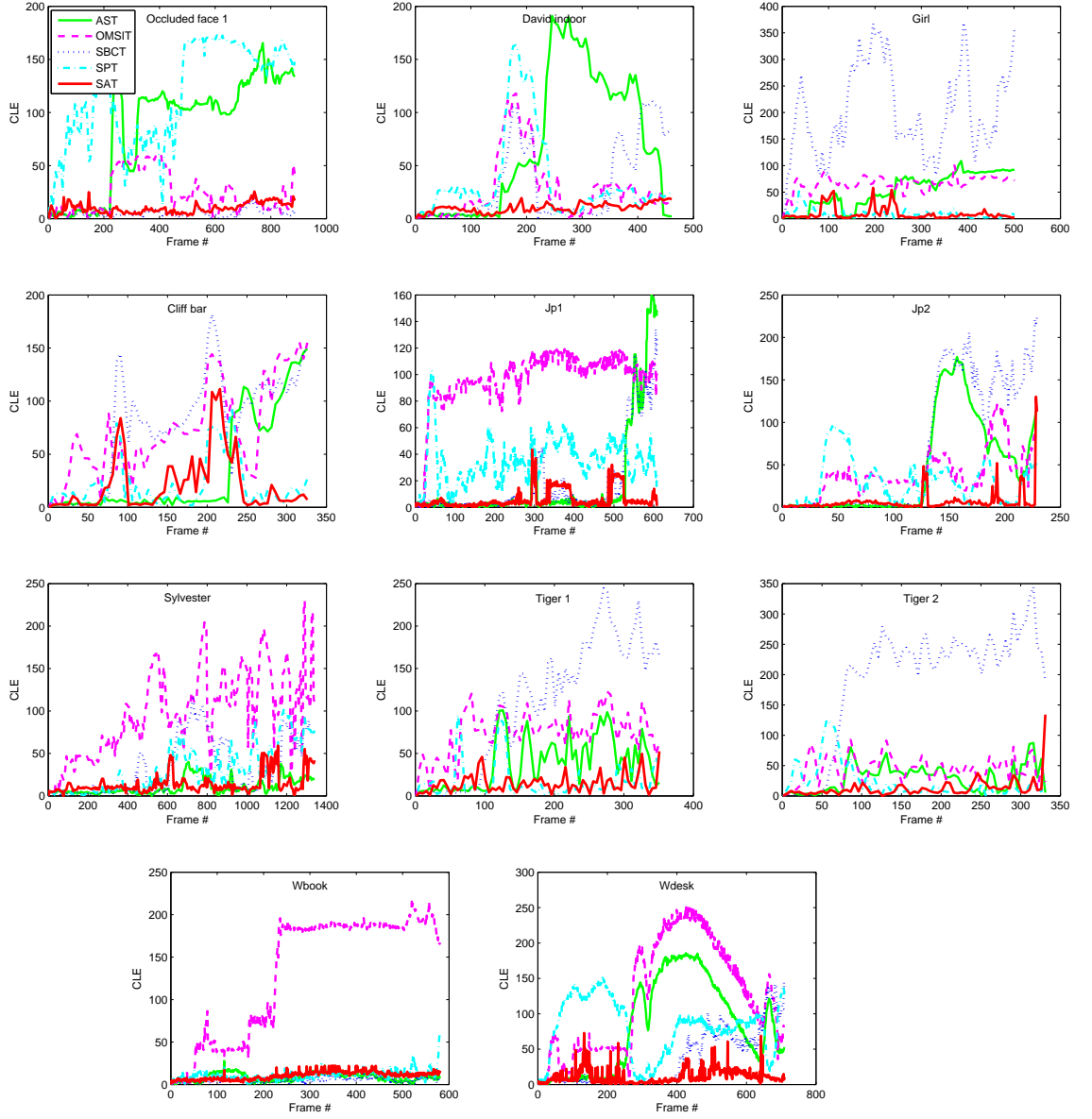


Figure 6: Center location error plots.

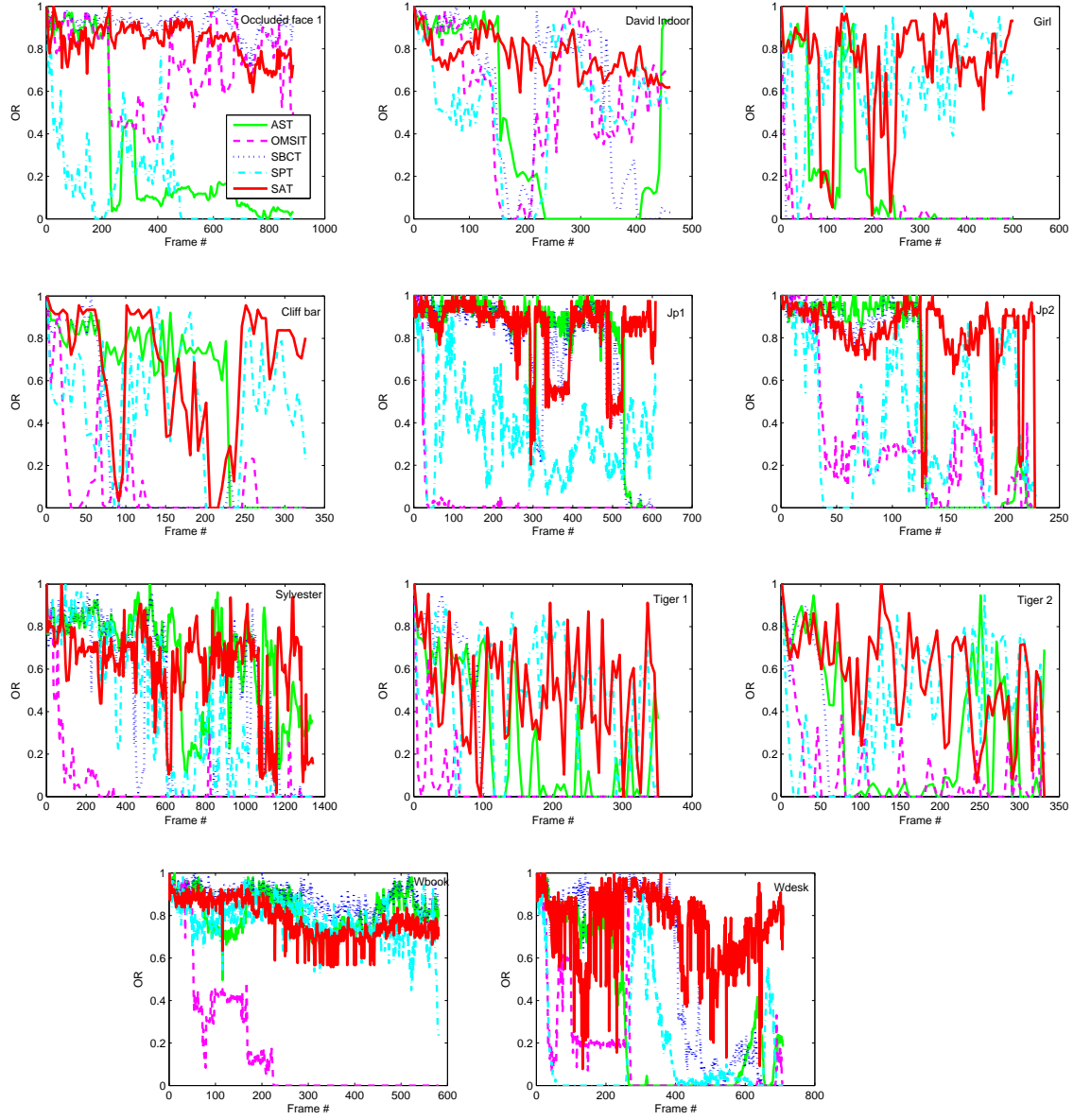


Figure 7: Overlap ratio plots.

413 for AST and SBCT indicate good performance in general, the two trackers
 414 drift twice, first at frame 530, and a second time at frame 570, to track other
 415 faces occluding or neighboring the target face. However, our tracker is not
 416 affected by the presence of similar objects around the target, even if partial
 417 occlusion occurs. This is mainly due to the distinctiveness of SIFT features
 418 compared to the local patches used in AST and SBCT to characterize the
 419 target. In this manner, SIFT features allow our tracker to handle situations
 420 where multiple instances of the same target object co-occur. In the *jp2*
 421 sequence, we track a walking person in a moderately crowded scene with
 422 four randomly moving persons. Here, we track a person’s face that crosses
 423 in front or behind another walking person that may completely occlude the
 424 target for a short time. Except the proposed method, none of the trackers
 425 is able to relocate the target after full occlusion by another person. For
 426 example, SBCT confused the target with the occluding face like in the video
 427 sequence *Jp1*. In this situation, SAT detects a total occlusion (since no
 428 features are matched). Our tracker continues searching the target based on
 429 color similarity without updating the appearance model. Tracking is finally
 430 recovered as soon as a small part of the target face becomes visible and
 431 feature matching becomes possible again.

432 **Illumination change:** In the *David indoor* video, the illumination changes
 433 gradually as the person moves from a dark room to an illuminated area (see
 434 figure 4). While most of the trackers were able to keep track of the person in
 435 more than 60% of the frames, SAT was the only tracker to achieve a success
 436 rate of 100%. In addition, SAT had the best performance on the *Sylvester*
 437 sequence in which the target object appearance changes drastically due to

438 abrupt illumination change. These two experiments show the superiority of
 439 our appearance model, which is the only one among the five models, to in-
 440 clude keypoints that are robust against lighting variations. Note that every
 441 time we update the reservoir of features, we replace the descriptors of all
 442 matched keypoints by their latest version computed on the current frame.
 443 This technique helps also to reflect appearance changes of keypoint regions
 444 (caused by illumination, viewpoint change, etc.), which facilitates matching
 445 features.

446 **Background clutters:** In the *Cliff bar* video, the background (the book)
 447 and the target have similar textures. Figure 4 shows that SBCT and OMSIT
 448 drift away from the target in most video frames. AST, SPT, and the pro-
 449 posed tracker were able to achieve a better performance despite the difficulty
 450 of this sequence. In fact, the target undergoes drastic appearance changes
 451 due to high motion blur. This caused drifts for all trackers several times
 452 (*e.g.* see the corresponding CLE and OR plots at frame 80). In the *Tiger*
 453 *1* and *Tiger 2* sequences, the tracked object exhibits fast movements in a
 454 cluttered background with frequent and various occlusion level. Owing to
 455 our voting mechanism that predicts the exact position of the target from the
 456 visible keypoints, our SAT tracker overcomes the frequent occlusion problem
 457 outperforming the other methods. All the other methods fail to locate the
 458 stuffed animal, except SPT that achieved better results due to its discrimi-
 459 native appearance model that facilitates the distinction between the object
 460 and the background based on superpixel over-segmentation. Note that our
 461 method also presents a discriminative aspect, since it uses information on
 462 the background color distribution to evaluate the tracking quality (see the

463 update subsection under section 3.5).

464 **Abrupt motion and out of plane rotation:** The target object in
465 *Sylvester* undergoes out of plane rotation and sudden movements during
466 more than 1300 frames. Most of the trackers, except AST and ours do not
467 perform well. In the *girl* video, the tracked face undergoes both pose change
468 and 360 degrees rotations abruptly. Our method had the highest success rate
469 and was significantly more robust and accurate than most of the methods
470 as we can see in figure 5. SAT handled efficiently pose change and partial
471 occlusion and our tracking was successful as long as the girl’s face was at
472 least partly visible. The target was lost only during the frames where it is
473 completely turned away from the camera (see the OR plot, frames 87-116
474 and 187-250), but tracking is recovered as soon as the face reappears.

475 **Computational cost:** Our tracker was implemented using Matlab on
476 a PC with a Core i7-3770 CPU running at a 3.4 GHz. SAT algorithm is
477 designed to maintain a reasonable computational complexity. In fact, we
478 extract local features in a limited image region determined by particle fil-
479 tering, in order to reduce the computational cost of keypoint detection and
480 local descriptors creation. The particle filter generates $N = 400$ particles,
481 among which only N^* particles are considered as a reduced search space,
482 and for generating the N particles on the subsequent frame. In practice, the
483 computation time of our tracker is closely related to the number of detected
484 keypoints voting for the object position, which mainly depends on the object
485 size and texture. As an example, the video sequences *tiger 1* and *tiger 2*,
486 with a small target size, are processed at nearly one second per frame. On
487 the other hand, when the object size is larger such as in the *occluded face 1*,

	SPT	SBCT	AST	OMSIT	SAT
time/video	1854.31	1990.52	259.84	1327.23	707.41
time/frame	3.95	4.24	0.55	2.82	1.51
ranking	4	5	1	3	2

Table 3: Processing time comparison on the *David indoor* sequence. time/video: the total processing time (seconds), time/frame: the average processing time for one frame (seconds).

488 SAT requires up to 3 seconds to find the target on certain frames. The table
489 3 provides a computation time comparison for the five trackers on the face
490 tracking video *David indoor*. All the compared trackers were implemented in
491 Matlab by the authors, and run on the same described computer. According
492 to the performed measures, our algorithm requires in average 1.51 s to pro-
493 cess one frame, which is the second best execution time. We note that AST
494 achieved the shortest time, processing one frame in 0.55 s.

495 **Application constraints and risk of failure:** The proposed tracker
496 uses SIFT algorithm as an external mechanism to detect the target keypoints.
497 Generally, our method achieves high accuracy when a significant number of
498 keypoints are detected on the target object. On the other hand, the tracking
499 quality may decrease if the target region is not sufficiently textured, or if
500 it is too far from the camera (object details not visible). As an example,
501 we verified that the face tracking application requires a maximum distance
502 of 10 meters between the tracked person and the camera. At this distance,
503 SIFT allows detecting between two and four keypoints in most face tracking
504 scenarios. Furthermore, a drastic decrease in the number of visible target

keypoints increases the drifting risk, regardless of the target type. In practice,
 our tracker relies on keypoint matching only if at least three keypoints from
 the reservoir are matched on the current frame. Otherwise, SAT applies the
 particle filter (that we use to reduce the search space) to track the object
 based on its global color distribution. Another limitation may result from
 the use of a small number of particles to limit the keypoint detection region.
 Indeed, the target may undergo large displacements between consecutive
 frames due to fast movements or low frame rates (e.g. real-time tracking
 using a remote IP camera). As a result, the target object may be located
 outside the keypoint detection area, causing tracking failure. If this situation
 occurs, tracking can be recovered only if the target reappears in the reduced
 search space. Note that this problem can be solved at the cost of an additional
 computation time, by increasing the number of particles (N^*) forming the
 reduced search space.

5. Conclusion

In this paper, we proposed a robust tracking algorithm named SAT
 (Structure Aware Tracker). Our core idea is to exploit the structural prop-
 erties of the target, in a voting-based method, to provide accurate location
 prediction. The target is described by color distribution, keypoints, and their
 geometrical constraints encoding the object internal structure. This multi-
 features appearance model is learned during tracking and thus incorporates
 new structural properties in an online manner. Numerous experiments in a
 comparison with four state-of-the-art trackers, on eleven challenging video
 sequences, demonstrate the superiority of the proposed method in handling

multiple tracking perturbation factors. Our results also highlight the importance of encoding the object structure via keypoint regions, that are more invariant and stable than other types of patches (*e.g.* the local patches encoding the object spatial information in AST and SBCT).

Acknowledgements

This work was supported by a scholarship from FRQ-NT and partially supported by NSERC discovery grant No. 311869-2010.

References

- [1] B. Babenko, M.-H. Y. S. Belongie, Robust object tracking with online multiple instance learning, IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI).
- [2] Z. Kalal, J. Matas, K. Mikolajczyk, Pn learning: Bootstrapping binary classifiers by structural constraints, in: Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on, IEEE, 2010, pp. 49–56.
- [3] L. Zhang, L. van der Maaten, Structure preserving object tracking, in: Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on, 2013, pp. 1838–1845.
- [4] J. M. LukasCerman, V. Hlavac, Sputnik tracker: Having a companion improves robustness of the tracker, in: Image Analysis: 16Th Scandinavian Conference, Scia 2009, Oslo, Norway, June 15-18, Proceedings, Vol. 5575, Springer, 2009, p. 291.

- 550 [5] M. Yang, Y. Wu, G. Hua, Context-aware visual tracking, Pattern Anal-
551 ysis and Machine Intelligence, IEEE Transactions on 31 (7) (2009) 1195–
552 1209.
- 553 [6] L. Wen, Z. Cai, Z. Lei, D. Yi, S. Z. Li, Online spatio-temporal structural
554 context learning for visual tracking, in: Computer Vision–ECCV 2012,
555 Springer, 2012, pp. 716–729.
- 556 [7] H. Grabner, J. Matas, L. Van Gool, P. Cattin, Tracking the invisible:
557 Learning where the object might be, in: Computer Vision and Pattern
558 Recognition (CVPR), 2010 IEEE Conference on, IEEE, 2010, pp. 1285–
559 1292.
- 560 [8] A. Saffari, M. Godec, T. Pock, C. Leistner, H. Bischof, On-
561 line multi-class lpboost, in: Computer Vision and Pattern Recog-
562 nition (CVPR), 2010 IEEE Conference on, 2010, pp. 3570–3577.
563 doi:10.1109/CVPR.2010.5539937.
- 564 [9] W. Bouachir, G.-A. Bilodeau, Structure-aware keypoint tracking for
565 partial occlusion handling, IEEE Winter Conference on Applications
566 of Computer Vision (WACV 2014).
- 567 [10] G. D. Hager, M. Dewan, C. V. Stewart, Multiple kernel tracking with
568 ssd, in: Computer Vision and Pattern Recognition, 2004. CVPR 2004.
569 Proceedings of the 2004 IEEE Computer Society Conference on, Vol. 1,
570 IEEE, 2004, pp. I–790.
- 571 [11] S. Hare, A. Saffari, P. H. Torr, Struck: Structured output tracking with

- 572 kernels, in: Computer Vision (ICCV), 2011 IEEE International Confer-
573 ence on, IEEE, 2011, pp. 263–270.
- 574 [12] G. Shu, A. Dehghan, O. Oreifej, E. Hand, M. Shah, Part-based multiple-
575 person tracking with partial occlusion handling, in: Computer Vision
576 and Pattern Recognition (CVPR), 2012 IEEE Conference on, IEEE,
577 2012, pp. 1815–1821.
- 578 [13] H. Zhou, Y. Yuan, C. Shi, Object tracking using sift features and mean
579 shift, Computer Vision and Image Understanding 113 (3) (2009) 345–
580 352.
- 581 [14] S. Hare, A. Saffari, P. H. Torr, Efficient online structured output learning
582 for keypoint-based object tracking, in: Computer Vision and Pattern
583 Recognition (CVPR), 2012 IEEE Conference on, IEEE, 2012, pp. 1894–
584 1901.
- 585 [15] W. Bouachir, G.-A. Bilodeau, Visual face tracking: A
586 coarse-to-fine target state estimation, 2013 International Con-
587 ference on Computer and Robot Vision 0 (2013) 45–51.
588 doi:<http://doi.ieeecomputersociety.org/10.1109/CRV.2013.18>.
- 589 [16] S. Gu, C. Tomasi, Branch and track, in: Computer Vision and Pattern
590 Recognition (CVPR), 2011 IEEE Conference on, IEEE, 2011, pp. 1169–
591 1174.
- 592 [17] T. B. Dinh, N. Vo, G. Medioni, Context tracker: Exploring supporters
593 and distracters in unconstrained environments, in: Computer Vision and

- 594 Pattern Recognition (CVPR), 2011 IEEE Conference on, IEEE, 2011,
595 pp. 1177–1184.
- 596 [18] K. M. Yi, H. Jeong, B. Heo, H. J. Chang, J. Y. Choi, Initialization-
597 insensitive visual tracking through voting with salient local features,
598 2013 International Conference on Computer Vision (ICCV).
- 599 [19] W. Zhong, H. Lu, M.-H. Yang, Robust object tracking via sparsity-
600 based collaborative model, in: Computer Vision and Pattern Recognition
601 (CVPR), 2012 IEEE Conference on, IEEE, 2012, pp. 1838–1845.
- 602 [20] X. Jia, H. Lu, M.-H. Yang, Visual tracking via adaptive structural local
603 sparse appearance model, in: Computer Vision and Pattern Recognition
604 (CVPR), 2012 IEEE Conference on, IEEE, 2012, pp. 1822–1829.
- 605 [21] Y. Guo, Y. Chen, F. Tang, A. Li, W. Luo, M. Liu, Object tracking using
606 learned feature manifolds, Computer Vision and Image Understanding
607 118 (2014) 128–139.
- 608 [22] F. Yang, H. Lu, M.-H. Yang, Learning structured visual dictionary for
609 object tracking, Image and Vision Computing 31 (12) (2013) 992–999.
- 610 [23] V. Belagiannis, F. Schubert, N. Navab, S. Ilic, Segmentation based par-
611 ticle filtering for real-time 2d object tracking, Computer Vision–ECCV
612 2012 (2012) 842–855.
- 613 [24] H.-M. Sun, Skin detection for single images using dynamic skin color
614 modeling, Pattern recognition 43 (4) (2010) 1413–1420.

- 615 [25] D. G. Lowe, Distinctive image features from scale-invariant keypoints,
616 International journal of computer vision 60 (2) (2004) 91–110.
- 617 [26] A. Mian, M. Bennamoun, R. Owens, An efficient multimodal 2d-3d
618 hybrid approach to automatic face recognition, Pattern Analysis and
619 Machine Intelligence, IEEE Transactions on 29 (11) (2007) 1927–1943.
620 doi:10.1109/TPAMI.2007.1105.
- 621 [27] C. Geng, X. Jiang, Face recognition using sift features, in: Image Pro-
622 cessing (ICIP), 2009 16th IEEE International Conference on, 2009, pp.
623 3313–3316. doi:10.1109/ICIP.2009.5413956.
- 624 [28] A. Mian, M. Bennamoun, R. Owens, Keypoint detection and local fea-
625 ture matching for textured 3d face recognition, International Journal of
626 Computer Vision 79 (1) (2008) 1–12. doi:10.1007/s11263-007-0085-5.
627 URL <http://dx.doi.org/10.1007/s11263-007-0085-5>
- 628 [29] J. Heinly, E. Dunn, J.-M. Frahm, Comparative evaluation of binary
629 features, Computer Vision–ECCV 2012 (2012) 759–773.
- 630 [30] M. Isard, A. Blake, Condensation: conditional density propagation for
631 visual tracking, International journal of computer vision 29 (1) (1998)
632 5–28.
- 633 [31] L. Matthews, T. Ishikawa, S. Baker, The template update problem,
634 Pattern Analysis and Machine Intelligence, IEEE Transactions on 26 (6)
635 (2004) 810–815.
- 636 [32] K. Bernardin, F. Van De Camp, R. Stiefelhagen, Automatic person de-
637 tection and tracking using fuzzy controlled active cameras, in: Computer

- 638 Vision and Pattern Recognition, 2007. CVPR'07. IEEE Conference on,
639 IEEE, 2007, pp. 1–8.
- 640 [33] D. A. Ross, J. Lim, R.-S. Lin, M.-H. Yang, Incremental learning for
641 robust visual tracking, *International Journal of Computer Vision* 77 (1-
642 3) (2008) 125–141.
- 643 [34] S. Birchfield, Elliptical head tracking using intensity gradients and color
644 histograms, in: *Computer Vision and Pattern Recognition, 1998. Pro-*
645 *ceedings. 1998 IEEE Computer Society Conference on, IEEE, 1998, pp.*
646 *232–237.*
- 647 [35] A. Adam, E. Rivlin, I. Shimshoni, Robust fragments-based tracking us-
648 ing the integral histogram, in: *Computer Vision and Pattern Recogni-*
649 *tion, 2006 IEEE Computer Society Conference on, Vol. 1, IEEE, 2006,*
650 *pp. 798–805.*
- 651 [36] S. Wang, H. Lu, F. Yang, M.-H. Yang, Superpixel tracking, in: *Com-*
652 *puter Vision (ICCV), 2011 IEEE International Conference on, IEEE,*
653 *2011, pp. 1323–1330.*
- 654 [37] Q.-H. Zhou, H. Lu, M.-H. Yang, Online multiple support instance track-
655 ing, in: *Automatic Face & Gesture Recognition and Workshops (FG*
656 *2011), 2011 IEEE International Conference on, IEEE, 2011, pp. 545–*
657 *552.*

Algorithm 3 Predicting the target location

```
1: - initialize  $RF$ ,  $\hat{q}$ ,  $\hat{q}_{bg}$ 
2: for all  $frames$  do
3:   - reduce the search space: Alg. 1
4:   - predict a preliminary state: Alg. 2
5:   for all  $voting\_space\_positions$   $\mathbf{x}$  do
6:     for all  $matched\_features$  ( $f^{(i)} \in F_t$ ) do
7:       - estimate  $P(\mathbf{x}|f^{(i)})$ : (Eq. 5)
8:     end for
9:     - estimate location probability  $P(\mathbf{x})$ : (Eq. 6)
10:  end for
11:  -  $target\_location = select\_peak(voting\_space\_positions)$  {tracker's
    output for the current frame}
12:  if ( $update\_condition == true$ ) then
13:    - update  $\hat{q}_t$  and  $\hat{q}_{bg,t}$ : (Eq. 7 & 8)
14:    for all  $matched\_features$  ( $f^{(i)} \in F_t$ ) do
15:      - update  $\mu_t^{(i)}$  (Eq. 10 & 11)
16:      - update  $\Sigma_t^{(i)}$  (Eq. 12)
17:    end for
18:    - update  $w_t^{(i)}$  (Eq. 9) for the entire reservoir
19:    - remove non-persistent features (i.e.  $w_t^{(i)} \leq \theta_p$ )
20:    for all  $newly\_detected\_features$   $f^{(i)}$  do
21:      - add  $f^{(i)}$  to  $RF$ 
22:      -  $\mu_t^{(i)} = [\Delta_{x,new}^{(i)}, \Delta_{x,new}^{(i)}]$ ;  $\Sigma_t^{(i)} = \sigma_0^2 I_2$ ;  $w_t^{(i)} = w_0$ 
23:    end for
24:  end if
25: end for
```
